Refactoring Pattern Matching

Meng Wang^a, Jeremy Gibbons^b, Kazutaka Matsuda^c, Zhenjiang Hu^d

^a Computer Science and Engineering, Chalmers University of Technology 412 96 Göteborg, Sweden
 ^b Department of Computer Science, Oxford University Wolfson Building, Parks Road, Oxford OX1 3QD, UK
 ^c Graduate School of Information Sciences, Tohoku University Aramaki aza Aoba 6-3-09, Aoba-ku, Sendai-city, Miyagi-pref. 980-8579, Japan
 ^d GRACE Center, National Institute of Informatics 2-1-2 Hitotsubashi, Chiyoda-ku, Tokyo 101-8430, Japan

Abstract

Defining functions by pattern matching over the arguments is advantageous for understanding and reasoning, but it tends to expose the implementation of a datatype. Significant effort has been invested in tackling this loss of modularity; however, decoupling patterns from concrete representations while maintaining soundness of reasoning has been a challenge. Inspired by the development of invertible programming, we propose an approach to program refactoring based on a right-invertible language RINV—every function has a right (or pre-) inverse. We show how this new design is able to permit a smooth incremental transition from programs with algebraic datatypes and pattern matching, to ones with proper encapsulation, while maintaining simple and sound reasoning.

Keywords: functional programming, refactoring, pattern matching, invertible programming, abstract datatypes, fusion

1. Introduction

1.1. Program Development

Suppose that you are developing a program involving some data structure. You don't yet know which operations you will need on the data structure, nor what efficiency constraints you will impose on those operations. Instead, you want to prototype the program, and conduct some initial experiments on the prototype; on the basis of the results from those experiments, you will decide

whether a naive representation of the data structure suffices, or whether you need to choose a more sophisticated implementation. In the latter case, you do not want to have to conduct major surgery on your prototype in order to refactor it to use a different representation.

The traditional solution to this problem is to use data abstraction: identify (or evolve) an interface, program to that interface, and allow the implementation of the interface to vary without perturbing the program. However, that requires you to prepare in advance for the possible change of representation: it doesn't provide a smooth revision path if you didn't have the foresight to introduce the interface in the first place, but used a bare algebraic datatype as the representation.

Moreover, choosing a naive representation in terms of an algebraic datatype has considerable attractions. Programs that manipulate the data can be defined using *pattern matching* over the constructors of the datatype, rather than having to use 'observer' operations on a data abstraction. This leads to a concise and elegant programming style, which—being based on equations—is especially convenient for reasoning about program behaviour [1].

1.2. Pattern Matching

As a simple example, consider encoding binary numbers as lists of bits, most significant first:

```
data Bin = Zero \mid One

type Num = [Bin]
```

The above declarations introduce a new datatype (Bin) for binary bits, and define Num as a type synonym for lists of bits. Throughout this paper, we will use the syntax and standard prelude functions of Haskell [2] for illustration, although any language providing algebraic datatypes would work just as well; however, in contrast to Haskell, we assume a semantics based on sets and total functions rather than on complete partial orders.

Consider this function for normalizing binary numbers by eliding leading zeroes, defined by pattern matching.

The definition forms a collection of equations, which give a straightforward explanation of the operational behaviour of the function:

```
normalize [Zero, One, Zero]
\equiv \{ Clause (3) \}
normalize [One, Zero]
\equiv \{ Clause (2) \}
[One, Zero]
```

They are also convenient for reasoning; for example, here is one case of an inductive proof that *normalize* is idempotent:

```
normalize (normalize (Zero: num))

≡ { Clause (3) }
normalize (normalize num)

≡ { inductive hypothesis }
normalize num

≡ { Clause (3) }
normalize (Zero: num)
```

An equivalent definition without using pattern matching is harder to read:

```
normalize :: Num \rightarrow Num
normalize \ num = \mathbf{if} \ null \ num \lor one \ (head \ num) \ \mathbf{then}
num
\mathbf{else}
normalize \ (tail \ num)
```

It is also much less convenient for calculating with.

Pattern matching has accordingly been supported as a standard feature in most modern functional languages, since its introduction in Hope [3]. More recently, it has started gaining recognition from the object-oriented community too [4, 5, 6]. Unfortunately, the appeal of pattern matching wanes when we need to change the implementation of a data structure: function definitions are tightly coupled to a particular representation, and a change of representation has a far-reaching effect. As a result, it has been observed that the use of pattern matching "leads to a discontinuity in programming: programmers initially use pattern matching heavily, and are then forced to abandon the technique in order to regain abstraction over representations" [7].

In this stand-off between clarity and abstraction, functional languages usually lean towards the former while object-oriented languages prefer the latter. Can we hope to achieve the best of both worlds? With current technology, encapsulating datatypes in a functional language hinders pattern matching. Ad hoc approaches such as allowing user-defined computations to be embedded in the pattern matching and data construction processes threaten soundness of reasoning. For example, we rely on reasoning to perform knowncase elimination that reduces expression **case** ($Cons\ x\ xs$) **of** ($Cons\ a\ as$) $\rightarrow a\ to\ x$ based on the knowledge that the computation constructing ($Cons\ x\ xs$) is reversed by matching the pattern $Cons\ a\ as$. If user-defined computations are embedded in the pattern matching process such as in Wadler's view proposal [8], we no longer have the guarantee that the above mentioned round-trip property still holds, which invalidates reasoning of this kind. To recover soundness, a more disciplined framework is needed.

In an object-oriented setting, there is a concept of typical object [9], which is a simple and general data structure (often expressed in a way similar to an algebraic datatype) used as a model of the underlying more complicated and user-defined structure. For example, various kinds of linear structure can be modelled as a list datatype. The user-defined structures are used for execution while the typical object caters for specification. The function that maps an underlying structure to its typical object is called an abstraction function. An implementation using the underlying structure can be proven correct, with respect to its specification in terms of the typical object [10], by reasoning about the abstraction function. However, reasoning in this manner does not help with incrementally refactoring a program already written with bare algebraic datatypes.

1.3. Our Contribution

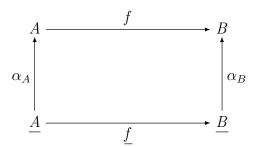
In this work, then, we strive to address the tension between the convenience of pattern matching and the flexibility of encapsulated data abstractions by proposing a mechanism to allow programs written with pattern matching to be refactored smoothly and incrementally into ones with encapsulation, without losing the benefits of simple equational reasoning. In particular, we propose a framework for incrementally refactoring and identify necessary and sufficient conditions for correctness of such refactoring. As part of the framework, we sketch a domain-specific language RINV that helps to guarantee these conditions by construction. We have implemented the

proposed system to demonstrate its feasibility. A prototype of our proposal can be found online [11] and is described in detail in the appendix.

The rest of the paper is structured as follows. Section 2 gives a brief introduction to data abstraction. Section 3 presents our proposed design for incremental refactoring of program with pattern matching, and Section 4 provides a formal definition of the right-invertible language RINV on which our design is based. We then evaluate the performance and explore alternative points in the design space of our system (Section 5), and apply our results to other problems in addition to refactoring (Section 6), before discussing related work (Section 7) and concluding (Section 8).

2. Data Abstraction

Choosing the right data structure is key to achieving an efficient program; data abstraction allows us to defer the choice of representation until after the uses of the data are fully understood. The idea is to firstly program with 'abstract' data, which is then replaced by a more efficient 'concrete' representation. Operations on abstract data are reimplemented on the concrete representation, with the original abstract operations serving as specifications. Milner [12] observed that the transition from abstract to concrete representation can be proven correct by showing the following square commutes,



This proof obligation is also known as the *promotion* condition [13]

$$\alpha_B \circ f \equiv f \circ \alpha_A$$

Throughout this paper, we use \equiv for equality between expressions, whereas we use = for definitions. As a notational convention, we write \underline{X} for the concrete counterpart of an abstract representation X. In the above, the α_X

family of functions are abstraction functions [10] that map the concrete representations to the abstract ones; we will simply write α if its type is clear from the context. In the above diagram, A and B are abstract representations, whereas \underline{A} and \underline{B} are concrete representations. The abstract program is $f: A \to B$, and its concrete version is $\underline{f}: \underline{A} \to \underline{B}$. In the sequel, we use the terms 'abstract program' and 'specification' interchangeably.

2.1. An Example

Consider queue structures, represented abstractly as lists.

```
type Queue a = [a]

emptyQ = []

first = head

isEmpty = null

enQ \ a \ q = q + [a]

deQ = tail
```

(Note that + is a Haskell infix operator that concatenates two lists.) For a more efficient definition of enQ, a plausible concrete representation reads:

type Queue
$$a = ([a], [a])$$

 $\alpha :: \underline{Queue} \ a \to Queue \ a$
 $\alpha \ (fq, bq) = fq + reverse \ bq$

The second list of the pair, representing the latter part of a queue, is reversed, so that enqueuing is simply list prefixing. The library operations can be implemented as follows:

$$\begin{array}{ll} \underline{emptyQ} & = ([],[]) \\ \underline{first} \; ([],bq) & = last \; bq \\ \underline{first} \; ((a:fq),bq) = a \\ \underline{isEmpty} \; ([],[]) & = True \\ \underline{isEmpty} \; q & = False \\ \underline{enQ} \; a \; (fq,bq) & = (fq,a:bq) \\ \underline{deQ} \; ([],bq)) & = \underline{deQ} \; (reverse \; bq,[]) \\ \underline{deQ} \; (a:fq,bq) & = \overline{(fq,bq)} \end{array}$$

Through standard equational reasoning, we can establish the correctness of the implementation by proving the promotion condition for each operation. For example,

```
\alpha \left( \underline{deQ} \left( a:fq,bq \right) \right)
\equiv \left\{ \begin{array}{l} \operatorname{definition of } \underline{deQ} \right\} \\ \alpha \left( fq,bq \right) \\ \end{array}
\equiv \left\{ \begin{array}{l} \operatorname{definition of } deQ \right\} \\ deQ \left( a:\alpha \left( fq,bq \right) \right) \\ \end{array}
\equiv \left\{ \begin{array}{l} \operatorname{definition of } \alpha \right\} \\ deQ \left( a:\left( fq + reverse \ bq \right) \right) \\ \end{array}
\equiv \left\{ \begin{array}{l} \operatorname{definition of } + \right\} \\ deQ \left( \left( a:fq \right) + reverse \ bq \right) \\ \end{array}
\equiv \left\{ \begin{array}{l} \operatorname{definition of } \alpha \right\} \\ deQ \left( \alpha \left( a:fq,bq \right) \right) \\ \end{array}
```

Other than the handful of 'library' operations above, we have other abstract programs benefitting from the simple list representation. For example, the map function on queues:

```
mapQ :: (a \rightarrow b) \rightarrow Queue \ a \rightarrow Queue \ b

mapQ \ f \ [] = []

mapQ \ f \ (a : q) = f \ a : mapQ \ f \ q
```

or a prioritisation function, which is essentially a stable sort based on element weight:

```
prioritise :: Ord a \Rightarrow Queue \ a \rightarrow Queue \ a
prioritise [] = []
prioritise (a:q) = insert \ a \ (prioritise \ q)
where insert \ b \ [] = [b]
insert \ b \ (a:q) = if \ b \leqslant a \ then \ b : (a:q)
else a:(insert \ b \ q)
```

To maintain executability, all uses of the abstract representation have to be changed at once, even though some of the old definitions may not gain from the refactoring. One has to (re)implement all the functions either by pattern matching on the new representation (discouraging further refactoring) or by the use of library operations (at the cost of losing convenient equational reasoning). For example, a definition of map using only the library operations is likely to be more clumsy:

```
mapQ_{prim} :: (a \rightarrow b) \rightarrow Queue \ a \rightarrow Queue \ b
mapQ_{prim} f \ q = accum \ f \ q \ emptyQ
where accum \ f \ q \ aq = if \ isEmpty \ q \ then
aq
else
accum \ f \ (deQ \ q) \ (enQ \ (f \ (first \ q)) \ aq)
```

Being restricted to an explicit interface reduced our options here: it is not easy to add to the front of a queue, for which we need an accumulating parameter.

In the next section, we propose a framework free from the above pitfalls: refactoring can be done selectively; and at any point in the process, executability and reasoning are fully supported. We look into the details of the design by means of examples.

3. Selective Refactoring

Our purpose is to allow incremental refactoring of a program, replacing a specification in stages by a more sophisticated implementation. We have seen in the previous section that once an implementation is given, the promotion condition is sufficient to guarantee its correctness. However, the promotion condition does not suggest a way of integrating the new definition into its original context, which may still operate on the old representation. We need a computation law of the form

$$f \; \equiv \; \alpha \circ f \circ \alpha^{\circ}$$

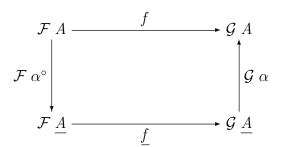
in order to replace specification f with its implementation $\alpha \circ f \circ \alpha^{\circ}$.

Here, the function α° does the 'opposite' of α , mapping an abstract value down to a concrete value. The abstract program f is replaced by the composition $\alpha \circ \underline{f} \circ \alpha^{\circ}$, which converts an abstract value to the concrete representation, executes the concrete program, then converts the concrete result back to the abstract representation. For example, the abstract program fragment deQ x is replaced with $(\alpha \circ \underline{deQ} \circ \alpha^{\circ})$ x, where the two functions α° and α convert a value of type \underline{Queue} into one of type \underline{Queue} and back again.

It is obvious that if α and α° are each other's inverses (a situation sometimes called *strong simulation* [12]), the computation law is equivalent to the promotion condition. However, requiring them to be inverses restricts the abstraction function to be bijective, which is rather a strong condition. In fact, one-sided invertibility is sufficient in this case. Specifically, the α° function should be a right inverse of α – that is, $\alpha \circ \alpha^{\circ} \equiv id$; it is not necessary for α° also to be a left inverse of α .

3.1. The Computation Law

Before going into the details of our results, we firstly generalise the notation used. The abstract operation f will not always have a type as simple as $Queue\ a \to Queue\ a$, like deQ does. Suppose we have datatypes A and \underline{A} , and conversion functions $\alpha::\underline{A}\to A$ and $\alpha^\circ::A\to \underline{A}$. In the general case, an abstract operation will take not just a single value in the abstract representation (of type A), but some combination of abstract values and other arguments. We capture this combination in terms of an operation $\mathcal F$ on datatypes A. Similarly, the operation will return a different combination $\mathcal G$ of abstract values and other results. The operations $\mathcal F$ and $\mathcal G$ are functors; they lift the conversion functions in the obvious way to $\mathcal F$ $\alpha^\circ::\mathcal F$ $A\to \mathcal F$ A and $\mathcal G$ $\alpha::\mathcal G$ $A\to \mathcal G$ A, acting pointwise on combinations of values, and respecting identity and composition. Then an abstract operation f will have type $\mathcal F$ $A\to \mathcal G$ A, and the corresponding concrete operation f:: $\mathcal F$ $A\to \mathcal G$ A should satisfy the computation law $f\equiv \mathcal G$ $\alpha\circ f\circ \mathcal F$ α° , as shown in the following commuting diagram.



For example, for the operation first: Queue $a \to a$, the input context \mathcal{F} is the identity functor, because the source type Queue a of first consists of a single occurrence of the datatype in question, and the output context \mathcal{G} is a constant functor, because the target type a of first has no occurrences of the datatype. The operation must satisfy the following computation law:

$$first \equiv first \circ \alpha^{\circ}$$

The computation equations for the rest of the operations are listed below.

```
isEmpty \equiv \underline{isEmpty} \circ \alpha^{\circ}
emptyQ \equiv \alpha \circ \underline{emptyQ}
enQ \ a \equiv \alpha \circ \underline{enQ} \ a \circ \alpha^{\circ}
deQ \equiv \alpha \circ \overline{deQ} \circ \alpha^{\circ}
```

The computation law provides a way of replacing a specification with its implementation. But since it involves both α and α° , it seems more difficult to prove than the promotion condition. However, the following result allows us to derive the computation law from the promotion condition, using the right-invertibility property of α and α° .

Lemma 1 (Computation). Given operation specifications $f :: \mathcal{F} A \to \mathcal{G} A$, and their implementations $\underline{f} :: \mathcal{F} \underline{A} \to \mathcal{G} \underline{A}$, for all conversion functions $\alpha :: \underline{A} \to A$ satisfying the promotion condition $\overline{\mathcal{G}} \alpha \circ \underline{f} \equiv f \circ \mathcal{F} \alpha$, and their right inverses $\alpha^{\circ} :: A \to A$, we can deduce the computation law $f \equiv \mathcal{G} \alpha \circ f \circ \mathcal{F} \alpha^{\circ}$.

Proof.

```
\mathcal{G} \ \alpha \circ \underline{f} \circ \mathcal{F} \ \alpha^{\circ}
\equiv \{ \text{promotion } \} 
f \circ \mathcal{F} \ \alpha \circ \mathcal{F} \ \alpha^{\circ}
\equiv \{ \mathcal{F} \text{ respects composition } \} 
f \circ \mathcal{F} \ (\alpha \circ \alpha^{\circ})
\equiv \{ \alpha \circ \alpha^{\circ} \equiv id \} 
f \circ \mathcal{F} \ id
\equiv \{ \mathcal{F} \text{ respects identity } \} 
f
```

Applying the computation law requires additional infrastructure (the right inverse α° of the abstraction function and its properties) in addition to that needed for a data abstraction framework based on the promotion condition. To minimise this extra obligation, we will explore a correctness-by-construction technique: in the next section, we will present a combinator-based language RINV implemented as a library, in which every definable function gets a right inverse for free. That is to say, the programmer carrying

out a refactoring writes only α , in RINV, and the corresponding α° is automatically generated. For the sake of completeness, in the case of the queue example presented above, a possible definition in RINV reads:

$$\alpha = app \circ (id \times reverse)$$

The details of the language are completely orthogonal to the discussion in this section, and can be safely ignored for the time being.

3.2. Refactoring by Translation

In our proposal, a programmer wishing to migrate to the concrete representation has the choice of keeping the original abstract definitions, or of refactoring them into the style using only the library operations, or of having a mixture of the two. For example, consider a queue that is read circularly for a certain amount of time, say repeatedly playing a piece of music.

$$\begin{array}{ccc} play :: Time \rightarrow Queue \; (IO \; ()) \rightarrow IO \; () \\ play \; 0 & (a: _) = a \\ play \; (n+1) \; (a:q) = \mathbf{do} \; a \\ & play \; n \; (q + [a]) \end{array}$$

We use Haskell's **do** notation above to sequence IO actions: the action a will be followed by the actions of play n (q + [a]). To refactor this function, we firstly have to rewrite it with library operations.

$$play_{\text{prim}} :: Time \rightarrow Queue \ (IO \ ()) \rightarrow IO \ ()$$
 $play_{\text{prim}} \ 0 \qquad q = first \ q$
 $play_{\text{prim}} \ (n+1) \ q = \mathbf{do} \ hd$
 $play_{\text{prim}} \ n \ (enQ \ hd \ tl)$
 $\mathbf{where} \ hd = first \ q$
 $tl = deQ \ q$

There is no magic here. The programmer carrying out the refactoring has to write $play_{\text{prim}}$ and be responsible for its correctness: $play_{\text{prim}} \equiv play$.

The library operations are implemented using the concrete representation, allowing constant time performance of enQ.

$$\frac{play_{\text{prim}}}{play_{\text{prim}}} :: Time \to \underline{Queue} \ (IO \ ()) \to IO \ ()$$

$$= \underline{first} \ q$$

$$\frac{play_{\text{prim}}}{(n+1)} (n+1) q = \mathbf{do} hd$$

$$\frac{play_{\text{prim}}}{n} n (\underline{enQ} hd tl)$$

$$\mathbf{where} hd = \underline{first} q$$

$$tl = \overline{deQ} q$$

Re-implementing every function to its "prim" version is certainly laborious, and is not necessarily beneficial. We perform selective refactoring, requiring $\underline{play_{\text{prim}}}$ to do the job of play in the original abstract context. For example, we may want to use $\underline{play_{\text{prim}}}$ together with mapQ (defined in Section 2.1) in expressions like:

$$play_{prim} \ n \circ map Q f$$

The gap between different representations used by $\underline{play_{\text{prim}}}$ and mapQ can be bridged by a mechanical translation, following a rather straightforward scheme: each use of a library operation is replaced with its implementation precomposed with α and postcomposed with α° (subject to the appropriate functors).

The library operations that are defined on the concrete implementation require their inputs to be converted from the abstract representation before consumption, and the outputs converted back to the abstract representation. Effectively, all the refactored functions have the abstract representation as input and output types; the concrete representation remains only for intermediate structures. As an example, $\underline{play_{\text{prim}}}$ can be translated into the following.

$$\begin{array}{ll} play' :: Time \rightarrow Queue \; (IO \; ()) \rightarrow IO \; () \\ play' \; 0 \; q &= (\alpha \circ \underline{first} \circ \alpha^{\circ}) \; q \\ play' \; (n+1) \; q = \mathbf{do} \; hd \\ & play' \; n \; ((\alpha \circ \underline{enQ} \; hd \circ \alpha^{\circ}) \; tl) \\ \mathbf{where} \; hd = (\alpha \circ \underline{first} \circ \alpha^{\circ}) \; q \\ & tl \; = (\alpha \circ \underline{deQ} \circ \alpha^{\circ}) \; q \end{array}$$

Given the computation law, it is easy to conclude that play' is equivalent to $play_{\text{prim}}$, in the sense that exactly the same output is produced for each input.

The original abstract program such as the definition of play is turned into a specification and can be used for equational reasoning. For example, one can continue to show

$$play \ n \ xs \equiv play \ n \ (xs + xs)$$

on the abstract level, without worrying about the refactoring of play into play'.

3.3. Optimization

The above translation of $\underline{play_{\text{prim}}}$ is semantically correct, but inefficient. There are conversions everywhere in the program, which are unnecessary because there is no reference to the abstract representation by the function at all. One way to remove the redundant conversions is fusion. Consider the expression

$$(\alpha \circ enQ \ hd \circ \alpha^{\circ}) \ ((\alpha \circ deQ \circ \alpha^{\circ}) \ q)$$

Our target is to fuse the intermediate conversions to produce

$$(\alpha \circ enQ \ hd \circ deQ \circ \alpha^{\circ}) \ q$$

This would clearly follow from $\alpha^{\circ} \circ \alpha \equiv id$, but this is not a property that we guarantee—for good reason, since requiring it in addition to the existing right inverse property $\alpha \circ \alpha^{\circ} \equiv id$ entails isomorphic abstract and concrete representations, which is too restrictive to be practically useful. Instead, using the promotion condition, we can prove a weaker property that is sufficient for fusion.

Theorem 1 (Fusion Soundness). Given operation specifications $f::\mathcal{F} A \to \mathcal{G} A$ and $g::\mathcal{G} A \to \mathcal{H} A$, and their implementations $\underline{f}::\mathcal{F} \underline{A} \to \mathcal{G} \underline{A}$ and $\underline{g}::\mathcal{G} \underline{A} \to \mathcal{H} \underline{A}$, for all conversion functions $\alpha::\underline{A} \to \overline{A}$ satisfying the promotion condition $\mathcal{H} \alpha \circ \underline{g} \equiv g \circ \mathcal{G} \alpha$, and their right inverses $\alpha^{\circ}::A \to \underline{A}$, we have the following fusion law:

$$\mathcal{H}\;\alpha\,\circ\,g\,\circ\,\mathcal{G}\;\alpha^{\circ}\,\circ\,\mathcal{G}\;\alpha\,\circ\,f\,\circ\,\mathcal{F}\;\alpha^{\circ}\,\equiv\,\mathcal{H}\;\alpha\,\circ\,g\,\circ\,f\,\circ\,\mathcal{F}\;\alpha^{\circ}$$

Proof.

$$\mathcal{H} \alpha \circ \underline{g} \circ \mathcal{G} \alpha^{\circ} \circ \mathcal{G} \alpha \circ \underline{f} \circ \mathcal{F} \alpha^{\circ}$$

$$\equiv \{ \text{promotion} \}$$

$$g \circ \mathcal{G} \alpha \circ \mathcal{G} \alpha^{\circ} \circ \mathcal{G} \alpha \circ \underline{f} \circ \mathcal{F} \alpha^{\circ}$$

$$\equiv \{ \mathcal{G} \text{ is a functor; } \alpha \circ \alpha^{\circ} \equiv id \}$$

$$\begin{array}{l} g \, \circ \, \mathcal{G} \, \, \alpha \, \circ \, \underline{f} \, \circ \, \mathcal{F} \, \, \alpha^{\circ} \\ \equiv \, \quad \big\{ \, \text{promotion} \, \big\} \\ \mathcal{H} \, \, \alpha \, \circ \, \underline{g} \, \circ \, \underline{f} \, \circ \, \mathcal{F} \, \, \alpha^{\circ} \end{array}$$

Basically, this theorem states that although the input to \underline{g} may differ from the output of \underline{f} , due to the $\alpha^{\circ} \circ \alpha$ conversions, nevertheless the post-conversion of \underline{g} 's output brings possibly different results into the same value in the abstract representation.

It is clear theoretically that repeatedly applying the above fusion law will eliminate all intermediate conversions until the point where the abstract representation is used. The absence of abstract representation is reflected in typing: for example the function $\underline{play_{\text{prim}}}$ is well-typed without translation, as we never have to unify the types \overline{Queue} and \underline{Queue} . In this case, $\underline{play_{\text{prim}}}$ is no different from library functions such as \underline{deQ} translation-wise, where we lift the conversion out of the recursion. As a result, our implementation directly produces

$$play'\ n = \alpha \circ play_{\text{prim}}\ n \circ \alpha^{\circ}$$

which is free from intermediate conversions.

It is also clear that not all conversions can be eliminated. In expressions like

$$\alpha \circ \underline{play_{\text{prim}}} \ n \circ \alpha^{\circ} \circ mapQ \ f$$

some performance overhead is inevitable. Such overheads may or may not be acceptable, depending on the situation. The programmer has the option of refactoring mapQ as well to remove the conversion. We won't discuss this engineering trade-off further: although we advocate selective refactoring, we are not necessarily against complete refactoring; we simply provide the options for programmers to choose from.

3.4. More Examples

We now look at a few additional examples making use of the refactoring framework developed in this section.

3.4.1. Join Lists

Suppose we have the following definition of the reverse function:

```
reverse :: [a] \rightarrow [a]

reverse [] = []

reverse (x : xs) = append (reverse xs) [x]
```

It is well known that the above definition suffers from poor run-time performance due to the linear-time left-biased-list concatenation (append). One way to resolve such inefficiency is to transform individual concatenation-intensive definitions to reduce the calls to append, a technique that is systematised by Wadler [14]. Here we take a different approach by refactoring append itself based on a different representation of list, instead of trying to eliminate its usage.

As an alternative to the biased linear list structure, the *join* representation of lists has been proposed for program elegance [15, 16], efficiency [17], and more recently, parallelism [18]. It can be defined as:

$$\mathbf{data} \ \mathit{List} \ a = \mathit{Empty} \mid \mathit{Unit} \ a \mid \mathit{Join} \ (\mathit{List} \ a) \ (\mathit{List} \ a)$$

With this representation, a constant-time append function can be defined as

append
$$l1 l2 = Join l1 l2$$

At the same time, we don't want to give up on the familiar notions of [] and (:). Instead, they can serve as a specification of the join representation.

```
type List \ a = [a]

append [] \quad y = ys

append \ (x : xs) \ ys = x : append \ xs \ ys
```

Now suppose a programmer would like to refactor certain list functions to make use of the constant-time append function; she will need to define an abstraction function $\alpha :: \underline{List} \ a \to List \ a$ (a corresponding α° will be automatically constructed by RINV, as we will see shortly in the next section), and verify its correctness by proving the promotion condition. The translation outlined in Section 3.2 will complete the refactoring, by converting calls of append into calls of append. At the same time, many other list functions, including some that are yet to be defined, may still use pattern matching

on, and be manipulated using, the original *List* representation. For example, retrieving the head of a list can be defined as:

```
head(x:xs) = x
```

and calculated in the following manner, oblivious to the refactoring:

```
head\ (append\ (x:xs)\ ys) \equiv head\ (x:append\ xs\ ys) \equiv x
```

As we have seen, selective refactoring provides the option of keeping the widely used left-biased list, while taking advantage of alternative representations when needed. This flexibility frees programmers from the dilemma of committing to a single representation that is unlikely to be appropriate universally. In the sequel, we will see more examples where different representations showing distinct merits in different applications.

3.4.2. Binary Numbers

In the introduction, we showed a representation of binary numbers as lists of digits with the most significant bit first (MSB). This representation is intuitive, and offers good support for most operations; however, for incrementing a number, having the least significant bit (LSB) first is better. In Haskell, it is an idiom to introduce a type synonym to document the different usages of the same type.

```
type LSB = [Bin]

incr :: LSB \rightarrow LSB

incr [] = [One]

incr (Zero : num) = One : num

incr (One : num) = Zero : (incr num)
```

Effectively, in order to use the above definition with any other operations, we need to reverse the MSB representation. However, this implicitness of data representation is risky, as any incorrect usage won't be picked up by a compiler. At the same time, handling the two representations as different types can be cumbersome.

With our proposal, we program with only one representation (Num) and selectively refactor certain operations (such as incr), which effectively eliminates any possibility of misuse. We firstly create a new type for the LSB

representation (in Haskell the **newtype** declaration is a more efficient alternative to **data** when there is exactly one constructor with exactly one field inside it):

```
newtype Num = L Num
```

and implement the increment function with the new representation

```
\frac{incr}{incr} :: \underbrace{Num}_{loc} \to \underbrace{Num}_{loc} = L [One]
\underline{incr}_{loc} (L (Zero : num)) = L (One : num)
\underline{incr}_{loc} (L (One : num)) = L (Zero : (incr_{loc} num))
```

As a result, only one representation is exposed to the programmer, and all the conversions between representations are handled implicitly by the refactoring translation.

3.4.3. Sized Trees

A size-annotated binary tree is suitable for fast indexing, as we can traverse it quickly by not entering any left subtree that has a smaller size than the index.

(Note that the use of nested pairs in SFork's parameter purely is a pragmatic decision, to avoid having to deal with triples in programming the conversion functions.)

```
\begin{array}{ll} \underline{index} :: Int \rightarrow STree \ a \rightarrow a \\ \underline{index} \ 0 \ (SLeaf \ a) &= a \\ \underline{index} \ n \ (SFork \ (s,(lt,rt))) \mid n \geqslant ls &= \underline{index} \ (n-ls) \ rt \\ \mid otherwise = \underline{index} \ n \ lt \\ \\ \mathbf{where} \ ls &= getSize \ lt \\ rs &= getSize \ rt \\ \\ getSize \ Empty &= 0 \\ getSize \ (SLeaf \ \_) &= 1 \\ getSize \ (SFork \ (s, \_)) &= s \end{array}
```

This feature of fast indexing makes sized trees an attractive alternative to lists, when access to elements in the middle is required, with the following specification.

```
index \ 0 \ [x] = x

index \ n \ (\_: xs) = index \ (n-1) \ xs
```

Again, once an abstraction function is defined and verified, we can enjoy smooth transitions between the two representations, and reap the benefit of having both.

4. The Right-Invertible Language RINV

Our approach to refactoring discussed above relies on the existence of a right inverse, α° , of the user-defined conversion function α . One can of course implement both functions α° and α separately and prove their consistency. For example, a definition of α° for the two representations of queues in Section 2.1 reads:

```
\alpha^{\circ} :: [a] \rightarrow ([a], [a])

\alpha^{\circ} ls = (fq, bq)

where n = length ls 'div' 2

fq = take \ n \ ls

bq = reverse \ (drop \ n \ ls)
```

An explicit proof verifying that α° is a right inverse of α will be several times longer than the definition above, and make non-trivial uses of properties of the functions involved. Such proofs have to be hand-crafted, and generally cannot be reused.

Instead, we take a linguistic approach to the problem of generating a definition and a consistency proof for α° , by writing α in a right-invertible language that automatically determines such a right inverse for each function constructed. As a result, the task of defining two functions and proving their consistency is reduced to writing just one of them. This approach to invertible languages is not new. As a matter of fact, there have been a number of proposals for similar purposes [19, 20, 21, 22, 23]. On the other hand, reversible languages are usually domain-specific; it is hard to directly reuse them for different applications other than the ones for which they were designed. Inspired by the literature, in this section we sketch a language,

namely RINV, for our purpose of refactoring pattern matching; with this language, the α in the above example can be programmed as

```
\alpha = app \circ (id \times reverse)
```

and its right inverse is automatically generated.

The language RINV is defined as a combinator library; its syntax is as below. (Non-terminals are indicated in small capitals.)

The language is similar in flavour to the point-free 'algebra of programming' style [24], but with the additional feature that a right inverse can be automatically generated for each function that is defined. As a result, a definition $f :: s \hookrightarrow t$ in RINV actually represents a pair of functions (hence the notation $f :: s \hookrightarrow t$): the forward function $f :: s \to t$, and its right inverse $f :: t \to s$, which together satisfy $f :: t \to s$. For convenience, when clear from context, we do not distinguish between $f :: t \to s$ and its forward reading $f :: t \to s$.

The generated right inverses are intended to be total, so the forward functions have to be surjective; this property holds of the primitive functions, and is preserved by the combinators.

There is an extensible set of primitives (cf. Section 4.1) defining the basic non-terminal building blocks of the language. In principle, any surjective function could be made a primitive in RINV. All primitives are uncurried; this suits an invertible framework, where a clear distinction between 'input' and 'output' is required. For the sake of demonstration, we present a small but representative collection of primitive functions above: swap, assocl, and assocr rearrange the components of an input pair; id is the identity operation; app is the uncurried append function on lists; and fst_g and snd_g are projection functions on pairs – discarding one element of a pair that can be regenerated by applying function g to the other one. As we will show, with just these few primitives we can define many interesting functions.

The set of constructor functions (cf. Section 4.2) is also extensible, via new datatypes. We use lowercase names for the uncurried versions of constructors. In contrast to the primitives, constructor functions are inherently non-surjective, and so require some special treatment that will be discussed shortly. We take advantage of this relaxation of the surjectivity rule to admit other non-surjective functions which are crucial for expressiveness. For example, in addition to the left-biased list constructor *cons*, coming with the usual datatype declaration, we also include its right-biased counterpart *snoc*, which adds an element at the end; it can be defined in Haskell as

$$snoc = \lambda(x, xs) \rightarrow xs + [x]$$

Another additional constructor for lists is wrap, which creates a singleton list.

$$wrap \ x = [x]$$

Although this might seem ad hoc, it is by no means arbitrary. One should only use functions that truly model a different representation of the datatype. For example, *snoc* and *nil* form the familiar backwards representation of lists, while *wrap*, *nil* and the primitive function *app* correspond to the join list representation [16]. As we will show in the sequel, this ability to admit non-surjective functions as constructors functions that do not directly arise from a datatype declaration gives us much freedom in altering structures without having to invent new types.

Since constructor functions are exceptions to the surjectivity rule, we additionally require that lone constructors must be combined with other functions by the 'junc' combinator ∇ , which dispatches to one of two functions according to the result of matching on a sum. When one of the operands of ∇ is surjective, or the two operands cover both constructors of a two-variant datatype, the result is surjective; such pairs of operands are called *jointly surjective* [25]. For example, $nil \nabla cons$ and $nil \nabla id$ are both surjective, but $cons \nabla snoc$ is not. Since ∇ can be nested, this result extends to datatypes with more than two constructors. There are built-in annotations in RINV that group constructors into constructor sets, and are used in static checking of surjectivity. For example, for the language presented above, we will have the following set of constructor groups:

$$\{\{nil, cons\}, \{nil, snoc\}, \{nil, wrap, app\}\}$$

The annotations are provided by the language designer, and extended when adding new constructors, but they do not appear in RINV programs. Constructor functions can be composed with other functions as well, using the

standard function composition combinator \circ , but only to the left: once a non-surjective function appears in a chain of compositions other than in the leftmost position, it is difficult to analyse the exact range of the composition. Both the above requirements can be enforced by a rather straightforward syntactic check. The checking algorithm can be found in Appendix C.

Other than the two already mentioned combinators, \times is the cartesian product of two functions (cf. Section 4.3.1), and $fold_X$ f is the operator for regular structural recursion, which decomposes a structure of type X and replaces the constructors with its body f (cf. Section 4.3.2); very often we omit the subscript X, when it can be understood from the context. In combination with swap, assocl and assocr, \times is able to define all functions that rearrange the components of a tuple, while ∇ is useful in constructing the body of a fold. We don't include Δ , the dual of ∇ , in RINV, because of surjectivity, as will be explained shortly.

With the language RINV, we can state the following property.

Theorem 2 (Right invertibility). Given a function f in RINV consisting of forward and right-inverse functions $[\![f]\!]$ and $[\![f]\!]^{\circ}$, we have $[\![f]\!]^{\circ} \equiv id$.

The correctness of this theorem will become evident by the end of this section, as we discuss in detail the various constructs of RINV and their properties.

4.1. The Primitive Functions

The function id is the identity; functions assocr, assocl and swap manipulate pairs.

```
assocr :: ((a, b), c) \leftrightarrows (a, (b, c))
[assocr] = \lambda((a, b), c) \rightarrow (a, (b, c))
[assocl] \circ = [assocl]
assocl :: (a, (b, c)) \leftrightarrows ((a, b), c)
[assocl] = \lambda(a, (b, c)) \rightarrow ((a, b), c)
[assocl] \circ = [assocr]
swap :: (a, b) \leftrightarrows (b, a)
[swap] = \lambda(a, b) \rightarrow (b, a)
[swap] \circ = [swap]
```

Function app is the uncurried append function, which is not injective. The admission of non-injective functions allows us to break away from the

isomorphism restriction. In this case, the right inverse for *app* is not unique; and the RINV programmer will need to choose the appropriate behaviour for the inverse among multiple variants of the function provided by a language implementation. This user-controlled flexibility is standard in invertible languages that admit non-injective functions [21, 23]. For the purposes of demonstration, we pick the following definition for the right-inverse of *app*.

$$\llbracket app \rrbracket^{\circ} = \lambda xs \rightarrow splitAt \ ((length \ xs + 1) \ 'div' \ 2) \ xs$$

Pure projection functions such as fst and snd are difficult to handle in a total function setting; in a setting based on complete partial orders, one could replace the discarded element with an undefined value (for example, by defining $\llbracket fst \rrbracket^{\circ} = \lambda x \to (x, \bot)$), but when types are sets, no analogue of \bot is available. Besides, although such a definition would satisfy right invertibility, it would not be practically very useful. Therefore, we restrict the use of projection functions to cases in which the discarded values are recoverable, giving rise to the following:

$$\begin{aligned} &\mathit{fst}_g :: (a,b) \leftrightarrows a \\ & [\![\mathit{fst}_q]\!]^\circ = \lambda x \to (x,g\ x) \end{aligned}$$

and

$$snd_g :: (a, b) \leftrightarrows b$$

 $\llbracket snd_g \rrbracket^{\circ} = \lambda y \to (g \ y, y)$

Thus, fst_g and snd_g remove duplications from a pair and reintroduce them in the other direction. Note that the function g itself is not necessarily in RINV; it just needs to be definable in the host language (in our case, Haskell).

4.2. The Constructors

The semantics of the constructor functions are simple: they follow directly from the corresponding constructors introduced by datatype declarations, except that they are uncurried. For example,

$$[\![nil]\!] = \lambda() \to [\!]$$
$$[\![cons]\!] = \lambda(x, xs) \to x : xs$$

Inverses of the primitive constructor functions are obtained simply by swapping the right- and left-hand sides of the definitions.

They are effectively partial 'guard' functions, succeeding when the input value matches the pattern.

Although *snoc* and *wrap* are not the primitive constructors for left-biased lists, they can be encoded:

$$\begin{bmatrix} snoc \end{bmatrix} = \lambda(xs, x) \to xs + [x]
\llbracket wrap \rrbracket = \lambda x \to [x]$$

The right inverses of *snoc* and *wrap* are

$$\begin{bmatrix} snoc \end{bmatrix}^{\circ} [x] = ([], x) \\
\llbracket snoc \rrbracket^{\circ} (x : xs) = \mathbf{let} (ys, y) = \llbracket snoc \rrbracket^{\circ} xs \mathbf{in} (x : ys, y) \\
\llbracket wrap \rrbracket^{\circ} [x] = x$$

The inverses of constructor functions are generally not case-exhaustive. For example, $[\![cons]\!]^{\circ}$ only accepts non-empty lists, while $[\![nil]\!]^{\circ}$ only accepts the empty list. As a result, in contrast to primitive functions, constructor functions cannot be composed arbitrarily.

4.3. The Combinators

The combinators in RINV are mostly standard.

4.3.1. Composition, Sum and Product

Combinator \circ sequentially composes two functions:

Its inverse is the reverse composition of the inverses of the two arguments.

Combinators \times and ∇ compose functions in parallel. The former applies a pair of functions component-wise to its input:

$$\begin{array}{l} (\times) :: (a \leftrightarrows b) \to (c \leftrightarrows d) \to ((a,c) \leftrightarrows (b,d)) \\ \llbracket f \times g \rrbracket &= \lambda(w,x) \to (\llbracket f \rrbracket \ w, \llbracket g \rrbracket \ x) \\ \llbracket f \times g \rrbracket^{\circ} &= \lambda(y,z) \to (\llbracket f \rrbracket^{\circ} \ y, \llbracket g \rrbracket^{\circ} \ z) \end{array}$$

Note that we have chosen not to define \times in terms of the more primitive combinator \triangle that executes both of its input functions on a single datum:

$$(\triangle) :: (a \leftrightarrows b) \to (a \leftrightarrows c) \to (a \leftrightarrows (b, c))$$
$$[\![f \triangle g]\!] = \lambda x \to ([\![f]\!] x, [\![g]\!] x)$$

In the inverse direction, $[\![f]\!]^\circ x$ and $[\![g]\!]^\circ y$ would have to agree, which is difficult to enforce statically. Indeed, functions constructed with Δ are generally not surjective, and so do not have total right inverses; for this reason, we exclude Δ from RINV.

The combinator ∇ constructs a function that consumes an element of a sum type ('*Either*' in Haskell).

$$(\triangledown) :: (a \leftrightarrows c) \to (b \leftrightarrows c) \to (Either\ a\ b \leftrightarrows c)$$

$$\llbracket f \triangledown g \rrbracket = \lambda x \to \mathbf{case}\ x\ \mathbf{of}\ \{Left\ a \to \llbracket f \rrbracket\ a\ ; Right\ b \to \llbracket g \rrbracket\ b\ \}$$

In the inverse direction, if both f and g are surjective, it doesn't matter which branch is chosen. However, the use of constructor functions deserves some attention, since they are not surjective in isolation. In contrast to the case of Δ , here the totality in the inverse direction can be recovered by choosing a non-failing branch. As a result, in the event that $[\![f]\!]^{\circ}$ fails on certain inputs, $[\![g]\!]^{\circ}$ should be applied. To model this failure handling, we lift functions in RINV into the Maybe monad (allowing an extra possibility for the return value), and handle a failure in the first function by invoking the second.

```
 \llbracket f \triangledown g \rrbracket^{\circ} = \lambda x \to (\llbracket f \rrbracket^{\circ} x) \text{ `choice' } (\llbracket g \rrbracket^{\circ} x)   \text{choice :: } Maybe \ a \to Maybe \ a \to Maybe \ a   \text{choice } (Just \ x) \ \_ = x   \text{choice } \ \_ (Just \ y) = y
```

This shallow backtracking is sufficient because the guards of conditionals are only pattern matching outcomes, which are completely decided at each level. With the introduction of the Maybe monad, the types of the functions in RINV have to be lifted. As a result, the invertible function type $s \leftrightarrows t$ is actually the pairing of $s \to Maybe$ t and $t \to Maybe$ s. For brevity, in the main text of this paper, we still use the non-monadic types for RINV functions to avoid the explicit handling of Maybe, with the understanding that all functions are lifted to the Maybe monad in the implementation. The details are spelt out in Appendix B.

4.3.2. Recursion

With the ground prepared, we are now ready to discuss a recursive combinator. We define

$$\llbracket fold_X f \rrbracket^{\circ} = unfold_X \llbracket f \rrbracket^{\circ}$$

The forward semantics of $fold_X f$ is the standard fold for a datatype X; its inverse semantics is defined by a corresponding $unfold_X$. Intuitively, fold disassembles a structure and replaces the constructors with applications of the body of the fold. Function unfold, on the other hand, takes a seed, splitting it with the body of the unfold into building blocks of a structure and new seeds, which are themselves recursively unfolded. In short, fold collapses a structure, whereas unfold grows one.

When an algebraic data type X is given, Haskell definitions of $fold_X$ and $unfold_X$ can be constructed mechanically. Our prototype automatically generates fold and unfold for each non-mutually recursive regular data type. For example, consider the data type of lists:

Another example is leaf-labelled binary trees. Note that the constructor *Fork* is uncurried, to fit better into the RINV framework.

```
 \begin{aligned} \mathbf{data} \ \mathit{LTree} \ a &= \mathit{Leaf} \ a \mid \mathit{Fork} \ (\mathit{LTree} \ a, \mathit{LTree} \ a) \\ \mathit{fold}_{\mathit{LTree}} &:: (\mathit{Either} \ a \ (b, b) \to b) \to \mathit{LTree} \ a \to b \\ \mathit{fold}_{\mathit{LTree}} \ f &= \lambda t \to \mathbf{case} \ t \ \mathbf{of} \\ \mathit{Leaf} \ a &\to f \ (\mathit{Left} \ a) \\ \mathit{Fork} \ (t_1, t_2) \to f \ (\mathit{Right} \ (\mathit{fold}_{\mathit{LTree}} \ f \ t_1, \mathit{fold}_{\mathit{LTree}} \ f \ t_2)) \\ \mathit{unfold}_{\mathit{LTree}} &:: (b \to \mathit{Either} \ a \ (b, b)) \to b \to \mathit{LTree} \ a \\ \mathit{unfold}_{\mathit{LTree}} \ f &= \lambda b \to \mathbf{case} \ f \ b \ \mathbf{of} \\ \mathit{Left} \ a &\to \mathit{Leaf} \ a \\ \mathit{Right} \ (b_1, b_2) \to \mathit{Fork} \ (\mathit{unfold}_{\mathit{LTree}} \ f \ b_1, \mathit{unfold}_{\mathit{LTree}} \ f \ b_2) \end{aligned}
```

We use unfold to construct fold's right inverse. The following lemma is well known from the literature [26].

Lemma 2.
$$fold \ \llbracket f \rrbracket \circ unfold \ \llbracket f \rrbracket^{\circ} \sqsubseteq id.$$

The use of \sqsubseteq in the above lemma states that the left-hand side might be less defined than the right-hand side. Since both fold and unfold are case-exhaustive when their bodies are case-exhaustive, the only issue is the termination of unfold: when a body does not split a seed into 'smaller' seeds, unfolding that seed creates an infinite structure. It is well known that a function constructed by unfold terminates if the seed transformation is well-founded (that is, the 'leads to' ordering on seeds induced by the body of the unfold should not admit an infinite descending chain), at least on polynomial datatypes. Such well-founded seed transformations are known as recursive coalgebras [27], and allow the above composition of fold and unfold to be well-defined in a setting of total functions. Static termination checkers exist in the literature [28, 29] and are orthogonal to the discussion here.

4.4. Programming in RINV

We are now ready to look into the kinds of function we can define with RINV.

To start with, let's look first at a derived combinator map that can be defined in terms of fold. For example, map on lists, map_{List} , is defined as follows.

$$\begin{array}{l} \operatorname{map}_{\operatorname{List}} :: (a \leftrightarrows b) \to (\operatorname{List}\ a \leftrightarrows \operatorname{List}\ b) \\ \operatorname{map}_{\operatorname{List}} f = \operatorname{fold}_{\operatorname{List}} \left(\operatorname{nil} \triangledown \left(\operatorname{cons} \circ (f \times \operatorname{id})\right)\right) \end{array}$$

Function $map_{List} f$ applies argument f uniformly to all the elements of a list, without modifying the list structure. Since nil and cons form a complete set of constructors for lists, we know they are jointly surjective.

Similarly, map on leaf-labeled trees, map_{LTree} , is defined as follows.

$$\begin{aligned} & \operatorname{map}_{\operatorname{LTree}} :: (a \leftrightarrows b) \to (\operatorname{Tree}\ a \leftrightarrows \operatorname{Tree}\ b) \\ & \operatorname{map}_{\operatorname{LTree}}\ f = \operatorname{fold}_{\operatorname{LTree}}\ ((\operatorname{leaf} \circ f) \triangledown \operatorname{fork}) \end{aligned}$$

The function reverse on lists can be defined as a fold:

$$\begin{array}{ll} reverse &= fold_{\mathit{List}} \; (nil \; \triangledown \; snoc) \\ \llbracket reverse \rrbracket^{\circ} &= unfold_{\mathit{List}} \; \llbracket nil \; \triangledown \; snoc \rrbracket^{\circ} \\ \end{array}$$

Function reverse is one of many functions for which the use of derived constructors (such as snoc) is essential. If we restrict ourselves to the primitive constructors, we can only reverse a list into a different type. In the forward direction, a list is taken apart and the first element is appended to the rear of the output list by snoc. This process terminates on reaching an empty list, when an empty list is returned as the result. Function $[snoc]^{\circ}$ extracts the last element in a list and adds it to the front of the result list by unfold, which terminates when $[nil]^{\circ}$ can be successfully applied (i.e when the input is the empty list). Since nil and snoc form a complete set of constructors for lists, they are jointly surjective.

Function reverse is also used to construct the apprev function that reverses a list and appends it.

$$apprev :: ([a], [a]) \leftrightarrows [a]$$

 $apprev = app \circ (id \times reverse)$

For example, we have:

$$[apprev] ([1,2],[3,4,5,6,7]) \equiv [1,2,7,6,5,4,3]$$

The companion $[apprev]^{\circ}$ function is

$$\llbracket apprev \rrbracket^{\circ} :: [a] \to ([a], [a])$$

$$\llbracket apprev \rrbracket^{\circ} = \llbracket app \circ (id \times reverse) \rrbracket^{\circ}$$

which can be reduced to

$$(\llbracket id \rrbracket^{\circ} \times \llbracket reverse \rrbracket^{\circ}) \circ \llbracket app \rrbracket^{\circ}$$

In the inverse direction, a list is split into two, and functions $[id]^{\circ}$ and $[reverse]^{\circ}$ are applied to the two parts. For example, we have

On the other hand,

$$[apprev]^{\circ} ([apprev]] ([1,2],[3,4,5,6,7])) \equiv [apprev]^{\circ} ([1,2,7,6,5,4,3])$$

$$\equiv ([1,2,7,6],[3,4,5])$$

It is clear from above that $[apprev]^{\circ}$ is not a left inverse of [apprev], and it is not intended to be a term in the language RINV.

The other function we looked at in Section 3.4 yields the list view of a size-annotated binary tree. This abstraction function is defined as

$$\alpha = fold_{STree} \ (nil \ \forall \ wrap \ \forall \ (app \circ snd_{\lambda(x,y) \to size \ x+size \ y}))$$

Function α drops the size annotations from a tree and flattens it into a list. In the α° direction, the function $\lambda(x, y) \to size \ x + size \ y$ is used to recompute the lost size information.

As a remark, the primitive function app can be defined in Haskell with $foldr :: (a \to b \to b) \to b \to [a] \to b$ as:

$$app = uncurry (flip (foldr (:)))$$

which effectively partially applies foldr and awaits an input as the base case. This idiom of taking an extra argument to form the base case constructs the fold body during execution, whereas in RINV, fold bodies are constructed syntactically so that they can be checked and inverted separately from the recurisve combinator. As a result, we include app as a primitive function in RINV.

4.5. Extending RINV

We use RINV for the purpose of constructing conversion functions between different representations of data. As a result, RINV has to be extended from time to time, when new abstract representations are introduced. (New concrete representations only require corresponding versions of fold and unfold, already captured by the parameterized $fold_X$ and $unfold_X$.) At the very least, a new abstract representation brings in a new set of primitive constructors. For example, a binary tree abstract representation might extend RINV with the following:

$$\begin{bmatrix} tip \end{bmatrix} = \lambda x & \to Tip \ x \\
 \llbracket fork \end{bmatrix} = \lambda(lt, rt) \to Fork \ lt \ rt$$

The inverses of primitive constructors are always straightforward; and the annotation that records the constructor sets will be extended with the new set $\{tip, fork\}$.

$$\begin{aligned}
&[tip]^{\circ} = \lambda(Tip \ x) \to x \\
&[fork]^{\circ} = \lambda(Fork \ lt \ rt) \to (lt, rt)
\end{aligned}$$

We can now start using binary trees, for example, as the abstract representation of sized trees.

$$\alpha = fold_{STree} \left(tip \ \nabla \left(fork \circ snd_{\lambda(x,y) \to size \ x+size \ y} \right) \right)$$

Only primitive constructors (together with the existing data projection function snd_g) are used in this case. In general, we may need to extend the set of primitive functions as well. For example, consider using natural numbers as an abstract representation:

$$\mathbf{data} \ Nat = Zero \\ | \ Succ \ Nat$$

which are then implemented as binary numbers with the least significant bit first:

type
$$Binary = [Nat]$$

Since we are representing binary numbers, only *Zero* and *Succ Zero* can appear in the list. We do not use a specialized datatype for binary bits, as found in the example in Section 1.2, because of the need to add binary bits to natural numbers in the process of converting between the two representations.

Again, primitive constructors arise from the datatype definition:

$$[zero] = \lambda() \rightarrow Zero$$

 $[succ] = \lambda n \rightarrow Succ n$

and these can easily be inverted:

$$[\![zero]\!]^\circ = \lambda Zero \longrightarrow ()$$

$$[\![succ]\!]^\circ = \lambda Succ \ n \longrightarrow n$$

To construct the conversion function from Binary to Nat, we need an additional primitive function

$$\begin{array}{l} plusTwice :: (Nat, Nat) \leftrightarrows Nat \\ \llbracket plusTwice \rrbracket = \lambda(c, x) \to c + 2 * x \\ \llbracket plusTwice \rrbracket^{\circ} = \lambda y \to (mod \ y \ 2, \ div \ y \ 2) \end{array}$$

Function plusTwice is a linear function with slope two, whereas its right inverse calculates the x-value and the y-intercept from a y-value. Similar to app, function plusTwice is surjective but not injective. Among the multiple choices, the user will have to choose one that suits the intended use. Our experience is that the appropriate choice is usually fairly obvious. For example, the above pattern of combining mod and div is quite common; it even appears as the function divMod in the Haskell standard prelude.

With the newly introduced primitive function, we can define the conversion from Binary to Nat as

$$\alpha = fold_{List} \ (zero \ \nabla \ plusTwice)$$

In a sense, trying to refactor a datatype that is not already supported by RINV requires an extension to the language. As we have seen in this section, such extensions of primitive functions and constructors typically only generate small proof obligations. On the other hand, we also know that the nature of abstract representations means that they are relatively simple and few in number. It is a known fact that functional programmers "too often reach for lists when an ADT would be more appropriate" [30]. As a result, we might expect a library that is sufficient for most common cases can be built up reasonably quickly.

5. Discussion

5.1. RINV Expressiveness

The most general constraint on α functions is surjectivity, in order to ensure the existence of right inverses: valid abstract values are bounded by the actual range of the user-defined α function; invertibility is not guaranteed for abstract values outside this range. In the current proposal, RINV faithfully enforces surjectivity, which explains its restricted expressiveness compared to the standard point-free programming framework. An already-mentioned example that shows this difference is the combinator Δ , which executes both of its input functions, and is defined as

$$(\Delta) :: (a \to b) \to (a \to c) \to a \to (b, c)$$
$$(f \Delta g) = \lambda x \to (f x, g x)$$

Since $f \triangle g$ is generally not surjective, it has no right inverse, despite the fact that we can easily guard against inconsistent input in the reverse direction:

```
\llbracket f \triangle g \rrbracket^{\circ} = \lambda(a, b) \rightarrow \mathbf{if} \ x == y \ \mathbf{then} \ x \ \mathbf{else} \ error \ "violation"

\mathbf{where} \ x = \llbracket f \rrbracket^{\circ} \ a \ ; \ y = \llbracket g \rrbracket^{\circ} \ b
```

Definitions like the one above are known as weak right inverses [31]. Another useful function is unzip, which can be defined as a fold.

```
unzip :: List (a, b) \leftrightarrows (List a, List b)

unzip = fold_{List} ((nil \triangle nil) \triangledown ((cons \times cons) \circ trans))

trans :: ((a, b1), (b2, c)) \leftrightarrows ((a, b2), (b1, c))

trans = assocl \circ (id \times subr) \circ assocr
```

This definition will be rejected in RINV, since $cons \times cons$ and $nil \triangle nil$ are not jointly surjective—indeed, unzip only produces pairs of lists of equal length.

If a value outside the range is constructed, the integrity of specification-level equational reasoning may be corrupted. On the other hand, it is valid to argue that the same invariant assumed for the original datatype prior to the refactoring applies to the specification too. For example, consider a program that requires balanced binary trees. An abstraction function that only produces balanced binary trees is safe if the invariant is correctly preserved in the original program. It remains an open question whether we should allow programmers to take some reasonable responsibilities, or should insist on enforcing control through the language.

We do not include unfold as a combinator in RINV as the dual of fold because there is no combinator in RINV that creates values of a sum type: the choice operator introduced for the right-inverse of ∇ does not preserve surjectivity. An option is to define the bodies of unfolds as primitive functions; but this does require the programmers to deal with Either types explicitly, something we have tried to avoid.

5.2. The Dual Story

In this paper, we have picked the α function to be user-provided; the design of RINV and the subsequent discussion of refactoring is based on this decision. However, this choice is not absolute. One can well imagine a programmer coming up with α° functions first, and a left-invertible language generating the corresponding α functions; this would give the same invertibility property $\alpha \circ \alpha^{\circ} \equiv id$. The promotion condition can be adapted to involve only α° , as in $\underline{f} \circ \mathcal{F} \alpha^{\circ} \equiv \mathcal{G} \alpha^{\circ} \circ f$. Nevertheless, the crucial computation law and fusion law that form the foundation of the translation and optimization are still derivable; for computation, we have

```
\begin{array}{l} \mathcal{G} \ \alpha \circ \underline{f} \circ \mathcal{F} \ \alpha^{\circ} \\ \equiv \quad \{ \underline{f} \circ \mathcal{F} \ \alpha^{\circ} \equiv \mathcal{G} \ \alpha^{\circ} \circ f \ \} \\ \mathcal{G} \ \alpha \circ \mathcal{G} \ \alpha^{\circ} \circ f \\ \equiv \quad \{ \mathcal{G} \ \text{respects composition and identity; } \alpha \circ \alpha^{\circ} \equiv id \ \} \\ f \end{array}
```

and for fusion:

$$\begin{array}{l} \mathcal{H} \ \alpha \circ \underline{g} \circ \mathcal{G} \ \alpha^{\circ} \circ \mathcal{G} \ \alpha \circ \underline{f} \circ \mathcal{F} \ \alpha^{\circ} \\ \equiv & \{\underline{f} \circ \mathcal{F} \ \alpha^{\circ} \equiv \mathcal{G} \ \alpha^{\circ} \circ f \} \\ \mathcal{H} \ \alpha \circ \underline{g} \circ \mathcal{G} \ \alpha^{\circ} \circ \mathcal{G} \ \alpha \circ \mathcal{G} \ \alpha^{\circ} \circ f \\ \equiv & \{\mathcal{G} \ \text{respects composition and identity;} \ \alpha \circ \alpha^{\circ} \equiv id \} \\ \mathcal{H} \ \alpha \circ \underline{g} \circ \mathcal{G} \ \alpha^{\circ} \circ f \\ \equiv & \{\underline{f} \circ \mathcal{F} \ \alpha^{\circ} \equiv \mathcal{G} \ \alpha^{\circ} \circ f \} \\ \mathcal{H} \ \alpha \circ \underline{g} \circ \underline{f} \circ \mathcal{F} \ \alpha^{\circ} \end{array}$$

If we were to develop a left-invertible language in a similar style to RINV's, we expect that many problems will dualize. A notable difference is that users of the language will have to use unfold (as ∇ is not injective), with destruction functions combined together by the backtracking operator choice, resulting in a rather unusual programming style.

6. Further Applications

6.1. Pattern Matching for ADTs

The refactoring technique proposed in this paper rewrites selected definitions using pattern matching into ones using library operations; the concrete implementations of the library operations are not used in reasoning, and thus are not exposed. In this case, a more structured way of organizing the refactored program is with abstract datatypes (ADTs) [32].

As a matter of fact, the data abstraction framework developed in this paper can be used as a way of introducing pattern matching to certain ADTs. For example, we can construct an ADT for queues:

```
adt Queue a = [a] where

emptyQ :: Queue \ a

enQ :: a \rightarrow Queue \ a \rightarrow Queue \ a

deQ :: Queue \ a \rightarrow Queue \ a
```

```
first :: Queue a \rightarrow a is Empty :: Queue a \rightarrow Bool
```

with the following abstract programs as specifications:

```
emptyQ = []

first = head

isEmpty = null

enQ \ a \ q = q + [a]

deQ = tail
```

This approach is known as *constructive specification* [10]: the semantics of operations are explicitly defined by expressing them in terms of a model. For example, the queue ADT is related to the list model. It is worth emphasising that the list datatype acts only as a model of the ADT: it may suggest but it does not imply a particular implementation. We also note that this constructive approach does not cover all ADTs: for example, unordered sets cannot be fully modelled by an algebraic datatype.

Once an implementation of the ADT is shown correct (by proving the promotion condition), users of the ADT can pattern match on the model when defining non-library functions, and reason about these functions at the level of models. A translation based on the computation law, similar to the one in Section 3.2, elaborates the semantics of programs using the ADT.

6.2. Stream Fusion

Streams are the dual of lists, and are a kind of *codata*. Instead of being built as data structures, streams encapsulate operations that can be unfolded to produce stream elements. Coutts et al. [33] introduce an ADT of streams in their work on fusion optimizations:

```
data Stream\ a = \exists s. Strm\ (s \rightarrow Step\ a\ s)\ s
data Step\ a\ s = Done
|\ Yield\ a\ s
|\ Skip\ s
```

To use the stream ADT as an implementation of lists, they have functions *stream* and *unstream* for converting lists to and from streams.

```
unstream :: Stream \ a \rightarrow [a]
unstream \ (Strm \ next \ s) = unfold \ s
```

```
where unfold\ s = \mathbf{case}\ next\ s of
Done \rightarrow []
Skip\ s' \rightarrow unfold\ s'
Yield\ x\ s' \rightarrow x : unfold\ s'
stream\ :: [a] \rightarrow Stream\ a
stream\ xs = Strm\ next\ xs
\mathbf{where}\ next\ [] = Done
next\ (x:xs) = Yield\ x\ xs
```

The *stream* function is non-recursive, corresponding to the non-recursive definition of *Stream*. The *unstream* function repeatedly calls the access operation of the stream, and produces a list by accumulating the elements this yields. The step *Skip* is unproductive, and therefore does not contribute to expressive power. However, it is crucial for the efficiency of their implementation, allowing intermediate stream/unstream conversions to be fused. For example, they want to reduce the expression

 $unstream \circ maps \ f \circ stream \circ unstream \circ maps \ g \circ stream$

to

 $unstream \circ maps \ f \circ maps \ g \circ stream$

This would follow from $stream \circ unstream \equiv id$; but this property is not satisfied by their conversion functions, and so their paper leaves open the question of soundness. Nevertheless, the result they want still holds; in fact, it is a consequence of our approach (Theorem 1), with unstream the abstraction function and stream its right inverse.

It is worth mentioning that despite our fusion theorem developed in Section 3.3 can be used to prove the correctness of Coutts et al.'s optimization, the implementation of the *unstream* and *stream* function pair is not able to benefit from RINV. The main problem is that *stream* is non-recursive; and the information needed for constructing *Stream* data is stored in its function component (*next*). At the moment, RINV is not able to handle such higher-order cases.

7. Related Work

The theoretical foundation of this work goes all the way back to the dawn of the program correctness era [12, 34, 35, 10]. The promotion condition and

computation law are standard tools for establishing the correctness of program development. In particular, full invertibility of the abstraction function has been shown [12] to be a sufficient (although not necessary) precondition for the computation law. As expected, none of the early works deal specifically with pattern matching and its implications in program development, since this language feature first appeared only a decade later.

7.1. Pattern Matching with Data Abstraction

Efforts to combine data abstraction and pattern matching started over twenty years ago with Wadler's *views* proposal [8]; and it is still a hot research topic [36, 37, 38, 39, 40, 41, 42, 43, 7, 44]. To avoid possible confusion over terminology, we always refer to Wadler's proposal [8] as "Wadler's views".

Wadler's views provide different ways of viewing data than their actual implementations. With a pair of conversion functions, data can be converted to and from a view. Consider the left- and right-biased representations of lists:

```
\begin{array}{lll} \mathbf{data} \ List \ a = Nil \mid Cons \ a \ (List \ a) \\ \mathbf{view} \ List \ a = Lin \mid Snoc \ (List \ a) \ a \\ & \text{to} \ Nil & = Lin \\ & \text{to} \ (Cons \ x \ Nil) & = Snoc \ Nil \ x \\ & \text{to} \ (Cons \ x \ (Snoc \ xs \ y)) & = Snoc \ (Cons \ x \ xs) \ y \\ & \text{from} \ Lin & = Nil \\ & \text{from} \ (Snoc \ Nil \ x) & = Cons \ x \ Nil \\ & \text{from} \ (Snoc \ (Cons \ x \ xs) \ y) & = Cons \ x \ (Snoc \ xs \ y) \end{array}
```

The **view** clause introduces two new constructors, Lin and Snoc, which may appear in both terms and patterns. The first argument to the view construction Snoc refers to the datatype List a, so a snoclist actually has a conslist as its child. The **to** and **from** clauses are similar to function definitions. The 'to' function **to** converts a conslist value to a snoclist value, and is used when Lin or Snoc appear as the outermost constructor in a pattern on the left-hand side of an equation. Conversely, the 'from' function **from** converts a snoclist into a conslist, when Lin or Snoc appear in an expression. Note that we are already making use of views in the definition above; for example, Snoc appears on the left-hand side of the third **to** clause; matching against this will trigger a recursive invocation of **to**.

Functions can now pattern match on and construct values in either the datatype or one of its views.

```
last (Snoc xs x) = x
rotLeft (Cons x xs) = Snoc xs x
rotRight (Snoc xs x) = Cons x xs
rev Nil = Lin
rev (Cons x xs) = Snoc (rev xs) x
```

Upon invocation, an argument is converted into the view by the to function; after completion of the computation, the result is converted back to the underlying datatype representation by the from function.

Just as with our proposal, this semantics can be elaborated by a straightforward translation into ordinary Haskell. First of all, view declarations are translated into data declarations:

```
data Snoc \ a = Lin \mid Snoc \ (List \ a) \ a
```

Note that the child of *Snoc* refers to the underlying datatype: view data is typically hybrid (whereas it is homogeneous with our approach). Now the only task is to insert the conversion functions at appropriate places in the program.

```
last xs = \mathbf{case} to xs of Snoc\ xs\ x \to x

rotLeft\ xs = \mathbf{case}\ xs of Cons\ x\ xs \to from\ (Snoc\ xs\ x)

rotRight\ xs = \mathbf{case} to xs of Snoc\ xs\ x \to Cons\ x\ xs

rev\ xs = \mathbf{case}\ xs of Nil\ \to from\ Lin\ (Cons\ x\ xs) \to from\ (Snoc\ (rev\ xs)\ x)
```

To take advantage of the hybrid representation of views, Wadler's applications mostly involves cute examples with frequent changes of views in one function definition. When the hybrid structures fit with the recursive pattern of a function, the shallow conversion pays off by improving efficiency. On the other hand, in the applications we are looking at, it is more common to have only one view involved in a recursion, which certainly benefits from having homogeneous data.

In contrast to our approach, Wadler exposes both a datatype and its views to programmers. To support reasoning across the different representations, the conversion clauses are used as axioms. It is expected for a view type to be isomorphic to a subset of its underlying datatype, and for the pair of conversions between the values of the two subsets to be each other's full inverses. This is certainly restrictive; and Wadler didn't suggest any way to enforce such an invertibility condition. As pointed out by Wadler himself [8], and followed up by several others [36, 37], this assumption is risky, and may lead to nasty surprises that threaten soundness of reasoning.

Inspired by Wadler's proposal, our work ties up the loose ends of views by hiding the implementation of selected primitives that are proven correct, and using only the view (our abstract representation) for pattern matching in user-defined functions. The language RINV for defining conversions guarantees right invertibility, a weaker condition that lifts the isomorphism restriction on abstract and concrete representations. However, in contrast to Wadler's views, our system only caters for linear refactoring—we cannot provide multiple views for the same implementation.

'Safe' variants of views [36, 37] have been proposed before. To circumvent the problem of equational reasoning, one typically restricts the use of view constructors to patterns, and does not allow them to appear on the right-hand side of a definition. As a result, expressions like $Snoc\ Lin\ 1$ become syntactically invalid. Instead, values are only constructed by 'smart constructors', as in $snoc\ lin\ 1$. In this setting, equational reasoning has to be conducted on the source level with explicit applications of to. A major motivation for such a design is to admit views and sources with conversion functions that do not satisfy the invertibility property. In another words, let $Constr\$ and $constr\$ be a constructor and its corresponding smart constructor; in general, we have $Constr\ x \not\equiv constr\ x$. This appears to hinder program comprehension, since the very purpose of the convention that the name of a smart constructor differs only by case from its 'dumb' analogue is to suggest the equivalence of the two.

More recently, language designers have started looking into more expressive pattern mechanisms. Active patterns [38, 39] and many of their variants [40, 41, 42, 43, 7] go a step further, by embedding computational content into pattern constructions. All the above proposals either explicitly recognise the benefit of using constructors in expressions, or use examples that involve construction of view values on the right-hand sides of function definitions. Nevertheless, none of them are able to support pattern constructors in expressions, due to the inability to reason safely. Knowing that there is an absence of good solutions for supporting constructors in expressions, some work focuses mainly on examples that are primarily data consumers, an escape that is expected to be limited and short-lived. Another common pitfall

of active patterns is the difficulty in supporting nested and overlapping patterns, because each active pattern is computed and matched independently.

Our proposal supports the full power of pattern matching while preserving safety. Yet, the additional expressiveness and guarantees come with a higher price tag: RINV is still rather restrictive, and can be somewhat awkward at times. Among alternative approaches, one can see higher-order recursion patterns such as hylomorphisms [45] as ways of imposing recursive views on possibly non-recursive abstract data. In the special case that the input is itself an inductively defined datatype, such recursive views can provide alternative decompositions to plain iteration—for example, paramorphisms or primitive recursion [46] and histomorphisms or course-of-value iteration [47]. Because RINV as presented here is first-order, it cannot capture these higherorder eliminators; we leave as future work the exploration of a higher-order version of RINV. McBride and McKinna [48] point out that view mechanisms providing case analysis are more informative in a dependently typed setting than a simply typed one, because the individual cases may have more refined types than the surrounding context; perhaps switching from polymorphic functional programming to dependent types will eliminate much of the need for the kind of language extension presented here.

7.2. Invertible Programming

The language RINV owes its origins to the rich literature on *invertible* programming [19, 22], a programming paradigm where programs can be executed both forwards and backwards. Mu et al. [19] concentrate their effort on designing a language that provides only injective functions. The resulting language Inv is a combinator library that syntactically rules out any non-injective functions. The most novel operator of Inv is dup f, which duplicates the input and applies f to one copy. In the backward direction, the two copies of the duplicated input are checked for consistency before being restored. It is shown that Inv is practically useful for maintaining consistency of structured data related by some transformations [20, 49]. The 2LT (Two-Level-Transformation) system [50] has a left-invertible combinator library at its core. Instead of addressing changes to individual values, the update is expressed as a format (type) evolution, which migrates a whole database to one in a new format. A typical usage of the system involves retrieving relational data into a hierarchical format, updating by enriching the hierarchical format, and putting back the updated data into a new relational format reflecting the enrichment. The whole process is specified in terms of formats; the corresponding value-level transformations are induced. Invertible arrows [22] extend the arrow framework [51] (a generalization of monads) with a combinator that encodes pairs of functions being each other's inverses; the paper introducing invertible arrows [22] recognizes that, when full invertibility is not achievable (due to the non-isomorphic nature of the two sides), left- or right-biased semi-invertibility is nevertheless a useful approximation.

Right inverses have been studied as a component of the much more elaborate bidirectional programming framework of lenses [21, 52, 53, 54]; in this context, right inverses are known as 'create' functions. Based on record types, the combinators of lenses have little similarity to those of RINV. A distinctive feature of the lenses framework is the use of semantic types [55] to give precise bounds to the ranges of forward functions (thus the domains of backwards functions). As a result, surjectivity now concerns the relationships between the domains and ranges of lenses connected by a combinator, instead of being a property between a function and its target datatype.

We have carefully chosen right invertibility as the central property our system is based on, to balance expressiveness and robustness. As we have seen, it has served our purpose well by being sufficient to establish the crucial computation and fusion laws.

On the other hand, RINV does not keep values on the concrete level stable: bringing a value on the concrete level up to the abstract level and putting it back without modification may induce changes. This has not been a problem with our design, but rules out the possibility of having multiple abstract representations for a single implementation. For example, as shown in Wadler's view proposal [8], it can be useful to have both *cons*-lists and *snoc*-lists as views at the same time.

Another omission of RINV is general data-projection functions. We can only discard information that is re-constructible, in other words redundant. At a glance, it may appear that both the above "flaws" would be fixed if we were to use a framework similar to that of bidirectional programming [49, 56, 21, 57, 23, 58], where a forward function $fwd :: A \to B$ is coupled with a "backwards" counterpart $bwd :: (B, A) \to A$. The idea is that fwd transforms a source into a view, which is then modified; the modification is to be saved in the form of an updated source, produced by bwd. Instead of being an inverse, bwd takes both the original source and a modified view as inputs, and tries to reinstate the source/view consistency by producing a new source. In this setting, stability of sources is achievable; and discarded values can be recovered by extracting them from the copied source. However, in

our case the modification operations, which are basically arbitrary abstract operations, are far more complicated than those in bidirectional programming, which assumes a single type for the view. As a result, an update from type $(List\ a, List\ a)$ to type $List\ a$ would not be allowed in bidirectional programming, but is perfectly sensible in programming with abstractions. If we use bwd on the resulting list, it is never clear what the source input is.

As we can see from the literature summarized above, building invertible languages using a combinatorial approach is by no means a new idea. Through coming from a common foundation, different proposals customize their designs to cater for specific applications. For example, lenses [21, 52, 53, 54] deal with data as records, whereas Inv [19], pointfree lenses [23], and RINV deal with data as algebraic datatypes. In a sense, any invertible language admitting algebraic datatypes can be used in place of RINV; our refactoring framework will still work, but not as well. For example, executing inverses in Inv may fail at run-time, which will cause failures in pattern matching, while point-free lenses only provide primitive constructors, which is limited in expressiveness. We have carefully designed RINV to balance safety and expressiveness for our application. The non-primitive constructor functions and annotated data projection functions are artifacts of such consideration; and both the new features are nicely integrated into the more standard combinator framework.

Program inversion, in the sense of generating an inverse of a program, is not the only way of obtaining an input from the output that it yields. Glück and Abramov [59] proposed a universal resolving algorithm, which, give an output and a program, enumerates the original inputs. The core technique they use, known as positive driving [60], is similar to the notion of needed narrowing [61] studied in the context of functional logic programming. Thus, their work can be seen as interpretation of a functional program by using the semantics of functional logic programs, a connection that is made more explicit in subsequent work [62]. Note that, if a function is surjective, the inverse computation by the universal resolving algorithm is always guaranteed to terminate.

Matsuda et al. [63] discussed right-inverse computation using grammar parsing. They construct a grammar from a program, so that the productions of the grammar abstract the evaluation of the program. Then, right-inverse computation is given by parsing with respect to the grammar. In their approach, an affine and treeless first-order functional program [64] is always surjective onto the language of the derived grammar; and its inverse compu-

tation has linear complexity with respect to the original program's output.

8. Conclusion

Algebraic datatypes and pattern matching offer great promise to programmers seeking simple and elegant programming, but the promise turns sour when modular changes are demanded. Our work tackles this long-standing problem by proposing a framework for refactoring programs written with pattern matching into ones with proper encapsulation: programmers are able to selectively reimplement original function definitions into primitive operations, and either rewrite the rest in terms of the primitive ones, or simply leave them unchanged. This migration is completely incremental: executability and proofs through equational reasoning are preserved at all times during the process.

At the heart of our proposal is the framework of data abstraction. When an abstraction function is verified by the promotion condition, the computation law is able to replace abstract function calls with concrete ones. The soundness of such refactoring is based on the right-inverse property of the conversion pairs that bridge the abstract and concrete representations, for which we have designed RINV to guarantee right-invertibility by construction.

At this stage, our focus is on supporting refactoring of programs written with datatypes and pattern matching, which automatically excludes some ADTs, such as unordered sets, that cannot be fully modelled by algebraic datatypes. We leave it as future work to investigate the applicability of our proposal in a more general setting.

Acknowledgements

We are grateful to the anonymous reviewers of the earlier conference version of this paper [65] and of this revised and expanded version; we are especially grateful to Ralf Hinze for his valuable comments on an early draft of the paper—the binary number example is due to him. This work was supported by the UK Engineering and Physical Sciences Research Council through the *Generic and Indexed Programming* project (EP/E02128X) and a PhD Plus award on *Bidirectional Programming* via the Doctoral Training Grant (EP/P503876/1), and was partly conducted during Wang's internship at National Institute of Informatics, Japan. Matsuda is supported by Grantin-Aid for Start-up 22800003, and part of the work was done while he was at

the University of Tokyo as JSPS Research Fellow supported by Grant-in-Aid for JSPS Fellows $20 \cdot 9584$.

References

- [1] P. Wadler, A critique of Abelson and Sussman: Why calculating is better than scheming, ACM SIGPLAN Notices 22 (3) (1987) 83–94.
- [2] S. Peyton Jones (Ed.), Haskell 98 Language and Libraries: The Revised Report, Cambridge University Press, 2003.
- [3] R. Burstall, D. MacQueen, D. Sannella, Hope: An experimental applicative language, in: Lisp and Functional Programming, ACM, 1980, pp. 136–143.
- [4] B. Emir, M. Odersky, J. Williams, Matching objects with patterns, in: European Conference on Object-Oriented Programming, Springer, 2007, pp. 273–298.
- [5] P.-E. Moreau, C. Ringeissen, M. Vittek, A pattern matching compiler for multiple target languages, in: Conference on Compiler Construction, Vol. 2622 of Lecture Notes in Computer Science (LNCS), Springer, 2003, pp. 61–76.
- [6] J. Liu, A. C. Myers, JMatch: Iterable abstract pattern matching for Java, in: Proceedings of the International Symposium on Practical Aspects of Declarative Languages, Springer, London, UK, 2003, pp. 110– 127.
- [7] D. Syme, G. Neverov, J. Margetson, Extensible pattern matching via a lightweight language extension, in: Proceedings of the International Conference on Functional Programming (ICFP), ACM, New York, NY, USA, 2007, pp. 29–40.
- [8] P. Wadler, Views: A way for pattern matching to cohabit with data abstraction, in: Proceedings of the Principles of Programming Languages (POPL), ACM, New York, NY, USA, 1987, pp. 307–313.
- [9] B. Liskov, J. Guttag, Program Development in Java: Abstraction, Specification, and Object-Oriented Design, Addison-Wesley, Boston, MA, USA, 2000.

- [10] C. A. R. Hoare, Proof of correctness of data representations, Acta Informatica 1 (1972) 271–281.
- [11] K. Matsuda, M. Wang, Z. Hu, J. Gibbons, prototype implementation of RINV. [link]. URL http://www.kb.ecei.tohoku.ac.jp/~kztk/viewpp/
- [12] R. Milner, An algebraic definition of simulation between programs, in: Proceedings of the 2nd international joint conference on Artificial intelligence, Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1971, pp. 481–489. URL http://portal.acm.org/citation.cfm?id=1622876.1622926
- [13] R. S. Bird, The promotion and accumulation strategies in transformational programming, ACM Transactions on Programming Languages and Systems 6 (4) (1984) 487–504.
- [14] P. Wadler, The concatenate vanishes, Tech. rep., University of Glasgow (1987).
- [15] L. G. L. T. Meertens, Algorithmics: Towards programming as a mathematical activity, in: CWI Symposium on Mathematics and Computer Science, no. 1 in CWI-Monographs, North-Holland, 1986, pp. 289–344.
- [16] R. S. Bird, An introduction to the theory of lists, in: M. Broy (Ed.), Logic of Programming and Calculi of Discrete Design, Springer-Verlag, 1987, pp. 3–42, nATO ASI Series F Volume 36. Also available as Technical Monograph PRG-56, from the Programming Research Group, Oxford University.
- [17] M. R. Sleep, S. Holmström, A short note concerning lazy reduction rules for append, Software: Practice and Experience 12 (11) (1982) 1082– 1084.
- [18] G. L. Steele, Jr., Organizing functional code for parallel execution or, foldl and foldr considered slightly harmful, in: Proceedings of the International Conference on Functional Programming (ICFP), ACM, New York, NY, USA, 2009, pp. 1–2.

- [19] S.-C. Mu, Z. Hu, M. Takeichi, An injective language for reversible computation, in: Mathematics of Program Construction, Vol. 3125 of Lecture Notes in Computer Science, Springer, 2004, pp. 289–313.
- [20] S.-C. Mu, Z. Hu, M. Takeichi, An algebraic approach to bi-directional updating, in: ASIAN Symposium on Programming Languages and Systems (APLAS), Vol. 3302 of Lecture Notes in Computer Science, Springer, 2004, pp. 2–18.
- [21] J. N. Foster, M. B. Greenwald, J. T. Moore, B. C. Pierce, A. Schmitt, Combinators for bidirectional tree transformations: A linguistic approach to the view update problem, ACM Transactions on Programming Languages and Systems 29 (3), preliminary version in POPL '05.
- [22] A. Alimarine, S. Smetsers, A. van Weelden, M. van Eekelen, R. Plasmeijer, There and back again: Arrows for invertible programming, in: Haskell Workshop, ACM, New York, NY, USA, 2005, pp. 86–97.
- [23] H. Pacheco, A. Cunha, Generic point-free lenses, in: C. Bolduc, J. Desharnais, B. Ktari (Eds.), Mathematics of Program Construction, Vol. 6120 of Lecture Notes in Computer Science, Springer Berlin / Heidelberg, 2010, pp. 331–352.
- [24] R. Bird, O. de Moor, Algebra of Programming, International Series in Computer Science, Prentice Hall, 1997.
- [25] S.-C. Mu, R. Bird, Theory and applications of inverting functions as folds, Science of Computer Programming 51 (2004) 87–116. URL http://portal.acm.org/citation.cfm?id=1007972.1007977
- [26] M. Fokkinga, E. Meijer, Program calculation properties of continuous algebras, Tech. Rep. CS-R9104, CWI, Amsterdam, Netherlands (Jan. 1991). URL http://dbappl.cs.utwente.nl/Publications/PaperStore/ db-utwente-0000003528.pdf
- [27] V. Capretta, T. Uustalu, V. Vene, Recursive coalgebras from comonads, Information and Computation 204 (2006) 437–468.

- [28] C. S. Lee, N. D. Jones, A. M. Ben-Amram, The size-change principle for program termination, in: Proceedings of the Principles of Programming Languages (POPL), ACM, New York, NY, USA, 2001, pp. 81–92.
- [29] D. Sereni, Termination analysis and call graph construction for higherorder functional programs, in: N. Ramsey (Ed.), Proceedings of the International Conference on Functional Programming (ICFP), ACM Press, 2007, pp. 71–84.
- [30] C. Okasaki, An overview of Edison, Electronic Notes in Theoretical Computer Science 41 (1) (2001) 60 - 73, 2000 ACM SIGPLAN Haskell Workshop (Satellite Event of PLI 2000). URL http://www.sciencedirect.com/science/article/pii/ S1571066105805468
- [31] K. Morita, A. Morihata, K. Matsuzaki, Z. Hu, M. Takeichi, Automatic inversion generates divide-and-conquer parallel programs, in: Proceedings of the ACM SIGPLAN Conference on Programming Language Design and Implementation, ACM, New York, NY, USA, 2007, pp. 146– 155.
- [32] B. Liskov, S. Zilles, Programming with abstract data types, in: ACM Symposium on Very High Level Languages, 1974, pp. 50–59.
- [33] D. Coutts, R. Leshchinskiy, D. Stewart, Stream fusion: from lists to streams to nothing at all, in: Proceedings of the International Conference on Functional Programming (ICFP), ICFP '07, ACM, New York, NY, USA, 2007, pp. 315–326.
- [34] N. Wirth, Program development by stepwise refinement, Communications of the ACM 14 (1971) 221–227.
- [35] O. J. Dahl, E. W. Dijkstra, C. A. R. Hoare (Eds.), Structured programming, Academic Press Ltd., London, UK, UK, 1972.
- [36] F. W. Burton, R. D. Cameron, Pattern matching with abstract data types, Journal of Functional Programming 3 (2) (1993) 171–190.
- [37] C. Okasaki, Views for Standard ML, in: ACM Workshop on ML, 1998, pp. 14–23.

- [38] M. Erwig, Active patterns, in: 8th Int. Workshop on Implementation of Functional Languages, Vol. 1268 of Lecture Notes in Computer Science, Springer, 1996, pp. 21–40.
- [39] P. Palao Gostanza, R. Peña, M. Núñez, A new look at pattern matching in abstract data types, in: Proceedings of the International Conference on Functional Programming (ICFP), ACM, New York, NY, USA, 1996, pp. 110–121.
- [40] M. Erwig, S. Peyton Jones, Pattern guards and transformational patterns, in: Haskell Workshop, ACM, New York, NY, USA, 2000.
- [41] M. Tullsen, First class patterns, in: Practical Aspects of Declarative Languages, Vol. 1753 of Lecture Notes in Computer Science, Springer, 2000, pp. 1–15.
- [42] C. B. Jay, The pattern calculus, ACM Transactions on Programming Languages and Systems 26 (6).
- [43] D. Licata, S. Peyton Jones, View patterns: lightweight views for Haskell, http://hackage.haskell.org/trac/ghc/wiki/ViewPatterns (2007).
- [44] P. Nogueira, J. J. Moreno-Navarro, Bialgebra views: A way for polytypic programming to cohabit with data abstraction, in: Workshop on Generic Programming, ACM, New York, NY, USA, 2008, pp. 61–73.
- [45] E. Meijer, M. Fokkinga, R. Paterson, Functional programming with bananas, lenses, envelopes and barbed wire, in: Functional Programming Languages and Computer Architecture, Vol. 523 of Lecture Notes in Computer Science (LNCS), Springer-Verlag, 1991, pp. 124–144.
- [46] L. Meertens, Paramorphisms, Formal Aspects of Computing 4 (1992) 413–424.
- [47] T. Uustalu, V. Vene, Primitive (co-)recursion and course-of-value (co-)iteration, categorically, Informatica 3 (1999) 5–26.
- [48] C. McBride, J. McKinna, The view from the left, Journal of Functional Programming 14 (1) (2004) 69–111.

- [49] Z. Hu, S.-C. Mu, M. Takeichi, A programmable editor for developing structured documents based on bidirectional transformations, in: Proceedings of the Workshop on Partial Evaluation and Program Manipulation (PEPM), ACM, New York, NY, USA, 2004, pp. 178–189.
- [50] The 2LT Team, 2LT two level transformation. URL http://code.google.com/p/2lt/
- [51] J. Hughes, Generalising monads to arrows, Science of Computer Programming 37 (1-3) (2000) 67–111.
- [52] A. Bohannon, J. N. Foster, B. C. Pierce, A. Pilkiewicz, A. Schmitt, Boomerang: Resourceful lenses for string data, in: Principles of Programming Languages, ACM, New York, NY, USA, 2008, pp. 407–419.
- [53] J. N. Foster, A. Pilkiewicz, B. C. Pierce, Quotient lenses, in: Proceedings of the International Conference on Functional Programming (ICFP), ACM, New York, NY, USA, 2008, pp. 383–396.
- [54] J. N. Foster, B. C. Pierce, S. Zdancewic, Updatable security views, in: CSF '09: Proceedings of the 2009 22nd IEEE Computer Security Foundations Symposium, IEEE Computer Society, Washington, DC, USA, 2009, pp. 60–74.
- [55] A. Frisch, G. Castagna, V. Benzaken, Semantic subtyping: Dealing set-theoretically with function, union, intersection, and negation types, Journal of the ACM 55 (4).
- [56] K. Matsuda, Z. Hu, K. Nakano, M. Hamana, M. Takeichi, Bidirectionalization transformation based on automatic derivation of view complement functions, in: Proceedings of the International Conference on Functional Programming (ICFP), ACM, New York, NY, USA, 2007, pp. 47–58.
- [57] J. Voigtländer, Bidirectionalization for free! (Pearl), in: Proceedings of the Principles of Programming Languages (POPL), ACM, New York, NY, USA, 2009, pp. 165–176.
- [58] J. Voigtländer, Z. Hu, K. Matsuda, M. Wang, Combining syntactic and semantic bidirectionalization, in: Proceedings of the International

- Conference on Functional Programming (ICFP), ACM, New York, NY, USA, 2010, pp. 181–192.
- [59] S. Abramov, R. Glück, The universal resolving algorithm and its correctness: inverse computation in a functional language, Science of Computer Programming 43 (2-3) (2002) 193 229.
- [60] M. H. Sørensen, R. Glück, N. D. Jones, A positive supercompiler, Journal of Functional Programming 6 (6) (1996) 811–838.
- [61] E. Albert, G. Vidal, The narrowing-driven approach to functional logic program specialization, New Gen. Comput. 20 (2002) 3-26. URL http://portal.acm.org/citation.cfm?id=586507.586509
- [62] S. Abramov, R. Glück, Y. Klimov, An universal resolving algorithm for inverse computation of lazy languages, in: Proceedings of the 6th international Andrei Ershov memorial conference on Perspectives of systems informatics, PSI'06, Springer-Verlag, Berlin, Heidelberg, 2007, pp. 27– 40. URL http://portal.acm.org/citation.cfm?id=1760700.1760708
- [63] K. Matsuda, S.-C. Mu, Z. Hu, M. Takeichi, A grammar-based approach to invertible programs, in: Proceedings of the European Symposium on Programming, 2010, pp. 448–467.
- [64] P. Wadler, Deforestation: transforming programs to eliminate trees, in: Proceedings of the Second European Symposium on Programming, North-Holland Publishing Co., Amsterdam, The Netherlands, The Netherlands, 1988, pp. 231–248.
 URL http://portal.acm.org/citation.cfm?id=80098.80104
- [65] M. Wang, J. Gibbons, K. Matsuda, Z. Hu, Gradual refinement: Blending pattern matching with data abstraction, in: C. Bolduc, J. Desharnais, B. Ktari (Eds.), Mathematics of Program Construction, Vol. 6120 of Lecture Notes in Computer Science, Springer Berlin / Heidelberg, 2010, pp. 397–426.
- [66] The GHC Team, Glasgow Haskell Compiler. URL http://www.haskell.org/ghc/

[67] T. Sheard, S. Peyton Jones, Template metaprogramming for Haskell, in: M. M. T. Chakravarty (Ed.), Haskell Workshop, ACM Press, 2002, pp. 1–16.

Appendix A. Prototype Implementation

Our system is implemented as a preprocessor for the Glasgow Haskell Compiler, taking an input script and producing a Haskell program. The running example from the paper—of refactoring queues to use a non-list representation—can be fed to our system as a script structured as shown in Figure A.1. We deliberately use **typewriter** font rather than proper typesetting, because at the moment we see the input as a plain text script instead of a program in a rigorously defined extended language. Some familiarity with Haskell [2, 66] and Template Haskell [67] is required to understand the details of this section.

Appendix A.1. Invoking Our System

The preprocessor can be applied from the command line as follows.

ghci -F -pgmF refactor QueueDemo.hs

Here, refactor is the name of the executable file of our prototype implementation, and QueueDemo.hs the input script, structured as shown in Figure A.1. This invokes the interactive version of the Glasgow Haskell Compiler ghci, and loads the translated code.

Appendix A.2. The Input File

As shown in Figure A.1, the input file consists mainly of three blocks, labelled view, abstract and concrete, specifying the intended refactoring; it also has some client code making use of the refactoring. In the following, we explain the various blocks in turn.

Appendix A.2.1. The view Block

The purpose of the view block, shown in Figure A.2, is to declare the abstract and concrete representations, and define the abstraction function. In the queue example, the abstract representation Queue a is implemented in terms of the concrete representation Qi a; the conversion between them is defined by RINV function abstraction. Note that the actual declarations of datatypes Queue and Qi are imported from separate modules Queue and QueueImpl, due to a restriction imposed by Template Haskell.

In the definition of abstraction, we use an ASCII encoding of standard RINV combinators: we write '<**>' for \times , '.' for \circ , and '<||>' for ∇ . We use Template Haskell to generate appropriate fold functions and constructor

```
{-# OPTIONS -XTemplateHaskell -XMultiParamTypeClasses #-}
import UpdatableViewStub
import Queue
import QueueImpl

view Qi a as Queue a by
   abstraction = ...
   completeSets = ...
abstract
   ...
concrete
   ...
play1 = ...
```

Figure A.1: An example input file: QueueDemo.hs

```
view Qi a as Queue a by
abstraction = $(foldT ''Qi) (121 . app . (id <**> reverse))
    where 121 = $(foldT ''[]) ( $(toC 'None) <||> $(toC 'More) )
    completeSets _ = [$(namesT ''Qi), $(namesT ''Queue)]
```

Figure A.2: An example input file: the view block

functions for given datatypes. For example, foldT ''Qi generates the fold for type Qi, and toC 'None generates the constructor function for constructor None. Note that ''' refers to the *constructor* name and '''' refers to the *type* name; for example, '[] means the name of the data constructor [] while '''[] means the name of the type constructor [].

The view block also contains the definition of completeSets, which specifies the constructor grouping introduced with the new datatypes.

Our preprocessor translates the view block into a RINV program, as described in Appendix B; the surjectivity check described in Appendix C guarantees the validity of the generated program, which is done statically through Template Haskell.

Appendix A.2.2. The abstract and concrete Blocks

The abstract and concrete blocks, shown in Figure A.3, contain the definitions of library functions on the specification and implementation level respectively. The specifications are simply copied to the output, while conversions are inserted into the implementations.

Instead of explicitly using functors, as in Section 3, we overload the abstraction functions so that the Haskell type class mechanism will implicitly choose the appropriate instances. This is the reason for the type annotations of the library functions (no type inference or type checking is performed in the prototype). Another technical detail is that we group the α/α° pair into a single function toView, so that toView f is equivalent to $\alpha \circ f \circ \alpha^{\circ}$

Appendix A.2.3. Refactoring Programs

The view, abstract and concrete blocks prepare a refactoring system, which can be used to implement specifications. In the client code part of the example input file (shown in Figure A.4), play1 represents the $play_{prim}$ function from Section 3.2. As explained in Section 3.3, $play_{prim}$ does not use the abstract representation at all; consequently, it is translated as if it were a library function. In Section 3.3 we explained how this absence of abstract representation can be inferred through type checking. As we perform no type checking in the prototype, such functions have to be annotated by their types and the libary functions they actually call.

Our system produces the code in Figure A.5 for play1. A new function play1_impl (\underline{play}_{prim} from Section 3.2) is the implementation of play1 using the concrete representation. Calls to play1 are then redirected to play1_impl

```
abstract
 emptyQueue :: Queue a
 emptyQueue = None
 enQ :: a -> Queue a -> Queue a
 enQ = More
 deQ :: Queue a -> Queue a
 deQ (More _ x) = x
 first :: Queue a -> a
 first (More a _{-}) = a
concrete
 emptyQueue :: Qi a
 emptyQueue = Qi [] []
 enQ :: a -> Qi a -> Qi a
 enQ a (Qi x y) = Qi (a:x) y
 deQ :: Qi a -> Qi a
 deQ (Qi (a:x) y) = Qi x y
 deQ (Qi [] y) = deQ (Qi (reverse y) [])
 first :: Qi a -> a
 first (Qi (a:x) y) = a
 first (Qi [] y) = first (Qi (reverse y) [])
```

Figure A.3: An example input file: the abstract and concrete blocks

Figure A.4: An example input file: client code

after proper conversions. The mapQueue definition is not refactored—it is copied straight to the output.

```
play1_impl =
    let
        play1 0
                 q = first q
        play1 (n+1) q = do hd
                           play1 n (enQ hd tl)
           where hd = first q
                 tl = deQ
    in play1
    where
        enQ = enQ_impl
        first = first_impl
        deQ = deQ_{impl}
play1 = (toView :: (Int -> Qi (IO ()) -> IO ()) ->
                   (Int -> Queue (IO ()) -> IO ()))
        play1_impl
```

Figure A.5: An example of refactored output

Appendix B. RINV Translation

This section summarizes the translation of RINV to Haskell that was described in Section 4. In the clauses below, the use of *Maybe* monad for failure handling is made explicit. As we can see, the only interesting use of the *Maybe* type is in the implementation of *choice*, where the right operand is chosen when the execution of the left operand fails (i.e., returning *Nothing*). The remaining clauses are simply lifted into the *Maybe* monad.

```
\llbracket \cdot \rrbracket :: (a \leftrightarrows b) \to a \to Maybe \ b\llbracket \cdot \rrbracket^{\circ} :: (a \leftrightarrows b) \to b \to Maybe \ a
```

Appendix B.1. Constructors

```
\begin{array}{ll} nil :: () \leftrightarrows [a] \\ \llbracket nil \rrbracket &= \lambda() \to return \ [] \\ \llbracket nil \rrbracket^{\circ} &= \lambda x \to \mathbf{case} \ x \ \mathbf{of} \ \{[] \to return \ () \ ; \_ \to mfail \ "Unmatched" \} \\ cons :: (a, [a]) \leftrightarrows [a] \\ \llbracket cons \rrbracket &= \lambda(x, xs) \to return \ (x : xs) \end{array}
```

Appendix B.2. Primitives

```
id :: a \leftrightarrows a
[id] = return
[id]^{\circ} = return
assocr :: ((a, b), c)) \leftrightarrows (a, (b, c))
[assocr] = \lambda((a,b),c) \rightarrow return(a,(b,c))
\llbracket assocr \rrbracket^{\circ} = \llbracket assocl \rrbracket
assocl :: (a, (b, c)) \leftrightarrows ((a, b), c)
\llbracket assocl \rrbracket = \lambda(a, (b, c)) \rightarrow return((a, b), c)
[assocl]^{\circ} = [assocr]
swap :: (a, b) \leftrightarrows (b, a)
\llbracket swap \rrbracket = \lambda(a,b) \rightarrow return(b,a)
\llbracket swap \rrbracket^{\circ} = \llbracket swap \rrbracket
fst :: (a \to b) \to (a, b) \leftrightarrows a
[\![\mathit{fst}]\!] = \lambda g \to \mathit{return} \circ \mathit{fst}
[fst]^{\circ} = \lambda g \to \lambda x \to return (x, g x)
snd :: (b \to a) \to (a, b) \leftrightarrows b
\llbracket snd \rrbracket = \lambda g \rightarrow return \circ snd
[snd]^{\circ} = \lambda g \rightarrow \lambda x \rightarrow return (g x, x)
```

Appendix B.3. Combinators

```
(\circ) :: (b \leftrightarrows c) \to (a \leftrightarrows b) \to (a \leftrightarrows c)
\llbracket f \circ g \rrbracket = \lambda x \to \llbracket f \rrbracket \ x \ggg \llbracket g \rrbracket
\llbracket f \circ g \rrbracket^{\circ} = \lambda z \to \llbracket g \rrbracket^{\circ} z \gg \llbracket f \rrbracket^{\circ}
(\times) :: (a \leftrightarrows b) \to (c \leftrightarrows d) \to ((a, c) \leftrightarrows (b, d))
\llbracket f \times g \rrbracket = \lambda(w, x) \rightarrow liftM2 \ (,) \ (\llbracket f \rrbracket \ w) \ (\llbracket g \rrbracket \ x)
\llbracket f \times g \rrbracket^{\circ} = \lambda(y, z) \rightarrow liftM2 \ (,) \ (\llbracket f \rrbracket^{\circ} y) \ (\llbracket g \rrbracket^{\circ} z)
(\nabla) :: (a \leftrightarrows c) \to (b \leftrightarrows c) \to (Either \ a \ b \leftrightarrows c)
\llbracket f \triangledown g \rrbracket = \lambda x \to \mathbf{case} \ x \ \mathbf{of} \ \{ Left \ a \to \llbracket f \rrbracket \ a \ ; Right \ b \to \llbracket g \rrbracket \ a \}
\llbracket f \triangledown g \rrbracket^{\circ} = \lambda x \rightarrow \llbracket f \rrbracket^{\circ} x \text{ `choice' } \llbracket g \rrbracket^{\circ} x
     where choice (Just x) _{-}
                      choice \ \_  (Just \ x) = Just \ x
                      choice _
                                                                      = Nothing
fold_{List} :: (Either () (a, b) \leftrightarrows b) \rightarrow \lceil a \rceil \leftrightarrows b
\llbracket fold_{List} f \rrbracket = foldM_{List} \llbracket f \rrbracket
\llbracket fold_{List} f \rrbracket^{\circ} = unfoldM_{List} \llbracket f \rrbracket^{\circ}
foldM_{List} :: (Either\ ()\ (a,b) \to Maybe\ b) \to [a] \to Maybe\ b
foldM_{List} f [] = f (Left ())
foldM_{List} f(a:x) = foldM_{List} f x \gg \lambda r \rightarrow f(Right(a,r))
unfoldM_{List} :: (b \rightarrow Maybe (Either () (a, b))) \rightarrow b \rightarrow Maybe [a]
unfoldM_{List} f b =
     f \ b \gg \lambda r \to \mathbf{case} \ r \ \mathbf{of}
          Left() \longrightarrow return[]
          Right(a,b) \rightarrow unfoldM_{List} f b \gg \lambda x \rightarrow return(a:x)
```

Appendix C. Pseudocode for Surjectivity Check

This section describes the algorithm for the surjectivity check in RINV. Function checkSurj checks whether a constructor or a set of constructors (collected from operands of ∇) 'covers' a complete constructor set. In the case of ∇ , when the two operands 'cover' a complete constructor set, we do not require them individually to be surjective. However, we still need to make sure that compositions with the operands are properly constructed (i.e., only surjective functions on the right of a composition). This is the reason for invoking checkComp, which in turn relies on checkSurj for checking the surjectivity of its component. One property of the two functions is that checkSurj x implies checkComp x.

```
checkSurj (constructor) = isComplete constructor completeSets
checkSurj (primitive)
                          = True
checkSurj (f \nabla g)
                          = (isComplete (f \nabla g) completeSets \vee
                              checkSurj \ f \lor checkSurj \ q) \land
                             checkComp\ f \land checkComp\ g
checkSurj\ (f \circ q)
                          = checkSurj f \wedge checkSurj g
checkSurj\ (f \times q)
                          = checkSurj f \wedge checkSurj g
checkSurj (fold f)
                          = checkSurj f
checkComp\ (constructor) = True
checkComp\ (primitive)
                            = True
                            = checkComp \ f \land checkComp \ g
checkComp\ (f \ \nabla \ g)
checkComp\ (f \circ q)
                            = checkComp \ f \land checkSurj \ g
                            = checkComp \ f \land checkComp \ g
checkComp\ (f \times g)
checkComp (fold f)
                            = checkComp f
```

In the above, $isComplete\ f\ cs$ checks whether the constructors provided by f cover one of the constructor sets in cs. We implement sets as lists, but still use the set notations for presentation. When constructors are paired into tuples, the behaviour of isComplete is slightly more complicated, as we need to compute the cartesian products of the constructor sets.

```
isComplete \ f \ cs =
any \ (\lambda s \rightarrow s \subseteq (flatten \ f)) \ complete Tuple Sets
where
```

```
\begin{aligned} \dim &= \operatorname{length} \$ \left( \operatorname{flatten} f \right) !! \, 0 \\ &\operatorname{completeTupleSets} \\ &= \operatorname{foldNat} \left( \lambda r \to \left\{ \left\{ a : b \mid a \leftarrow c, b \leftarrow d \right\} \mid c \leftarrow cs, d \leftarrow r \right\} \right) \\ &\left\{ \left\{ \left[ \right] \right\} \right\} \\ &\operatorname{dim} \\ &\operatorname{foldNat} s \ n \ 0 = n \\ &\operatorname{foldNat} s \ n \ k = s \left( \operatorname{foldNat} s \ n \ (k-1) \right) \end{aligned}
```

In the above, completeTupleSets is the pairwise cartesian products of the elements in cs up to a given dimension. For example, given a set of sets of constructors

```
\{\{nil, cons\}, \{nil, snoc\}, \{nil, wrap, app\}\}
```

a complete Tuple Sets with dimension 2 is

```
 \begin{split} & \{ \{[nil,nil],[nil,cons],[cons,nil],[cons,cons] \}, \\ & \{ [nil,nil],[nil,snoc],[cons,nil],[cons,snoc] \}, \\ & \{ [nil,nil],[nil,wrap],[nil,app],[cons,nil],[cons,wrap],[cons\ app] \}, \\ & \ldots \} \end{split}
```

The dimension represents the number of elements in a possibly nested tuple, for which we flatten nested tuples into lists.

```
\begin{array}{ll} \textit{flatten (constructor)} = \{ [\textit{constructor}] \} \\ \textit{flatten } (f \triangledown g) &= \textit{flatten } f \cup \textit{flatten } g \\ \textit{flatten } (f \times g) &= \{ \textit{nf} + \textit{ng} \mid \textit{nf} \leftarrow \textit{flatten } f, \textit{ng} \leftarrow \textit{flatten } g \} \\ \textit{flatten } (f \circ g) &= \textit{flatten } f \\ \textit{flatten} \_ &= \{ [] \} &-- \text{ primitive functions are ignored} \end{array}
```

The nesting of \times increases the length of individual elements (dimension) in a complete set, whereas the nesting of ∇ increases the size of the complete sets. For example, we have the following executions

```
flatten\ (nil \times cons) = \{[nil, cons]\} flatten\ (nil\ \nabla\ cons) = \{[nil], [cons]\}
```

and

```
\mathit{flatten}\ ((\mathit{nil} \times \mathit{cons}) \ \triangledown \ (\mathit{cons} \times \mathit{nil})) = \{[\mathit{nil}, \mathit{cons}], [\mathit{cons}, \mathit{nil}]\}
```

Among the above, $\{[nil], [cons]\}$ is a complete set, so $nil \nabla cons$ passes the isComplete test but the other ones fail. On the other hand, $(nil \times nil) \nabla (nil \times cons) \nabla (cons \times nil) \nabla (cons \times cons)$ passes the isComplete test, because it flattens into $\{[nil, nil], [nil, cons], [cons, nil], [cons, cons]\}$, which covers a complete set of dimension 2.