

Type-Inference Based Short Cut Deforestation (nearly) without Inlining — Work in Progress —

Olaf Chitil

Lehrstuhl für Informatik II, RWTH Aachen,
52056 Aachen, Germany

`chitil@informatik.rwth-aachen.de`

Abstract

Deforestation optimises a functional program by transforming it into another one that does not create certain intermediate data structures. In [Chi99] we presented a type-inference based deforestation algorithm which performs extensive inlining. However, across module boundaries only limited inlining is practically feasible. Furthermore, inlining is a non-trivial transformation which is therefore best implemented as a separate optimisation pass. To perform short cut deforestation (nearly) without inlining, Gill suggested to split definitions into workers and wrappers and inline only the small wrappers, which transfer the information needed for deforestation. We show that Gill's use of a function `build` limits deforestation and note that his reasons for using `build` do not apply to our approach. Hence we develop a more general worker/wrapper scheme without `build`. We give a type-inference based algorithm which splits definitions into workers and wrappers. Finally, we show that we can deforest more expressions with the worker/wrapper scheme than the algorithm with inlining.

1 Type-Inference-Based Short Cut Deforestation

In lazy functional programs two functions are often glued together by an intermediate data structure that is produced by one function and consumed by the other. For example, the function `any`, which tests whether any element of a list `xs` satisfies a given predicate `p`, may be defined as follows in Haskell [PH⁺99]:

```
any p xs = or (map p xs)
```

The function `map` applies `p` to all elements of `xs` yielding a list of boolean values. The function `or` combines these boolean values with the logical or operation (`||`).

Although lazy evaluation makes this modular programming style practicable [Hug89], it does not come for free. Each list cell has to be allocated, filled, taken apart and finally garbage collected. The following monolithic definition of `any` is more efficient.

```
any p []      = False
any p (x:xs) = p x || any p xs
```

It is the aim of *deforestation* algorithms to automatically transform a functional program into another one that does not create such intermediate data structures. We say that the producer and the consumer of the data structure are *fused*.

1.1 Short Cut Deforestation

The fundamental idea of short cut deforestation is to restrict deforestation to intermediate lists that are consumed by the function `foldr`. This higher-order function uniformly replaces the constructors `(:)` in a list by a given function `c` and the empty list constructor `[]` by a constant `n`:

$$\text{foldr } c \ n \ [x_1, \dots, x_k] = x_1 \ 'c' \ (x_2 \ 'c' \ (x_3 \ 'c' \ (\dots (x_k \ 'c' \ n) \dots)))$$

So if `foldr` consumes a list that is produced by an expression e , short cut deforestation replaces the list constructors already at compile time. We cannot, however, simply replace all list constructors in e . Consider $e = (\text{map } p \ [1,2])$. Here the constructors in `[1,2]` are not to be replaced but those in the definition of `map`, which is not even part of e .

Therefore we need the producer e in a form such that the list constructors that build the intermediate list are explicit and can easily be replaced. The convenient solution is to have the producer in the form $(\lambda c \ n \rightarrow e') \ (:) \ []$ where the abstracted variables `c` and `n` mark the intermediate list constructors `(:)` and `[]`. Then fusion is performed by the simple rule:

$$\text{foldr } e_{(:)} \ e_{[]} \ ((\lambda c \ n \rightarrow e') \ (:) \ []) \rightsquigarrow (\lambda c \ n \rightarrow e') \ e_{(:)} \ e_{[]}$$

The rule removes the intermediate list constructors. A subsequent β -reduction puts the consumer components $e_{(:)}$ and $e_{[]}$ into the places that were before hold by the list constructors.

We observe that generally $e_{(:)}$ and $e_{[]}$ have different types from `(:)` and `[]`. Hence for this transformation to be type correct the function in the producer must be polymorphic. This can be expressed in Haskell with the help of a special function `build` with a second-order type:

```
build :: (forall b. (a -> b -> b) -> b -> b) -> [a]
build g = g (:) []
```

$$\text{foldr } e_{(:)} \ e_{[]} \ (\text{build } e_p) \rightsquigarrow e_p \ e_{(:)} \ e_{[]}$$

Here e_p will always have the form $\lambda c \ n \rightarrow e'$, but this is not necessary for the semantic correctness of the transformation. Strikingly, the polymorphic type of e_p already guarantees the correctness. Intuitively, e_p can only build its value of type `b` from its two term arguments, because only these have the right types. The correctness has been proved in [GLP93, Gil96].

1.2 Derivation of Producers through Type Inference

Whereas using `foldr` for defining list consumers is generally considered as good, modular programming style, programmers can hardly be demanded to use `build`. The idea of the first works on short cut deforestation is that all list-manipulating functions in the standard libraries are defined in terms of `foldr` and `build`. However, thus deforestation is confined to producers that are defined in terms of these standard list functions.

On the other hand we see that, if we can transform a producer e of type $[\tau]$ into the form `build (\lambda c \ n -> e')`, then the type system guarantees that we have abstracted exactly those

list constructors that build the intermediate list. Based on this observation we presented in [Chi99] a type-inference based algorithm which abstracts the intermediate list type and its constructors from a producer to obtain a `build` form.

For the producer `map p [1,2]` for example, this algorithm observes, that the intermediate list is constructed by the function `map`. Therefore it inlines the body of `map` to be able to proceed. Afterwards the algorithm decides that the list constructors in the body of `map` have to be abstracted whereas the list constructors in `[1,2]` remain unchanged. With this result the algorithm terminates successfully. In general, the algorithm recursively inlines all functions that are needed to be able to derive the desired form of the producer, only bounded by an arbitrary code size limit. We recapitulate the algorithm in more detail in Section 3.

1.3 The Problem of Inlining

It is neat that the algorithm determines exactly those functions which need to be inlined, but nonetheless it causes problems in practise. Inlining across module boundaries is usually implemented by saving some small function definitions in a file when a module is compiled. This file is read and used when a module is compiled which imports the former module. Extensive inlining across module boundaries would defeat the idea of separate compilation. In general, inlining, although trivial in principal, is in practise “a black art, full of delicate compromises that work together to give good performance without unnecessary code bloat” [PM99]. It is best implemented as a separate optimisation pass. Consequently, we would like to use our list abstraction algorithm without it having to perform inlining itself.

Gill already suggested a method to separate short cut deforestation from inlining[Gil96]: Each list-producing function definition is split into a worker and a wrapper. The latter is small enough to be inlined everywhere, also across module boundaries, and transfers enough information to permit short cut deforestation. Inlining can be performed subsequently to improve the result of deforestation further. However, Gill’s use of `build` in wrappers limits the expressibility of his worker/wrapper scheme. There are functions that cannot be defined in terms of `build` and hence cannot be deforested within Gill’s worker/wrapper scheme, but that can be deforested by our type-inference based algorithm with inlining. Fortunately, the function `build` turns out to be unnecessary for our type-inference based deforestation algorithm. Hence we developed a more general worker/wrapper scheme without `build`. We present it in Section 4.

Afterwards we show in Section 5 how the list-producing function definitions of a program are split into wrappers and workers by an algorithm that is based on our type-inference based algorithm without using inlining.

In Section 6 we present a class of recursively defined functions which cannot be deforested by the algorithm with inlining but can be deforested within the worker/wrapper scheme. However, to be able to split these function definitions into the required workers and wrappers we need to extend our worker/wrapper split algorithm. As basis we use Mycroft’s extension of the Damas-Milner type inference algorithm by polymorphic recursion [Myc84].

2 The second-order typed language

We use a small functional language with second-order types, which is similar to the intermediate language Core used inside the Glasgow Haskell compiler [GHC]. The syntax is defined in Figure 1 and the type system in Figure 2. The language is essentially the second-order

Type constructors	C	$::=$	$[\] \mid \mathbf{Int} \mid \dots$
Type variables	α, β, γ		
Types	τ	$::=$	$C\bar{\tau} \mid \alpha \mid \tau_1 \rightarrow \tau_2 \mid \forall\alpha.\tau$
Term variables	x, c		
Terms	e	$::=$	$x \mid \lambda x : \tau.e \mid e_1 e_2 \mid \mathbf{case} e \mathbf{of} \{c_i \bar{x}_i \rightarrow e_i\}_{i=1}^k \mid \mathbf{let} \{x_i : \tau_i = e_i\}_{i=1}^k \mathbf{in} e \mid \lambda\alpha.e \mid e\tau$

Figure 1: Terms and types of the language

$$\begin{array}{c}
\frac{}{\Gamma + x : \tau \vdash x : \tau} \text{VAR} \\
\\
\frac{\Gamma + x : \tau_1 \vdash e : \tau_2}{\Gamma \vdash \lambda(x : \tau_1).e : \tau_1 \rightarrow \tau_2} \text{TERM ABS} \qquad \frac{\Gamma \vdash e_1 : \tau_2 \rightarrow \tau \quad \Gamma \vdash e_2 : \tau_2}{\Gamma \vdash e_1 e_2 : \tau} \text{TERM APP} \\
\\
\frac{\forall i = 1..k \quad \Gamma + \{x_j : \tau_j\}_{j=1}^k \vdash e_i : \tau_i \quad \Gamma + \{x_j : \tau_j\}_{j=1}^k \vdash e : \tau}{\Gamma \vdash \mathbf{let} \{x_i : \tau_i = e_i\}_{i=1}^k \mathbf{in} e : \tau} \text{LET} \\
\\
\frac{\Gamma \vdash e : C\bar{\rho} \quad \Gamma(c_i) = \forall\bar{\alpha}.\bar{\rho}_i \rightarrow C\bar{\alpha} \quad \Gamma + \{\bar{x}_i : \bar{\rho}_i[\bar{\rho}/\bar{\alpha}]\} \vdash e_i : \tau \quad \forall i = 1..k}{\Gamma \vdash \mathbf{case} e \mathbf{of} \{c_i \bar{x}_i \mapsto e_i\}_{i=1}^k : \tau} \text{CASE} \\
\\
\frac{\Gamma \vdash e : \tau \quad \alpha \notin \text{freeTyVar}(\Gamma)}{\Gamma \vdash \lambda\alpha.e : \forall\alpha.\tau} \text{TYPE ABS} \qquad \frac{\Gamma \vdash e : \forall\alpha.\tau}{\Gamma \vdash e\rho : \tau[\rho/\alpha]} \text{TYPE APP}
\end{array}$$

Figure 2: Type system

typed λ -calculus augmented with **let** for arbitrary mutual recursion and **case** for decomposition of algebraic data structures. We view a type environment Γ as both a mapping from variables to types and a set of tuples $x : \tau$. The operator $+$ combines two type environments under the assumption that their domain is disjoint. We abbreviate $\Gamma + \{x : \tau\}$ by $\Gamma + x : \tau$. Data constructors c are just special term variables. The language does not have explicit definitions of algebraic data types like **data** $C\bar{\alpha} = c_1\bar{\tau}_1 \mid \dots \mid c_k\bar{\tau}_k$. Such a definition is implicitly expressed by having the data constructors in the type environment: $\Gamma(c_i) = \tau_{1,i} \rightarrow \dots \rightarrow \tau_{n_i,i} \rightarrow C\bar{\alpha} = \bar{\tau}_i \rightarrow C\bar{\alpha}$. Hence for the polymorphic type list, which we write $[\alpha]$ instead of $[\] \alpha$, we have $\Gamma((:)) = \forall\alpha.\alpha \rightarrow [\alpha] \rightarrow [\alpha]$ and $\Gamma([\]) = \forall\alpha.[\alpha]$. The functions **foldr** and **build** are defined as follows

$$\begin{aligned}
\mathbf{foldr} &: \forall\alpha.\forall\beta.(\alpha \rightarrow \beta \rightarrow \beta) \rightarrow \beta \rightarrow [\alpha] \rightarrow \beta \\
&= \lambda\alpha.\lambda\beta.\lambda c:\alpha \rightarrow \beta \rightarrow \beta.\lambda n:\beta.\lambda xs:[\alpha].\mathbf{case} \mathbf{xs} \mathbf{of} \{ \\
&\quad [\] \quad \rightarrow n \\
&\quad y:\mathbf{ys} \rightarrow c \ y \ (\mathbf{foldr} \ \alpha \ \beta \ c \ n \ \mathbf{ys}) \} \\
\\
\mathbf{build} &: \forall\alpha.(\forall\beta.(\alpha \rightarrow \beta \rightarrow \beta) \rightarrow \beta \rightarrow \beta) \rightarrow [\alpha] \\
&= \lambda\alpha.\lambda g:\forall\beta.(\alpha \rightarrow \beta \rightarrow \beta) \rightarrow \beta \rightarrow \beta.g \ [\alpha] \ ((:)) \ \alpha \ ([\] \ \alpha)
\end{aligned}$$

and the `foldr/build` rule takes the form:

$$\text{foldr } \tau_1 \tau_2 e_{(:)} e_{[]} (\text{build } \tau_1 e_p) \rightsquigarrow e_p \tau_2 e_{(:)} e_{[]}$$

3 List Abstraction through Type Inference

To understand the mode of operation of the list abstraction algorithm of [Chi99] we study an example. We have to start with the typing of the producer:¹

$$\begin{aligned} & \{\text{map}' : (\text{Int} \rightarrow \text{Int}) \rightarrow [\text{Int}] \rightarrow [\text{Int}], \text{g} : \text{Int} \rightarrow \text{Int}, 1 : \text{Int}, 2 : \text{Int}, \\ & \quad (:) : \forall \alpha. \alpha \rightarrow [\alpha] \rightarrow [\alpha], [] : \forall \alpha. [\alpha]\} \\ & \vdash \text{map}' \text{g} ((:) \text{Int} 1 ((:) \text{Int} 2 ([] \text{Int}))) : [\text{Int}] \end{aligned}$$

The algorithm replaces every list constructor application $(:) \tau$, respectively $[] \tau$, by a different variable c_i , respectively n_i . To use the existing type annotations as far as possible, we just replace every list type² in the expression and in the type environment by a new type inference variable. Furthermore, we add $c_i : \tau \rightarrow \gamma_i \rightarrow \gamma_i$, respectively $n_i : \gamma_i$, to the type environment, where γ_i is a new type inference variable for every variable c_i , respectively n_i .

$$\begin{aligned} & \{\text{map}' : (\text{Int} \rightarrow \text{Int}) \rightarrow \gamma_1 \rightarrow \gamma_2, \text{g} : \text{Int} \rightarrow \text{Int}, 1 : \text{Int}, 2 : \text{Int}, \\ & \quad \text{n}_1 : \gamma_3, \text{c}_1 : \text{Int} \rightarrow \gamma_4 \rightarrow \gamma_4, \text{c}_2 : \text{Int} \rightarrow \gamma_5 \rightarrow \gamma_5\} \\ & \vdash \text{map}' \text{g} (\text{c}_1 1 (\text{c}_2 2 \text{n}_1)) : \gamma \end{aligned}$$

This typing with type variables is the input to a modified version of the Damas-Milner type inference algorithm [DM82, LY98]. On the one hand the algorithm was extended to cope with explicit type abstraction and application. On the other hand the type generalisation step (type closure) at `let` bindings was dropped. The type inference algorithm replaces some of the type variables so that the typing is again derivable from the type inference rules, that is, the expression is well-typed in the type environment. Note that type inference cannot fail, because the typing we start with is derivable. We just try to find a more general typing.

$$\begin{aligned} & \{\text{map}' : (\text{Int} \rightarrow \text{Int}) \rightarrow \gamma_1 \rightarrow \gamma, \text{g} : \text{Int} \rightarrow \text{Int}, 1 : \text{Int}, 2 : \text{Int}, \\ & \quad \text{n}_1 : \gamma_1, \text{c}_1 : \text{Int} \rightarrow \gamma_1 \rightarrow \gamma_1, \text{c}_2 : \text{Int} \rightarrow \gamma_1 \rightarrow \gamma_1\} \\ & \vdash \text{map}' \text{g} (\text{c}_1 1 (\text{c}_2 2 \text{n}_1)) : \gamma \end{aligned}$$

The type of the expression is a type variable which can be abstracted, but this type variable also appears in the type of the function `map'`. So the the definition of `map'` has to be inlined, all lists types and list constructors be replaced by new variables and type inference be continued.

$$\begin{aligned} & \{\text{g} : \text{Int} \rightarrow \text{Int}, 1 : \text{Int}, 2 : \text{Int}, \text{n}_1 : [\text{Int}], \text{n}_2 : \gamma, \\ & \quad \text{c}_1 : \text{Int} \rightarrow [\text{Int}] \rightarrow [\text{Int}], \text{c}_2 : \text{Int} \rightarrow [\text{Int}] \rightarrow [\text{Int}], \text{c}_3 : \text{Int} \rightarrow \gamma \rightarrow \gamma\} \\ & \vdash \text{let map}' : (\text{Int} \rightarrow \text{Int}) \rightarrow [\text{Int}] \rightarrow \gamma \\ & \quad = \lambda f : \text{Int} \rightarrow \text{Int}. \text{foldr Int } \gamma (\lambda v : \text{Int}. \lambda w : \gamma. \text{c}_3 (f v) w) \text{n}_3 \\ & \quad \text{in map}' \text{g} (\text{c}_2 1 (\text{c}_3 2 \text{n}_2)) : \gamma \end{aligned}$$

¹We only consider a monomorphic version of `map`. Polymorphic functions that may be inlined require an additional instantiation step which we skip here (see [Chi99], Section 4.2).

²We only need to replace the type `[Int]`, because it is the type to be abstracted from. We are not permitted to replace list types which contain locally bound type variables.

Now the type of the expression is still a type variable that, however, does not occur in the type environment except in the types of the c_i and n_i . Hence the algorithm terminates successfully. The type environment tells us that c_1 and n_1 construct the result of the producer whereas c_2 , c_3 , and n_2 have to construct lists that are internal to the producer. So the producer can be expressed with abstracted list type and list constructors as follows, suitable as argument for `build`:

```

λγ. λc:Int → γ → γ. λn:γ.
  let map' : (Int→Int)→[Int]→γ
      = λf:Int→Int. foldr Int γ (λv:Int.λw:γ. c (f v) w) n
  in map' g ((:) Int 1 ((:) Int 2 ([] Int)))

```

4 The Worker/Wrapper Scheme

To apply the `foldr/build` rule we do not need to fully inspect the consumer and the producer. We only have to see that the consumer is in `foldr` form and the producer is in `build` form. The other arguments of `foldr` and the argument of `build` are of no interest, they are just rearranged by the transformation. This observation is the basis for obtaining producers in `build` form without inlining of large expressions.

4.1 Gill's Worker/Wrapper Scheme

Gill ([Gil96], Section 7.4) takes the definition of a list producer function from which the list constructors are abstracted

```
f x1 ... xn = build (\c n -> e)
```

and splits it into a definition of two functions, a worker `fW` and a wrapper `f`:

```
fW x1 ... xn c n = e
f x1 ... xn = build (fW x1 ... xn)
```

The wrapper just contains the `build` and a call to the worker. It is small enough to be inlined at all its call sites. For example, the definition

```
map f xs = build (\c n -> foldr (c . f) n xs)
```

is split up as follows:

```
mapW :: (a -> b) -> [a] -> (b -> c -> c) -> c -> c
mapW f xs c n = foldr (c . f) n xs
```

```
map f xs = build (mapW f xs)
```

Consider deforestation of the definition body of `any` that was given first in the introduction:

```

or (map p xs)
↔ foldr (||) True (build (mapW f xs))
↔ mapW f xs (||) True

```

First `or` and `map` are inlined. Then fusion is performed by the `foldr/build` rule. Afterwards it is left to the standard inliner, if `mapW` is inlined. This will be the case for this simple example, but not in general. In any case the expression has been deforested.

4.2 The Limitations of `build`

The use of `build` unfortunately limits the scope of Gill’s worker/wrapper scheme. In [Chi99] we showed that the type-inference based algorithm can abstract the list constructors from the producer `fst (unzip zs)`, where the function `unzip` maps a list of pairs to a pair of lists. To abstract the list the algorithm inlines the definition of `unzip`. The function `unzip` cannot be expressed in terms of `build`, because `build` can only wrap a producer that returns a single list. So if we replace the inlining of the type abstraction algorithm by using Gill’s worker/wrapper scheme we lose the ability to fuse a consumer with `fst (unzip zs)`.

There is also a second problem which is related to our list abstraction algorithm. We still want to use it — without inlining — to abstract the intermediate list constructors from the producer, because not all producers will already be in `build` form after inlining of wrappers. However, a `build` in the producer hinders the type-inference algorithm. For example from the producer `build (mapW f xs)` no list constructors can be abstracted, because they are hidden by `build`. We have to inline `build` to proceed with list constructor abstraction.

So, because `build` is both inflexible and is not needed but even disturbs the abstraction of list constructors, we draw the obvious conclusion to do without it.

4.3 Short Cut Deforestation without `build`

The purpose of `build` is to express in Haskell that the producer has the required polymorphic type and to prevent this special form from being destroyed by another compiler optimisation. Because we perform our transformation in a second-order typed language with explicit type abstraction and application we can enforce the type requirement directly. Instead of transforming a producer into the form

$$\text{build } \tau \ (\lambda\beta. \lambda c:\tau \rightarrow \beta \rightarrow \beta. \lambda n:\beta. e')$$

we transform it into

$$(\lambda\beta. \lambda c:\tau \rightarrow \beta \rightarrow \beta. \lambda n:\beta. e') \ [\tau] \ ((:) \ \tau) \ (\square \ \tau)$$

The short cut fusion rule now looks as follows:

$$\begin{aligned} &\text{for all } e_{(:)}:\tau \rightarrow \tau' \rightarrow \tau', \quad e_{\square}:\tau', \quad e_p:\forall\gamma.(\tau \rightarrow \gamma \rightarrow \gamma) \rightarrow \gamma \rightarrow \gamma \\ &\text{foldr } \tau \ \tau' \ e_{(:)} \ e_{\square} \ (e_p \ [\tau] \ ((:) \ \tau) \ (\square \ \tau)) \rightsquigarrow e_p \ \tau' \ e_{(:)} \ e_{\square} \end{aligned}$$

The `foldr/build` rule looks more simple, because `build` is easier to recognise than the producer form of the new rule. However, we neither have to construct nor to search for such a producer form. Our short cut deforestation algorithm searches for occurrences of `foldr`, abstracts the result list from the producer and then directly applies the short cut fusion rule. Therefore we also do not have to worry that another compiler optimisation might β -reduce the special producer form.

4.4 Wrappers without `build`

Similarly we do not use `build` in wrappers. As an example for our new worker/wrapper scheme consider the worker and wrapper for the function `map`:

$$\begin{aligned} \text{mapW} &: \forall\alpha.\forall\beta.\forall\gamma. (\beta \rightarrow \gamma \rightarrow \gamma) \rightarrow \gamma \rightarrow (\alpha \rightarrow \beta) \rightarrow [\alpha] \rightarrow \gamma \\ &= \lambda\alpha. \lambda\beta. \lambda\gamma. \lambda c: \beta \rightarrow \gamma \rightarrow \gamma. \lambda n: \gamma. \lambda f: \alpha \rightarrow \beta. \\ &\quad \text{foldr } \alpha \ \gamma \ (\lambda v: \alpha. \lambda w: \gamma. c \ (f \ v) \ w) \ n \end{aligned}$$

$$\begin{aligned} \text{map} &: \forall\alpha.\forall\beta. (\alpha \rightarrow \beta) \rightarrow [\alpha] \rightarrow [\beta] \\ &= \lambda\alpha. \lambda\beta. \text{mapW } \alpha \ \beta \ [\beta] \ ((:)\ \beta) \ ([\] \ \beta) \end{aligned}$$

The wrapper calls the worker directly with a list type and respective list constructors as arguments. The worker is almost identical to Gill’s worker. Note that we insert the abstraction from list type and its respective list constructors between the type abstractions and the term abstractions. We cannot insert the new term abstractions in front of the original type abstractions, because the list type $[\beta]$, from which we abstract, contains a type variable β which is bound in the type of the function. To insert the new abstractions before the original term abstractions has several minor advantages. First, we thus do not require that all term arguments are λ -abstracted at the top of the definition body, the definition may also use partial application. Second, the wrapper can be inlined and β -reduced at all call sites, even where it is only partially applied, because we can use partial application in the definition of the wrapper.

To see the gained expressiveness consider the function `unzip`, which produces two lists. The worker abstracts from both list types and their respective sets of list constructors.

$$\begin{aligned} \text{unzipW} &: \forall\alpha.\forall\beta.\forall\gamma_1.\forall\gamma_2. (\alpha \rightarrow \gamma_1 \rightarrow \gamma_1) \rightarrow \gamma_1 \rightarrow (\beta \rightarrow \gamma_2 \rightarrow \gamma_2) \rightarrow \gamma_2 \rightarrow [(\alpha, \beta)] \rightarrow (\gamma_1, \gamma_2) \\ &= \lambda\alpha. \lambda\beta. \lambda\gamma_1. \lambda\gamma_2. \lambda c_1: \alpha \rightarrow \gamma_1 \rightarrow \gamma_1. \lambda n_1: \gamma_1. \lambda c_2: \beta \rightarrow \gamma_2 \rightarrow \gamma_2. \lambda n_2: \gamma_2. \\ &\quad \text{foldr } (\alpha, \beta) \ (\gamma_1, \gamma_2) \\ &\quad (\lambda y: (\alpha, \beta). \lambda u: (\gamma_1, \gamma_2). \text{case } y \text{ of } \{(v, w) \rightarrow \text{case } u \text{ of } \{(vs, ws) \rightarrow \\ &\quad (\,) \ \gamma_1 \ \gamma_2 \ (c_1 \ v \ vs) \ (c_2 \ w \ ws) \} \}) \\ &\quad ((,)\ \gamma_1 \ \gamma_2 \ n_1 \ n_2) \end{aligned}$$

$$\begin{aligned} \text{unzip} &: \forall\alpha.\forall\beta. [(\alpha, \beta)] \rightarrow ([\alpha], [\beta]) \\ &= \lambda\alpha. \lambda\beta. \text{unzipW } \alpha \ \beta \ [\alpha] \ [\beta] \ ((:)\ \alpha) \ ([\] \ \alpha) \ ((:)\ \beta) \ ([\] \ \beta) \end{aligned}$$

4.5 Deforestation (nearly) without Inlining

The new wrappers transfer the information needed for the list abstraction algorithm. Let us reconsider the expression `fst $[\tau_1]$ $[\tau_2]$ (unzip τ_1 τ_2 zs)`. After the wrapper `unzip` is inlined

$$\text{fst } [\tau_1] \ [\tau_2] \ (\text{unzipW } \tau_1 \ \tau_2 \ [\tau_1] \ [\tau_2] \ ((:)\ \tau_1) \ ([\] \ \tau_1) \ ((:)\ \tau_2) \ ([\] \ \tau_2) \ \text{zs})$$

all result list types and constructors are exposed and the type-inference based algorithm can abstract them without any further inlining:

$$\lambda\gamma. \lambda c: \tau_1 \rightarrow \gamma \rightarrow \gamma. \lambda n: \gamma. \text{fst } \gamma \ [\tau_2] \ (\text{unzipW } \tau_1 \ \tau_2 \ \gamma \ [\tau_2] \ c \ n \ ((:)\ \tau_2) \ ([\] \ \tau_2) \ \text{zs})$$

So fusion with a consumer `foldr τ_1 τ_3 $e_{(:)}$ $e_{[\]}$` results in

$$(\lambda\gamma. \lambda c: \tau_1 \rightarrow \gamma \rightarrow \gamma. \lambda n: \gamma. \text{fst } \gamma \ [\tau_2] \ (\text{unzipW } \tau_1 \ \tau_2 \ \gamma \ [\tau_2] \ c \ n \ ((:)\ \tau_2) \ ([\] \ \tau_2) \ \text{zs})) \ \tau_3 \ e_{(:)} \ e_{[\]}$$

4.6 More Wrappers

The list append function (`++`) is notorious for being difficult to fuse with. The expression `(++) τ xs ys` does not produce the whole result list itself, because only `xs` is copied but not `ys`. Therefore `(++)` cannot be defined in terms of `build`. Gill defines a further second-order typed function `augment` to solve this problem. However, we can easily define a worker for `(++)` by abstracting not just from the result list but simultaneously from the type of the second argument:

$$\begin{aligned} \text{appW} &: \forall \alpha. \forall \gamma. (\alpha \rightarrow \gamma \rightarrow \gamma) \rightarrow \gamma \rightarrow [\alpha] \rightarrow \gamma \rightarrow \gamma \\ &= \lambda \alpha. \lambda \gamma. \lambda c: \alpha \rightarrow \gamma \rightarrow \gamma. \lambda n: \gamma. \lambda xs: [\alpha]. \lambda ys: \gamma. \\ &\quad \text{foldr } \alpha \ \gamma \ c \ ys \ xs \end{aligned}$$

$$\begin{aligned} (++) &: \forall \alpha. [\alpha] \rightarrow [\alpha] \rightarrow [\alpha] \\ &= \lambda \alpha. \text{appW } \alpha \ [\alpha] \ ((:) \ \alpha) \ ([\] \ \alpha) \end{aligned}$$

The type of `appW` implies, that we can only abstract from the result list constructors of an application of `(++)`, if we can abstract from the result list constructors of its second argument. We believe that this will seldom restrict deforestation in practise. For example the definition

$$\begin{aligned} \text{concat} &: \forall \alpha. [[\alpha]] \rightarrow [\alpha] \\ &= \lambda \alpha. \text{foldr } [\alpha] \ [\alpha] \ ((++) \ \alpha) \ ([\] \ \alpha) \end{aligned}$$

can be split into worker and wrapper thanks to the wrapper `appW`:

$$\begin{aligned} \text{concatW} &: \forall \alpha. \forall \gamma. (\alpha \rightarrow \gamma \rightarrow \gamma) \rightarrow \gamma \rightarrow [[\alpha]] \rightarrow \gamma \\ &= \lambda \alpha. \lambda \gamma. \lambda c: \alpha \rightarrow \gamma \rightarrow \gamma. \lambda n: \gamma. \text{foldr } [\alpha] \ \gamma \ (\text{appW } \alpha \ \gamma \ c \ n) \ n \end{aligned}$$

$$\begin{aligned} \text{concat} &: \forall \alpha. [[\alpha]] \rightarrow [\alpha] \\ &= \lambda \alpha. \text{concatW } \alpha \ [\alpha] \ ((:) \ \alpha) \ ([\] \ \alpha) \end{aligned}$$

None of the previous function definitions was recursive. The recursion was hidden by `foldr`. The recursively defined function `tails` returns the list of all final segments of `xs`, longest first:

$$\begin{aligned} \text{tails} &: \forall \alpha. [\alpha] \rightarrow [[\alpha]] \\ &= \lambda \alpha. \lambda xs: [\alpha]. \text{case } xs \text{ of } \{ \\ &\quad [\] \rightarrow (:) \ [\alpha] \ ([\] \ \alpha) \ ([\] \ [\alpha]) \\ &\quad y:ys \rightarrow (:) \ [\alpha] \ xs \ (\text{tails } \alpha \ ys) \} \end{aligned}$$

A list can also be abstracted from this recursive definition:

$$\begin{aligned} \text{tailsW} &: \forall \alpha. \forall \gamma. [\alpha] \rightarrow \gamma \\ &= \lambda \alpha. \lambda \gamma. \lambda c: [\alpha] \rightarrow \gamma \rightarrow \gamma. \lambda n: \gamma. \lambda xs: [\alpha]. \text{case } xs \text{ of } \{ \\ &\quad [\] \rightarrow c \ ([\] \ \alpha) \ n \\ &\quad y:ys \rightarrow c \ xs \ (\text{tailsW } \alpha \ \gamma \ c \ n \ ys) \} \end{aligned}$$

$$\begin{aligned} \text{tails} &: \forall \alpha. [\alpha] \rightarrow [[\alpha]] \\ &= \lambda \alpha. \text{tailsW } \alpha \ [[\alpha]] \ ((:) \ [\alpha]) \ ([\] \ [\alpha]) \end{aligned}$$

Note that to abstract the list the recursive calls in the definition must be to the worker itself. It is not possible to abstract the list first and to inline calls to the wrapper in the definition body later.

4.7 Effects on Performance

As Gill already noticed for his worker/wrapper scheme, there is a substantial performance difference between calling a function as originally defined ($\text{map } \tau' \tau$) and calling the worker with list constructors as arguments ($\text{mapW } \tau' \tau [\tau] ((:) \tau) ([\] \tau)$). Constructing a list with list constructors that are passed as arguments is more expensive than constructing the list directly. After deforestation all calls of workers that were not needed still have list constructors as arguments. So, as Gill suggested, we must have for each worker a version which is specialised to the list constructors and replace the call of each unused worker by a call of its specialised version. We could use the original, unsplit definition of the function, but by specialising the worker we can profit from any optimisations, especially deforestation, which were performed inside the definition of the worker. Note that we only derive one specialised definition for every worker definition.

The worker/wrapper scheme increases code size through the introduction of wrappers and specialised workers. However, this increase is clearly limited in contrast to the code increase that is caused by our original list abstraction algorithm with inlining. An implementation will show if the code size increase is acceptable. Note that the definition of workers which are not needed for deforestation can be removed by standard dead code elimination after worker specialisation has been performed.

5 The Worker/Wrapper Split Algorithm

For the worker/wrapper scheme each list-producing function definition has to be split into a worker and a wrapper definition.

5.1 Derivation of Workers through Type Inference

A worker is easily derived from a non-recursive function definition by application of the list abstraction algorithm. Consider the `foldr` definition of `map`. Only the preceding type abstractions have to be removed to form the input for the list abstraction algorithm:

$$\{\text{foldr} : \forall \alpha. \forall \beta. (\alpha \rightarrow \beta \rightarrow \beta) \rightarrow \beta \rightarrow [\alpha] \rightarrow \beta, (:) : \forall \alpha. \alpha \rightarrow [\alpha] \rightarrow [\alpha], [] : \forall \alpha. [\alpha]\} \\ \vdash \lambda f : \alpha \rightarrow \beta. \text{foldr } \alpha \ [\beta] \ (\lambda v : \alpha. \lambda w : [\beta]. (:) \ \beta \ (f \ v) \ w) \ ([\] \ \beta)$$

The algorithm returns:

$$\lambda \gamma. \lambda c : \beta \rightarrow \gamma \rightarrow \gamma. \lambda n : \gamma. \lambda f : \alpha \rightarrow \beta. \text{foldr } \alpha \ \gamma \ (c \ (f \ v) \ w) \ n$$

So the result list can be abstracted. The readdition of the abstraction from α to obtain the worker and the construction of the wrapper is straightforward. In the case that no list can be abstracted, no worker/wrapper split takes place.

Because all list types in the the type of the processed function are replaced by type variables — not only the list types in the result — also the worker of `(++)` is derived by this algorithm. The derivation of the workers `unzipW` and `concatW` poses no problem either.

5.2 Derivation of Workers of Recursively Defined Functions

For recursive definitions we have to modify the list abstraction algorithm slightly. Consider the definition of `tails` given in Section 4.6. The input typing for the type inference algorithm

must contain a type assignment for the recursive call(s). Because we remove the abstraction from the type variable α , we have to replace `tails` α by a new identifier `tails'`. The type environment assigns the same type to this identifier as is given for the whole definition body. The latter corresponds to the processing of recursive `lets` in the Damas-Milner type inference algorithm.

$$\begin{aligned} & \{\text{tails}' : \gamma_1 \rightarrow \gamma_2, \mathbf{n}_1 : \gamma_3, \mathbf{n}_2 : \gamma_4, \mathbf{c}_1 : \gamma_5 \rightarrow \gamma_6 \rightarrow \gamma_6, \mathbf{c}_2 : \gamma_7 \rightarrow \gamma_8 \rightarrow \gamma_8\} \\ & \vdash \lambda \mathbf{xs} : [\alpha]. \text{ case } \mathbf{xs} \text{ of } \{ \\ & \quad [] \rightarrow \mathbf{c}_1 \ \mathbf{n}_1 \ \mathbf{n}_2 \\ & \quad \mathbf{y} : \mathbf{ys} \rightarrow \mathbf{c}_2 \ \mathbf{xs} \ (\text{tails}' \ \mathbf{ys}) \} : \gamma_1 \rightarrow \gamma_2 \end{aligned}$$

Type inference yields:

$$\begin{aligned} & \{\text{tails}' : [\alpha] \rightarrow \gamma, \mathbf{n}_1 : [\alpha], \mathbf{n}_2 : \gamma, \mathbf{c}_1 : [\alpha] \rightarrow \gamma \rightarrow \gamma, \mathbf{c}_2 : [\alpha] \rightarrow \gamma \rightarrow \gamma\} \\ & \vdash \lambda \mathbf{xs} : [\alpha]. \text{ case } \mathbf{xs} \text{ of } \{ \\ & \quad [] \rightarrow \mathbf{c}_1 \ \mathbf{n}_1 \ \mathbf{n}_2 \\ & \quad \mathbf{y} : \mathbf{ys} \rightarrow \mathbf{c}_2 \ \mathbf{xs} \ (\text{tails}' \ \mathbf{ys}) \} : [\alpha] \rightarrow \gamma \end{aligned}$$

The remaining construction of the worker and wrapper is again straightforward.

For the worker/wrapper split of several mutually recursive definitions the type inference algorithm processes all definitions together.

5.3 Traversal Order

The worker/wrapper split algorithm splits each `let` defined function individually. The example of `concat` in Section 4.6 shows us, that the list could only be abstracted after the wrapper of `(++)` had been inlined. Hence the split algorithm must traverse the program in top-down order and inline wrappers in the remaining program directly after they were derived.

Additionally, definitions can be nested, that is, the right-hand-side of a `let` binding can contain another `let` binding. Here the inner definition has to be split first. Its wrapper can then be inlined in the body of the outer definition and thus enable the abstraction of more lists from the outer definition.

6 Functions that Consume their own Result

There are list functions which consume their own result. The most simple example is the definition of the function that reverses a list in quadratic time:

$$\begin{aligned} \text{reverse} & : \forall \alpha. [\alpha] \rightarrow [\alpha] \\ & = \lambda \alpha. \lambda \mathbf{xs} : [\alpha]. \text{ case } \mathbf{xs} \text{ of } \{ \\ & \quad [] \rightarrow [] \ \alpha \\ & \quad \mathbf{y} : \mathbf{ys} \rightarrow (++) \ \alpha \ (\text{reverse} \ \alpha \ \mathbf{ys}) \ ((:) \ \alpha \ \mathbf{y} \ ([] \ \alpha)) \} \end{aligned}$$

This definition can be split into the following worker and wrapper:

$$\begin{aligned} \text{reverseW} & : \forall \alpha. \forall \gamma. (\alpha \rightarrow \gamma \rightarrow \gamma) \rightarrow \gamma \rightarrow [\alpha] \rightarrow \gamma \\ & = \lambda \alpha. \lambda \gamma. \lambda \mathbf{c} : \alpha \rightarrow \gamma \rightarrow \gamma. \lambda \mathbf{n} : \gamma. \lambda \mathbf{xs} : [\alpha]. \text{ case } \mathbf{xs} \text{ of } \{ \\ & \quad [] \rightarrow \mathbf{n} \\ & \quad \mathbf{y} : \mathbf{ys} \rightarrow \text{appW} \ \alpha \ \gamma \ \mathbf{c} \ \mathbf{n} \ (\text{reverseW} \ \alpha \ [\alpha] \ ((:) \ \alpha) \ ([] \ \alpha) \ \mathbf{ys}) \ (\mathbf{c} \ \mathbf{y} \ \mathbf{n}) \} \end{aligned}$$

$$\begin{aligned} \text{reverse} &: \forall \alpha. [\alpha] \rightarrow [\alpha] \\ &= \lambda \alpha. \text{reverseW } \alpha \text{ } [\alpha] \text{ } ((:) \alpha) \text{ } ([] \alpha) \end{aligned}$$

In the definition of `reverseW` the worker `appW` can be inlined. Then short cut fusion yields:

$$\begin{aligned} \text{reverseW} &: \forall \alpha. \forall \gamma. (\alpha \rightarrow \gamma \rightarrow \gamma) \rightarrow \gamma \rightarrow [\alpha] \rightarrow \gamma \\ &= \lambda \alpha. \lambda \gamma. \lambda c: \alpha \rightarrow \gamma \rightarrow \gamma. \lambda n: \gamma. \lambda xs: [\alpha]. \text{case } xs \text{ of } \{ \\ &\quad [] \quad \rightarrow n \\ &\quad y: ys \rightarrow \text{reverseW } \alpha \ \gamma \ c \ (c \ y \ n) \ ys \} \end{aligned}$$

The deforested version performs list reversal in linear time. The worker argument for the abstraction from the list constructor `[]` is used as an accumulator.

The list abstraction algorithm with inlining cannot achieve this transformation of the quadratic version into the linear version. To abstract the intermediate list, that algorithm would inline the definition of `reverse`. Then the intermediate list would be eliminated successfully, but the inlined definition of `reverse` would contain a new starting point for deforestation which would lead to new inlining of `reverse` ... The quadratic version creates at run time an intermediate list between each recursive call. To remove all these intermediate lists through a finite amount of transformation the worker/wrapper scheme is required.

6.1 Worker Derivation with Polymorphic Recursion

Unfortunately, the worker `reverseW` cannot be derived by the algorithm described in Section 5. Compare the recursive definition of `reverseW` with the recursive definition of `tailsW`. The former is polymorphically recursive, that is, a recursive call uses type arguments different from the abstracted type variables. Obviously, functions that consume their own result need such polymorphically recursive workers.

Typability in the Damas-Milner type system with polymorphic recursion is semi-decidable, that is, there are algorithms which do infer the most general type of an expression within the Damas-Milner type system with polymorphic recursion if it is typable and may diverge otherwise. Fortunately, the input of the worker/wrapper split algorithm is typable, we only try to find a more general type than we have.

To derive a possibly polymorphically recursive worker definition, we build on Mycroft's extension of the Damas-Milner type inference algorithm [Myc84]. We start with the most general worker type possible, which is obtained from the original type by replacing every list type by a new type variable and abstracting from it and respective list constructors.

$$\begin{aligned} \{ &\text{reverseW}: \forall \alpha. \forall \gamma_1. \forall \gamma_2. (\alpha \rightarrow \gamma_1 \rightarrow \gamma_1) \rightarrow \gamma_1 \rightarrow (\alpha \rightarrow \gamma_2 \rightarrow \gamma_2) \rightarrow \gamma_2 \rightarrow (\gamma_1 \rightarrow \gamma_2), \\ &\mathbf{n}_1: \gamma_3, \mathbf{n}_2: \gamma_4, \mathbf{n}_3: \gamma_5, \mathbf{n}_4: \gamma_6, \mathbf{n}_5: \gamma_7, \\ &\mathbf{c}_1: \alpha \rightarrow \gamma_8 \rightarrow \gamma_8, \mathbf{c}_2: \alpha \rightarrow \gamma_9 \rightarrow \gamma_9, \mathbf{c}_3: \alpha \rightarrow \gamma_{10} \rightarrow \gamma_{10}, \mathbf{c}_4: \alpha \rightarrow \gamma_{11} \rightarrow \gamma_{11} \} \\ \vdash &\lambda xs: \gamma_{12}. \text{case } xs \text{ of } \{ \\ &\quad [] \quad \rightarrow \mathbf{n}_1 \\ &\quad y: ys \rightarrow \text{appW } \alpha \ \gamma_{13} \ \mathbf{c}_1 \ \mathbf{n}_2 \ (\text{reverseW } \alpha \ \gamma_{14} \ \gamma_{15} \ \mathbf{c}_2 \ \mathbf{n}_3 \ \mathbf{c}_3 \ \mathbf{n}_4 \ ys) \ (\mathbf{c}_4 \ y \ \mathbf{n}_5) \} \\ &: \gamma_{16} \end{aligned}$$

We perform type inference to obtain a first approximation of the type of the worker:

$$\begin{aligned}
& \{\mathbf{reverseW}: \forall \alpha. \forall \gamma_1. \forall \gamma_2. (\alpha \rightarrow \gamma_1 \rightarrow \gamma_1) \rightarrow \gamma_1 \rightarrow (\alpha \rightarrow \gamma_2 \rightarrow \gamma_2) \rightarrow \gamma_2 \rightarrow (\gamma_1 \rightarrow \gamma_2), \\
& \quad \mathbf{n}_1: \gamma, \mathbf{n}_2: \gamma, \mathbf{n}_3: [\alpha], \mathbf{n}_4: [\alpha], \mathbf{n}_5: \gamma, \\
& \quad \mathbf{c}_1: \alpha \rightarrow \gamma \rightarrow \gamma, \mathbf{c}_2: \alpha \rightarrow [\alpha] \rightarrow [\alpha], \mathbf{c}_3: \alpha \rightarrow [\alpha] \rightarrow [\alpha], \mathbf{c}_4: \alpha \rightarrow \gamma \rightarrow \gamma\} \\
& \vdash \lambda \mathbf{xs}: [\alpha]. \text{case } \mathbf{xs} \text{ of } \{ \\
& \quad [] \rightarrow \mathbf{n}_1 \\
& \quad \mathbf{y}: \mathbf{ys} \rightarrow \text{appW } \alpha \ \gamma \ \mathbf{c}_1 \ \mathbf{n}_2 \ (\mathbf{reverseW} \ \alpha \ [\alpha] \ [\alpha] \ \mathbf{c}_2 \ \mathbf{n}_3 \ \mathbf{c}_3 \ \mathbf{n}_4 \ \mathbf{ys}) \ (\mathbf{c}_4 \ \mathbf{y} \ \mathbf{n}_5) \} \\
& : [\alpha] \rightarrow \gamma
\end{aligned}$$

Subsequently we infer anew the type of the definition body, this time under the assumption that `reverseW` has the type $\forall \alpha. \forall \gamma. (\alpha \rightarrow \gamma \rightarrow \gamma) \rightarrow \gamma \rightarrow [\alpha] \rightarrow \gamma$, the result of the first type inference. This process iterates until the inferred type is stable, that is input and output type are identical. For our example the second iteration already shows that the result of the first iteration is correct. In general, worker derivation stops latest after $n + 1$ iterations, where n is the number of list types in the type of the original function.

6.2 Further Workers with Polymorphic Recursion

Similar to the example `reverse` are definitions of functions which traverse a tree to collect all node entries in a list. A straightforward quadratic time definition can be split into a polymorphically recursive worker and a wrapper and then be deforested to obtain a linear time definition which uses an accumulating argument.

A different, fascinating example is the definition of the function `inits`, which determines the list of initial segments of a list with the shortest first.

$$\begin{aligned}
& \mathbf{inits}: \forall \alpha. [\alpha] \rightarrow [[\alpha]] \\
& = \lambda \alpha. \lambda \mathbf{xs}: [\alpha]. \\
& \quad \text{case } \mathbf{xs} \text{ of } \{ \\
& \quad \quad [] \rightarrow (:) [\alpha] ([]) \alpha ([]) [\alpha] \\
& \quad \quad \mathbf{y}: \mathbf{ys} \rightarrow (:) [\alpha] ([]) \alpha (\text{map } [\alpha] \ [\alpha] ((:) \ \alpha \ \mathbf{y}) (\mathbf{inits} \ \alpha \ \mathbf{ys})) \}
\end{aligned}$$

It is split into the following polymorphically recursive worker and wrapper:

$$\begin{aligned}
& \mathbf{initsW}: \forall \alpha. \forall \gamma_1. \forall \gamma_2. (\alpha \rightarrow \gamma_1 \rightarrow \gamma_1) \rightarrow \gamma_1 \rightarrow (\gamma_1 \rightarrow \gamma_2 \rightarrow \gamma_2) \rightarrow \gamma_2 \rightarrow [\alpha] \rightarrow \gamma_2 \\
& = \lambda \alpha. \lambda \gamma_1. \lambda \gamma_2. \lambda \mathbf{c}_1: \alpha \rightarrow \gamma_1 \rightarrow \gamma_1. \lambda \mathbf{n}_1: \gamma_1. \lambda \mathbf{c}_2: \gamma_1 \rightarrow \gamma_2 \rightarrow \gamma_2. \lambda \mathbf{n}_2: \gamma_2. \lambda \mathbf{xs}: [\alpha]. \\
& \quad \text{case } \mathbf{xs} \text{ of } \{ \\
& \quad \quad [] \rightarrow \mathbf{c}_2 \ \mathbf{n}_1 \ \mathbf{n}_2 \\
& \quad \quad \mathbf{y}: \mathbf{ys} \rightarrow \mathbf{c}_2 \ \mathbf{n}_1 \ (\text{mapW } \gamma_1 \ \gamma_1 \ \gamma_2 \ \mathbf{c}_2 \ \mathbf{n}_2 \ (\mathbf{c}_1 \ \mathbf{y}) \\
& \quad \quad \quad (\mathbf{initsW} \ \alpha \ \gamma_1 \ [\gamma_1] \ \mathbf{c}_1 \ \mathbf{n}_1 \ ((:) \ \gamma_1) ([]) \ \gamma_1) \ \mathbf{ys})) \}
\end{aligned}$$

$$\begin{aligned}
& \mathbf{inits}: \forall \alpha. [\alpha] \rightarrow [[\alpha]] \\
& = \lambda \alpha. \mathbf{initsW} \ \alpha \ [\alpha] \ [[\alpha]] \ ((:) \ \alpha) ([]) \ \alpha \ ((:) \ [\alpha]) ([]) \ [\alpha]
\end{aligned}$$

Note the abstraction from both (nested) result lists, which cannot be expressed with `build`. Fusion can be performed in the definition body of `initsW`:

$$\begin{aligned}
& \mathbf{initsW}: \forall \alpha. \forall \gamma_1. \forall \gamma_2. (\alpha \rightarrow \gamma_1 \rightarrow \gamma_1) \rightarrow \gamma_1 \rightarrow (\gamma_1 \rightarrow \gamma_2 \rightarrow \gamma_2) \rightarrow \gamma_2 \rightarrow [\alpha] \rightarrow \gamma_2 \\
& = \lambda \alpha. \lambda \gamma_1. \lambda \gamma_2. \lambda \mathbf{c}_1: \alpha \rightarrow \gamma_1 \rightarrow \gamma_1. \lambda \mathbf{n}_1: \gamma_1. \lambda \mathbf{c}_2: \gamma_1 \rightarrow \gamma_2 \rightarrow \gamma_2. \lambda \mathbf{n}_2: \gamma_2. \lambda \mathbf{xs}: [\alpha]. \\
& \quad \text{case } \mathbf{xs} \text{ of } \{ \\
& \quad \quad [] \rightarrow \mathbf{c}_2 \ \mathbf{n}_1 \ \mathbf{n}_2 \\
& \quad \quad \mathbf{y}: \mathbf{ys} \rightarrow \mathbf{c}_2 \ \mathbf{n}_1 \ (\mathbf{initsW} \ \alpha \ \gamma_1 \ \gamma_2 \ (\lambda \mathbf{v}: \gamma_1. \lambda \mathbf{w}: \gamma_2. \mathbf{c}_2 \ (\mathbf{c}_1 \ \mathbf{y} \ \mathbf{v}) \ \mathbf{w}) \ \mathbf{n}_2 \ \mathbf{ys})) \}
\end{aligned}$$

The n -queens function as defined in [Gil96], Section 5.1, is another example in the same spirit.

6.3 Inaccessible Recursive Arguments

Unfortunately, a function may consume its own result but not be defined recursively. For example, the function `reverse` should actually be defined in terms of `foldr`, to enable short cut deforestation with `reverse` as consumer.

```
reverse: ∀α. [α] → [α]
        = λα. foldr α [α] (λy:α.λr:[α]. (++) α r ((:) α y ([] α))) ([] α)
```

Here type inference with polymorphic recursion cannot help. To enable list abstraction we rewrite the definition as follows:

```
reverse: ∀α. [α] → [α]
        = λα. foldr α [α]
          (λy:α.λr:(α → [α] → [α]) → [α] → [α].
            (++) α (r ((:) α) ([] α)) ((:) α y ([] α)))
          (λc:α → [α] → [α].λn:[α]. n)
          ((:) α)
          ([] α)
```

It is, however, unclear when and how such a lifting of the result type of a function that encapsulates recursion can be done in general.

6.4 Deforestation Changes Complexity

Deforestation of the definition of `reverse` changes its complexity from quadratic to linear time. In case of the definition of `inits`, the change of complexity is more subtle. Both the original definition and the deforested definition take quadratic time to produce their complete result. However, to produce only the outer list of the result, with computation of the list elements still suspended, the original definition still takes quadratic time whereas the deforested version only needs linear time.

A polymorphically recursive worker will nearly always enable deforestation which changes the asymptotic time complexity of a function definition. This power is, however, a double-edged sword. A small syntactic change of a program (cf. previous subsection) may cause deforestation to be no longer applicable, and thus change the asymptotic complexity of the program. It can hence be argued that such far-reaching modifications should be left to the programmer.

7 Summary and Future Work

In this paper we presented an expressive worker/wrapper scheme to perform short cut deforestation (nearly) without inlining. An algorithm which is based on our type abstraction algorithm [Chi99] splits all definitions of list-producing functions of a program into workers and wrappers. The wrappers are small enough to be inlined unconditionally everywhere, also across module boundaries. They transfer the information needed for type abstraction in the split algorithm and the actual deforestation algorithm.

The actual deforestation algorithm itself searches for occurrences of `foldr`, abstracts the result list from the producer and then directly applies the short cut fusion rule. Further optimisations may be obtained by a subsequent standard inlining pass.

The deforestation algorithm is separate from the worker/wrapper split algorithm. The algorithms may be integrated, but we intend to perform deforestation after worker/wrapper splitting. Because deforestation and other optimisations may lead to new deforestation opportunities, it is useful to repeat deforestation several times.

Finally, we studied functions which consume their own result. Their definitions can automatically be deforested if the split algorithm is extended on the basis of Mycroft’s extension of Damas-Milner type inference to polymorphic recursion. Nonetheless they still raise interesting questions.

A Worker/Wrapper Scheme for Consumers We focused on how to derive a producer for short cut deforestation without requiring large-scale inlining. Dually the consumer must be a `foldr` and hence sufficient inlining must be performed in the consumer to expose the `foldr`. If the arguments of the `foldr` are large expressions, the standard inliner will refuse to inline the `foldr` expression. So it seems reasonable to also split consumers into `foldr` wrappers and separate workers for the arguments of `foldr`. This transformation, however, does not require any (possibly type-based) analysis but can be performed directly on the syntactic structure.

Other Data Types The wrapper/worker scheme and the type-inference based derivation of workers is not specific to lists but can be used for arbitrary algebraic data types. In fact, a single worker can abstract from the data constructors of several different data types simultaneously and type inference can derive all abstractions simultaneously as well. For fusion we additionally need a catamorphism, that is a `foldr`, for the intermediate data type. This catamorphism could be provided by the user through a compiler directive or be inferred automatically by other algorithms.

Integration with Type Inference Algorithm The worker/wrapper split algorithm is not as efficient as it could be. The list abstraction algorithm traverses a whole definition body once. Even if we ignore polymorphic recursion, if n `let` bindings are nested, then the body of the inner definition is traversed n times.

However, as stated in Section 2, the list abstraction algorithm uses a modified version of the Damas-Milner type inference algorithm. The abstraction from type variables that replace list types corresponds to the generalisation step of the Damas-Milner algorithm. The list abstraction algorithm just additionally abstracts from term variables that replace list constructors and inserts both type and term abstractions into the program. The Damas-Milner algorithm recursively traverses a program only once. So we plan to integrate explicit type and term abstraction at `let` bindings into this type inference algorithm to obtain a single pass split algorithm. To deal with polymorphic recursion as well, the type inference algorithm of Emms and Leiß, which integrates semiunification into the Damas-Milner algorithm, may provide a good basis [EL99].

Implementation We have a working prototype of the list abstraction algorithm with inlining. On this basis we are implementing a simple worker/wrapper split algorithm. The final goal is an implementation in the Glasgow Haskell compiler to apply type-inference based short cut deforestation to real-world programs.

Acknowledgements

I thank Simon Peyton Jones for several comments that inspired this paper. Especially, he drew my attention to producers that consume their own result.

References

- [Chi99] Olaf Chitil. Type inference builds a short cut to deforestation. In *ICFP'99, International Conference on Functional Programming*. ACM Press, 1999.
- [DM82] L. Damas and R. Milner. Principal type-schemes for functional programs. In *Conference Record of the Ninth Annual ACM Symposium on Principles of Programming Languages*, pages 207–212. ACM Press, January 1982.
- [EL99] Martin Emms and Hans Leiß. Extending the type checker of Standard ML by polymorphic recursion. *Theoretical Computer Science*, 212(1–2):157–181, February 1999.
- [GHC] The Glasgow Haskell compiler. Available from <http://research.microsoft.com/users/t-simonm/ghc/>.
- [Gil96] Andrew Gill. *Cheap Deforestation for Non-strict Functional Languages*. PhD thesis, Glasgow University, 1996.
- [GLP93] Andrew Gill, John Launchbury, and Simon L. Peyton Jones. A Short Cut to Deforestation. In *FPCA'93, Conference on Functional Programming Languages and Computer Architecture*, pages 223–232. ACM Press, 1993.
- [Hug89] J. Hughes. Why Functional Programming Matters. *Computer Journal*, 32(2):98–107, 1989.
- [LY98] Oukseh Lee and Kangkeun Yi. Proofs about a folklore let-polymorphic type inference algorithm. *ACM Transactions on Programming Languages and Systems*, 20(4):707–723, July 1998.
- [Myc84] A. Mycroft. Polymorphic type schemes and recursive definitions. In M. Paul and B. Robinet, editors, *Proceedings of the International Symposium on Programming*, LNCS 167, pages 217–228, Toulouse, France, April 1984. Springer.
- [PH⁺99] Simon L. Peyton Jones, John Hughes, et al. Haskell 98: A non-strict, purely functional language. <http://www.haskell.org>, February 1999.
- [PM99] Simon L. Peyton Jones and Simon Marlow. Secrets of the Glasgow Haskell compiler inliner. Submitted to IDL'99, July 1999.