

Parallelizing Imperative Functional Programs: the Vectorization Monad

JONATHAN M.D. HILL[†], KEITH M. CLARKE[‡] AND RICHARD BORNAT[‡]

[†]Oxford Parallel, Oxford University Computing Laboratory, U.K. [‡]Department of Computer Science, Queen Mary & Westfield College, U.K.

(Received 30 May 1995)

Traditionally a vectorizing compiler matches the iterative constructs of a program against a set of predefined templates. If a loop contains no dependency cycles then a map template can be used; other simple dependencies can often be expressed in terms of fold or scan templates. This paper addresses the template matching problem within the context of functional programming. A small collection of program identities are used to specify vectorizable for-loops. By incorporating these program identities within a monad, all well-typed for-loops in which the body of the loop is expressed using the vectorization monad can be vectorized. This technique enables the elimination of template matching from a vectorizing compiler, and the proof of the safety of vectorization can be performed by a type inference mechanism.

© 1996 Academic Press Limited

1. Introduction

"'Should declarative languages be mixed with imperative languages?', clearly has the answer that they should, because at the moment we don't know how to do everything in pure declarative languages." (Strachey, 1966).

It has long been known that some of the most common uses of for-loops in imperative programs can easily be expressed using the standard higher-order functions *fold* and *map*. With this correspondence as a starting point, we derive parallel implementations of various iterative constructs, each having a better complexity than their sequential counterparts, and explore the use of monads to guarantee the soundness of the parallel implementation.

As an aid to the presentation of the material, we use the proposed syntax for parallel Haskell (Nikhil et al., 1995), shown in Figure 1, as a vehicle in which imperative functional programs will be expressed. Incorporating imperative features into a purely functional language has become an active area of research within the functional programming community (Peyton Jones and Wadler, 1993; Launchbury, 1993; Talpin, 1993). One of the techniques gaining widespread acceptance as a model for imperative functional programming is monads (Moggi, 1989; Wadler, 1990). Typically monads are used to guarantee single threadedness, enabling side effects to be incorporated into a purely functional language without losing referential transparency. We take a different approach. First,

```
exp \mapsto for pat \leftarrow exp do \{next \ var = exp\}^+ finally \ exp
```

Figure 1. Proposed syntax extensions for pH.

for-loops are translated into a monadic framework. Next, by ensuring that the monad satisfies the programming identities usually associated with the successful vectorization of imperative constructs, all well-typed for-loops in which the body is expressed using the vectorization monad can be parallelized. The technique is not just restricted to a data-parallel environment. It could be implemented by a divide and conquer technique on a multi-processor platform, or by a parallel implementation of graph reduction.

2. Background

We use a model of parallel computation based upon the Bird Meertens Formalism (Bird, 1986). To make things simple, all parallelism is expressed in terms of the map, fold, and scan functions. Given a function f which has a $\mathcal{O}(1)$, then (map f) is $\mathcal{O}(1)$, while (scanPar f) and (foldPar f) are each $\mathcal{O}(\log N)$ if f is associative. Unlike our earlier work, we use the BMF as is—all the operations presented here can be interpreted as though they were being expressed on lists. This stance contradicts an earlier paper (Hill, 1993a) in which we said that lists were unsuitable for data-parallel evaluation in a non-strict language. We still believe this to be true, but we have tried to sugar the pill somewhat. By using a purely combinator approach to programming, it is possible to give the impression of using lists, whilst actually implementing these list-like objects on top of more suitable data-parallel data structures (Hill, 1993a, b, 1994)

3. Parallelizing Simple Loops

It is part of folklore that programs expressed as for-loops can be rewritten using tail recursion [for example, Landin (1966) or Henderson (1980)], although the functional programming community has concentrated on how the reverse translation can be used as an optimisation technique (Peyton Jones, 1987). As an example of such a translation, the iterative factorial expressed using the pH syntax shown on the left below:

```
fact n =let ac = 1
    in for i<-[1..n] do
        next ac = ac*i
    finally ac</pre>
fact n = f 1 [1..n]
where
f ac [] =ac
f ac (i:is)=f (ac*i) is
```

can be rewritten as the tail recursive definition on the right. Unfortunately, replacing an inherently sequential for-loop with a sequential tail recursive function does not provide the right foundations for parallelisation. However it is interesting to note that the resulting tail recursive function is an instance of a fold-left computation over the loop range.

If \ominus is a function that describes the computation that occurs at each iteration of a for-loop, then the definition of the factorial function can be rewritten using schema TPH of Figure 2 as a fold-left of \ominus ; where \ominus is multiply in this example:

TPH[next v] = v' where TPH[exp] is the identity translation for other expressions. TPAT[next v] = v' where TPAT[pat] is the identity translation for other patterns.

Figure 2. Translating a pH for-loop.

```
fact n = let ac = 1

in case (foldl (\ac i->let ac'=ac*i in ac') ac [1..n])

ac -> ac

\Rightarrow fact n = foldl (\ac i -> ac * i) 1 [1..n]
```

We term \ominus the next-function of the for-loop. As we have already mentioned, a parallel fold has $\mathcal{O}(\log N)$ time complexity, if the function being folded is $\mathcal{O}(1)$ and associative. Whenever \ominus and the initial state of the for-loop form a monoid (\ominus and x form a monoid if \ominus is associative and has x as its left and right identity), then by the *first duality theorem* (Bird and Wadler, 1988) the fold-left can be safely transformed into a fold-right, or more interestingly a parallel $\mathcal{O}(\log N)$ fold (Hill, 1993a).

If this theorem is used to transform a fold-left into a right or parallel fold, then the side condition requiring finite lists is ignored. The justification for this is that in all situations where a fold-left yields a non-bottom value, a fold-right will give the same result. However, there may be situations where fold-left diverges and fold-right terminates. We regard this as a good thing, in the same way that we believe normal order reduction to be superior to applicative order reduction.

4. Relaxing Associativity

When the next-function \ominus is associative, the left argument used to represent the state of the loop, the right argument that represents the loop counter, and the result of the next-function that represents the successive state of the loop body, all have the same type (by the definition of associativity). This is rather unfortunate, as the numerical algorithms we intend to parallelize typically range over a subset of the integers, and the loop's state will be a mixture of floating point and integer values. In the situation where the body of the loop contains a single state (i.e., there is only one next binding in the body of the loop), a solution is to decompose \ominus into a part that is specific to the

computation of the loop-counter, and a remainder that is specific to the loop's state. If \ominus has the structure $\lambda s i \to s \oplus f i$, where \oplus and the initial value of the for-loop now form a monoid, then the part of the computation that is specific to the loop counter can be moved outside the fold-left computation by using the *fold-map fusion* law[†] (Bird, 1989).

$$foldl\left(\lambda s\:i\to s\oplus fi\right)x\equiv foldl\left(\oplus\right)x\circ map\:f.$$
 Identity 2: "fold-map fusion law"

Comparing this law to a conventional optimisation on loops, the inverse of loop fusion (Aho et al., 1986) is being performed in a scenario that would be detrimental to performance in a sequential environment as the overheads associated with evaluation of a loop will be incurred twice. Because of the differing complexities of fold and map in a data-parallel environment, the law is an optimization technique.

5. The Leaky Fold-left Law

The aim of this section is to investigate the operational phenomenon of space-leaks (Peyton Jones, 1987; Sparud, 1993), and show how fold-left computations are particularly susceptible to leaking space. A collection of identities are used to highlight the leaky behaviour of fold-left. From these identities, a parallel fold is developed that is not just restricted to associative operators. Unfortunately fold computations parallelized in this way suffer from a similar problem to a space leak—a vectorization leak. The purpose of this section is to therefore highlight some of the problems of parallel fold, and the techniques developed here form the foundations for the vectorization based upon monads later in the paper.

Given the problem of finding the sum of a list of numbers, there is the opportunity to use a left or right fold, as the addition used in the sum is associative. A Scheme or ML programmer would probably side for a fold-left computation because of its tail recursive nature. A non-strict functional programmer has a dilemma! For each iteration of a fold-left, a closure is created in the accumulated parameter of the fold that represents the application of the folded function to the previous value of the accumulator. This closure remains unevaluated, growing whilst the spine of the folded list is unwound. Only when the end of the list is reached, and the closure's size is proportional to the length of the list, is the closure evaluated—this is the space leak. A clever compiler should be able to eliminate this 'dragging' if the folded function is known to be strict, because some real computation could occur in the accumulated parameter.

Putting such dilemmas to one side, what this detour into the operational characteristics of fold-left has exposed is yet another way of thinking of fold. As opposed to the normal space-leaking behaviour of fold-left being interwoven with the vagaries of a non-strict evaluation mechanism, it can be made a *feature* of fold-left such that an implementation leaks space in both a strict and non-strict language!

$$foldl \ (\oplus) \ x \ xs \equiv (foldr (\lambda s \ i \rightarrow (\lambda y \rightarrow y \oplus i) \circ s) id \ xs) x. \qquad \qquad Identity \ 3: \text{``Leaky fold-left law''}$$

The *leaky fold-left law* gives a correspondence between fold-left and an explicitly leaky version of the function. As this definition is an instance of the fold-map fusion law, it

 $^{^{\}dagger}$ There seems to be a strange anomaly that laws are used in the opposite direction when reasoning about parallel functions, compared to similar functions on lists. The names are borrowed from the list context, even though the term *fold-map fission* would be more appropriate in a data-parallel setting.

can be transformed into the following identity [as in Bird (1986), $\tilde{\oplus}$ is used to represent the flipped version of the operator \oplus , i.e., $\lambda x \, y \to y \oplus x$]:

$$foldl\left(\oplus\right)x\ xs\equiv\left(foldr\left(\circ\right)id(map\left(\overset{\circ}{\oplus}\right)xs\right)\right)x.$$
 Identity 4

The program identity is now an instance of the first duality theorem (because $\tilde{\circ}$ and the identity function form a monoid), and the fold-right can be safely replaced by a left or parallel fold. The identity can be understood by decomposing the right-hand side into three parts:

- 1. $\tilde{\oplus}$ is partially applied to all the elements of the list being folded;
- 2. the resulting list of function values is folded using $\tilde{\circ}$. The effect is to create a closure that is equivalent to the function produced by unrolling $\tilde{\oplus}$ at runtime. e.g., given [1,2,3], the closure $(\oplus 3) \circ (\oplus 2) \circ (\oplus 1)$ is created;
- 3. the unrolled closure is then applied to the initial state x, and only at this point can the real computation of all the \oplus 's commence.

What we have shown is that a for-loop is nothing more than syntactic sugar for a fold-left. By means of a series of program transformations, a fold-left can be transformed into a left or right fold of the compose function, and a map of the original function that we wanted to fold. In Hill (1994) we show how a $\mathcal{O}(1)$ map can be implemented, and because of the associativity of \circ , a parallel fold can be used in the implementation of a for-loop regardless of the associativity of the function modelling the body of the loop. Unfortunately, there is a rather large downside to this transformation. A program implemented in such a manner creates a closure with the same number of compositions as there are loop iterations. Evaluation of this closure has a linear complexity, therefore the complexity of the entire for-loop is $\mathcal{O}(N)$, and not the desired $\mathcal{O}(\log N)$. Fortunately, in a wide class of problems, monads can be used as a solution to this dilemma, and this is addressed in §8.

6. Example: Testing if a Number is Prime

The program identities developed so far are used in the parallelisation of a simple function that determines if a natural number is prime. Looking at the definition of isPrime (Step 0) below, the first step in the parallelisation process is to determine if the body of the for-loop can be expressed by an associative function. If the body is found to be associative, then by the first duality theorem parallelisation of the for-loop would be complete as it could be implemented by a parallel fold.

Unfortunately, the conditional in the body of the loop is rather off-putting. Faced with simple arithmetic operations such as addition and multiplication, the associativity

Step 0: The Start.

of the body is quite simple to determine. However, what is the associativity of an expression that contains a conditional? This issue is side-stepped, as in this context the conditional can be converted into the conjunction of logical operators whose associativity is easy to determine. Using translation scheme TPH, in conjunction with the equivalence if x then False else $y \equiv \neg x \land y$, the definition of isPrime (Step 0) is transformed into the following:

The function used to model the body of the loop now has the structure required by the fold-map fusion law, as shown by Eqn 6.1.

$$\lambda s i \to \neg (n \text{ rem } i = 0) \land s \equiv \lambda s i \to s \oplus f i \quad \text{where} \begin{array}{c} s \oplus t = t \land s \\ f i = \neg (n \text{ rem } i = 0). \end{array}$$
 (6.1)

Applying the program identity transforms isPrime (Step 1) into a form in which the monoid of the flipped && operator and boolean True are used in the fold. This fulfils the side conditions of the first duality theorem, and the definition of fold-left can be replaced by a parallel fold, completing the parallelisation of the function. Contrasting the complexities of the original and the transformed function, to test if the number N is prime, the pH version of isPrime has a $\mathcal{O}(\sqrt{N})$ complexity, compared to the $\mathcal{O}(\log \sqrt{N})$ complexity of the transformed function.

```
isPrime n = fold1 (\ s t -> t && s) True  (\text{map ($\setminus$i -> n 'rem' i /= 0$)} \\ [2 .. truncate (sqrt (fromInteger n))])
```

Step 2: Parallelized with help from the first duality theorem.

7. Example: Inverse Sine

Next a numerical problem is considered, where a series that approximates the trigonometric function $\sin^{-1} x$ is evaluated to any desired accuracy (Eqn 7.1). For a given accuracy, the idea is to convert a sequential algorithm that performs a linear number of expansions of the inverse sine series, into an algorithm with a logarithmic number of parallel expansions.

$$\sin^{-1} x = x + \frac{x^3}{2 \cdot 3} + \frac{1 \cdot 3 \cdot x^5}{2 \cdot 4 \cdot 5} + \frac{1 \cdot 3 \cdot 5 \cdot x^5}{2 \cdot 4 \cdot 6 \cdot 7} + \cdots$$

$$= \sum_{i=0}^{\infty} \frac{x^{2i+1} \prod_{j=1}^{i} (2j-1)}{(\prod_{j=1}^{i} 2j) \cdot (2i+1)}.$$
(7.1)

Although there is little point in parallelising this algorithm as the series rapidly converges, it forms an interesting case study. A naive $\mathcal{O}(N^2)$ sequential algorithm leads to a naive $\mathcal{O}(N\log N)$ parallel algorithm. The naive sequential algorithm may be converted using standard imperative style program transformations into a $\mathcal{O}(N)$ algorithm. Unfortunately the improved sequential algorithm fails to be parallelized with the identities developed so far. This is remedied by a collection of additional laws, and a $\mathcal{O}(\log N)$ parallel algorithm is finally derived.

Step 0: Original Program.

The function sinI (Step 0) defines an implementation that performs a fixed number of expansions of the series of Eqn 7.1 in terms of a pH for-loop. The recursive let-bindings top and bot in the body of the for-loop have a structure that mimics the numerator and denominator of Eqn 7.1. As the body is an instance of the fold-map fusion law, the function can be transformed using the relationship of Eqn 7.2.

$$\lambda s \, i \to s \oplus f \, i \equiv \lambda s \, i \to s + \frac{x^{2i+1} \prod_{j=1}^{i} (2j-1)}{\left(\prod_{j=1}^{i} 2j\right) \cdot (2i+1)}$$

$$s \oplus t = s + t$$

$$\text{iff} \qquad f \, i = \frac{x^{2i+1} \prod_{j=1}^{i} (2j-1)}{\left(\prod_{j=1}^{i} 2j\right) \cdot (2i+1)}.$$
(7.2)

Because floating-point addition and 0.0 form a monoid[†], the first duality theorem can be used to transform the fold-left into a parallel fold.

Step 1: TPH and fold-map fusion

This implementation of inverse sine, perhaps surprisingly, is not $\mathcal{O}(\log N)$. The original function can be thought of as a nested $\mathcal{O}(N^2)$ for-loop, as the computations of both the numerator and denominator are $\mathcal{O}(N)$. The derived parallel algorithm has an outer loop with $\log N$ parallel iterations, each of which contains a $\mathcal{O}(N)$ product, resulting in a $\mathcal{O}(N\log N)$ parallel algorithm. It would be expected that if the inner loop were parallelized in a similar manner to the outer loop, then an $\mathcal{O}(\log^2 N)$ algorithm could be derived. This highlights a fundamental problem with these program identities, and the data-parallel model of computation we assume—parallelisation can only occur at a single level in a program.

Close examination of the series reveals each term to be similar to its predecessor. By taking advantage of this similarity, a common portion of the expression can be retained from one iteration of the loop to the next, and the linear complexity associated with

[†] Floating point addition is not associative because of rounding errors. However, like HPF's (Rice University, 1993) fold functions we are willing to pay the price of potential instabilities in this numerical algorithm, and assume addition is associative.

[‡] Some languages like NESL (Blelloch, 1993) are based around such forms of nested parallelism. Opposed to making a choice at which level parallelisation should occur in a nested expression, we take the conventional approach in such situations of flattening the computation into a single loop.

the product can be eliminated. Starting with the original pH definition sinI (Step 0), a conventional imperative-style program transformation can be performed to produce an improved $\mathcal{O}(N)$ sequential algorithm.

```
sinI x n = let pow = x; s = 0.0; top = 1; bot = 1
    in for i <- [0..n] do
        next pow=pow * x * x
        next s =s+((pow* toFloat top) / (toFloat ((2*i+1)*bot)))
        next top=top * (2 * (i+1) -1)
        next bot=bot * (2 * (i+1))
        finally s</pre>
```

Step 2: Imperative munging.

A cursory investigation of the types of the loop-range and the body of the loop reveals that none of the program identities developed so far are applicable as they rely upon an associative operator, which by definition must have a type of the form $\alpha \to \alpha \to \alpha$. The integer type used as the loop-counter and the tuple of integer and floating point numbers used in the body of the loop means that no associative next-function exists for loops of this form as the types don't 'fit' together.

One solution to this problem is to abstract part of the state of the body of the loop outwards. A way of achieving this is to perform induction-variable elimination (Aho et al., 1986). The idea is to infer if any of the changes to a subset of the states in the body of the loop occur in a lock-step manner (the induction variables). When there are two or more induction variables in a loop, it may be possible to remove all but one of them by abstracting part of the loop's state outwards. We use a generalized scheme that abstracts part of the state outside the loop, whether the state is an induction variable or not. The trick is to split the loop in two, but ensure that the abstracted loop remembers all of its intermediary states—a parallel $\mathcal{O}(\log n)$ scan accomplishes this. All the values of the states in the scanned list are then 'zipped' together with the original loop-range so that each of the values can be used by an expression inside the remaining part of the fold-left.

Although the idea is quite simple, we need to identify which part of the original loop's state should be abstracted out. We are trying to make the types of the loop range and the state of the body of the loop the same. This is achieved by abstracting just those parts which make the zipped range, and the new state in the body of the loop the same; we have to be aware that variables may become free. A topological sort of the variables is used to ensure strongly connected groups of states are lifted out as a whole. However, after applying the fold-map fusion law the type of the body of the loop may change. For example, looking back at step 2, the type of the states in the body of the loop is (Float, Float, Int, Int), and the loop range is an Int. Abstracting all but the state s out of the loop produces a zipped loop range that is susceptible to further simplification.

The body of the loop now has the form required by the fold-map fusion law. Equation 7.3

defines the relationship between the body of the loop and the function required by the law.

$$\lambda s \, i \to s + \frac{pow \times top}{(2 \times i + 1) * bot} \equiv \lambda s \, i \to s \oplus f \, i \qquad \text{iff} \quad \begin{array}{c} s \oplus t = s + t \\ f \, i = \frac{pow \times top}{(2 \times i + 1) * bot}. \end{array}$$
 (7.3)

Applying the law transforms the fold-left into a form where the first duality theorem is applicable. Parallelisation isn't finished as the scans introduced by moving part of the state outside the loop need to be transformed into a parallelisable form. Luckily, the first duality theorem and the fold-map fusion law are applicable to scan computations as well as folds. However, the first use of scan poses a problem as the identity of multiplication is not used as the starting value of the scan. Program identity 5 is taken from Bird and Wadler (1988), and follows from the definition of fold-left. This identity forms the basis of the analogous identity 6 on scans.

The first use of scan in sinI (Step 4) is an instance of program identity 6, where the expression $x \times 1.0$ has the form $x \oplus y$. As 1.0 and \times form a monoid, identity 6 and the fold-map fusion law can be used to transform the function into the form shown in sinI (Step 5), which completes the parallelisation of inverse sine. The original version of the problem is $\mathcal{O}(N^2)$, but can be trivially made linear. As zipWith4 is similar to map and has a $\mathcal{O}(1)$ complexity, the improved linear algorithm can then be transformed into a parallel $\mathcal{O}(\log N)$ algorithm.

8. Category Theory and Monads

Although the program identities used in the prior sections are straightforward, as was seen to be the case in the inverse sine example, the application of the identities in even a relatively simple example can be troublesome. Instead of writing programs which are inferred to be vectorizable, the aim of the rest of the paper is to provide a 'sub-programming' language in which all well-formed and well-typed programs are vectorizable. In the next section a model inspired by monads is presented that enables the program identities developed so far to be incorporated into an object that looks like a monad (for more worked examples of the program identities, and an extended introduction to category theory see Hill and Clarke (1993)).

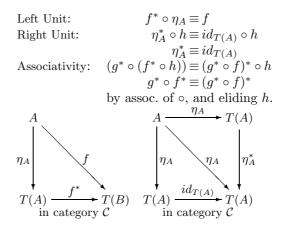
The principle underlying the work of Moggi (1989) on monads and the computational λ -calculus is the distinction between simple data-valued functions and functions that perform computations. A data-valued function is one in which the value returned by the function is determined solely by the values of its arguments. In contrast, a function that performs a *computation* can encompass ideas such as side-effects or non-determinism,

which implicitly produce more results as a consequence of an application of the function than the result explicitly returned.

Given the objects A in the category C, and the endofunctor $T: C \to C$, Moggi's work on monads views the endofunctor T as a mapping between all the objects Obj(C) of the category C which are to be viewed as the set of all values of type τ , to a corresponding set of objects T(Obj(C)) which are to be interpreted as the set of computations of type τ .

Now consider for each morphism $f:A\to T(B)$ a new morphism $f^*:T(A)\to T(B)$, where \bot^* is the "extension" operator that lifts the domain A of the morphism f to a computation T(A). In the context of computations, f is a function from values to computations, whereas f^* is a function from computations to computations. The expression $g^*\circ f$, where $f:A\to T(B)$ and $g:B\to T(C)$ is interpreted as applying f to some value a to produce some computation f a; this computation is evaluated to produce some value b, and g is applied to b to produce a computation as a result.

The Kleisli triple (T_{obj}, η, \bot^*) is defined as the restriction of the endofunctor T to objects, the extension operator \bot^* , and a natural transformation $\eta : id_{\mathcal{C}} \to T$, where $id_{\mathcal{C}} : \mathcal{C} \to \mathcal{C}$ is the identity functor for category \mathcal{C} . In the context of computations, η can be thought of as an operator that includes values into a computation. For the triple to be well formed, the following laws are required to hold; a pictorial presentation is given in the commuting diagrams below, where $h : A \to T(A)$.



The Kleisli triple can be thought of as a different syntactic presentation of a classical category theory monad, as there is a one-to-one correspondence between a Kleisli triple and a monad [see Mac-Lane (1971) for a proof].

The use of monads in the functional programming community bears a closer resemblance to Kleisli triples, than classical monads. Wadler (1990) adapted Moggi's ideas of using monads to structure the semantics of computations into a tool for structuring functional programs. The Kleisli triple $(T_{obj}, \eta, _^*)$ can be defined in a functional language by the Wadler-ized triple $(M, unit, \star)$, where M is a type constructor, \star is a function of type $M\alpha \to (\alpha \to M\beta) \to M\beta$, and $f \star g$ is the same as $g^* \circ f$ of the Kleisli triple. The natural transformation η can be modelled as a polymorphic function as it can be thought of as a family of morphisms from each object in a category to objects in another

 $^{^{\}dagger}\,$ i.e., a functor with a mapping to and from the same category.

category (the components of the natural transformation). A natural transformation is therefore similar to a polymorphic function, and as a consequence η is written as the polymorphic function unit of type $\alpha \to M \alpha$. The laws required by the Kleisli triple can now be recast as:

```
Left Unit: unit \ a \star (\lambda b \to n) \equiv n[a/b]

Right Unit: m \star (\lambda b \to unit \ b) \equiv m

Associativity: m \star ((\lambda a \to n) \star (\lambda b \to o)) \equiv (m \star (\lambda a \to n)) \star (\lambda b \to o).
```

9. A Vectorization Monad

A model of stream-based output can be defined in terms of Moggi (1989) side-effecting monad (example 3.6, page 11). In a simplified scenario where the output stream is a list of characters (a Landin-stream 1965), the monad (IO, unit, bind) is used, where bind is the Haskell identifier that represents \star of the previous section:

The monad operations are augmented with a print computation that outputs its string argument as a 'side-effect' onto the output stream and delivers (). In the context of category theory, given the Kleisli triple (T, η, \bot^*) in the category \mathcal{C} , functions like print are the morphisms $f: A \to T(B)$, and the set of all such morphisms is the homset $\mathcal{C}(\mathcal{A}, \mathcal{T}(\mathcal{B}))$. If the monad is to be interpreted as an abstract data-type, then this hom-set of morphisms forms an interface to the type as it requires 'inside knowledge' of the representation of the monad.

A monadic for-loop is defined to have the structure for $ctr \leftarrow range$ do body, where body is a computation of type IO (). The result returned by a monadic loop expression will always be (), but the state of the for-loop is changed by side-effects during successive iterations of the loop. The state of a pH for-loop is represented by next bindings. In contrast, a monadic loop hides the state, which makes it possible to straightjacket any interactions with the state such that the first duality theorem and the fold-map fusion laws can be satisfied, therefore making vectorization possible. The function helloWorlds shown in loop 1 is an example of a monadic for-loop. The function takes a numeric argument n, and performs n iterations of a for-loop printing the string "Hello World" followed by printing the value of the loop counter by side-effecting the output stream:

In a similar vein to the translation of a pH for-loop, a monadic loop is translated into a fold-left computation. The translation is relatively straightforward: (1) the body of the loop is modelled by a function that is parameterized on the loop counter; (2) a fold-left of the function $\lambda s \, i \to (s \, 'bind' \, (\lambda() \to bodyFn \, i))$ is performed where bodyFn is bound by the function that represents the body of the original loop; (3) during the first iteration of the loop, the 's' argument of the folded lambda expression represents a computation that encapsulates the state on entry to the loop; (4) during the kth iteration of the loop, the argument 's' encapsulates the state of all the previous k-1 iterations; (5) the initial state of the for-loop is a 'do-nothing' computation represented by unit ().

9.1. From monads to monoids

The flaw in the translation of a monadic for-loop into a fold-left is that the monad operation bind used as the folded function has the type $M\alpha \to (\alpha \to M\beta) \to M\alpha$ and is therefore not associative. One solution to this problem is not to use monads, but use a structure that is very similar, and can be used to achieve the same operational behaviour—re-enter monoids!

As each monad operation in the running example returns the same value (), and there is never a situation in which anything other than the computation delivering () can be returned (i.e., the dreaded bottom †), then all the ()s can be elided from the program. Removing the parameterization from the type constructor IO, and changing the monad operations accordingly produces:

This transformation is possible because the extension operator $_^*$ of the Kleisli triple isn't needed in this example. As $f \star g$ is syntactic sugar for $g^* \circ f$, we can note that all the functions used with the monad belong to the hom-set $\mathcal{C}(\mathcal{A}, \mathcal{T}(\{()\}))$. Because we are not interested in the result of f (because we know it is ()), there is no need to use the Kleisli star to lift the argument of the function g so that it picks up the known object

[†] A non-terminating function whose type is IO () is different from a *computation* that returns \bot . The distinction arises because of the lifted type used in the implementation of the monad. Because \bot is different from a tuple containing \bot , then all we are guaranteeing is that the 'result' part of the tuple will never be \bot .

(). Given the Wadler-ized triple $(T, unit, \lambda f \ g \to f \ \tilde{\circ} \ g^*)$, we create the monoid $(T', \tilde{\circ}, id)$, where T' is a perturbation of the type T which is no longer parameterized.

This monoid is very similar to the monoid $(ShowS, \circ, id)$ that is used for printing in Haskell (Hudak et al., 1992). This output monoid is used to ensure a $\mathcal{O}(1)$ when printing data-structures such as trees, and not the quadratic complexity usually associated with using + as the compositional printing operator [see Hudak and Fasel (1992)]. The monoid used here utilizes the opposite behaviour as print injects the string to be printed onto the end of the output stream—the complexity of the helloworlds function is therefore $\mathcal{O}(N^2)$ where N is the number of iterations of the for-loop. Using the monoid shown in Step 1, the definition of helloworlds can be transformed into:

This can be translated into a fold-left computation in which the lambda expression $(\lambda s \, i \to s \, \text{`bind'} \, body Fn \, i)$ is folded down all the values of the loop range. As this lambda expression is now an instance of the lambda expression required by the *fold map fusion law*, then helloworlds shown in loop 3 can be transformed into:

It would seem that vectorization is complete as bind and unit form a monoid. Unfolding the definitions of bind and unit into helloworlds shown in loop 4 reveals the fold-left to be nothing more than an instance of the *leaky fold-left law*, rendering vectorization futile—a further transformation (see Section 9.3) is required!

9.2. AN HISTORICAL ASIDE

A monoid used in step 1 is the essence of the state monad. Monads are typically equated with single-threadedness, and are therefore used as a technique for incorporating imperative features into a purely functional language. Category theory monads have little to do with single-threadedness; it is the sequencing imposed by composition that ensures single-threadedness. In a Wadler-ized monad this is a consequence of bundling the Kleisli star and flipped compose into the bind operator. There is nothing new in this connection. Peter Landin in his Algol 60 paper (Landin, 1965) used functional composition to model semicolon. Semicolon can be thought of as a state transforming operator that threads the state of the machine throughout a program. The work of Wadler (1993) has turned full circle back to Landin's earlier work as their use of Moggi's sequencing monad enables real side-effects to be incorporated into monad operations such as print. This is similar to Landin's implementation of his sharing machine where the assignandhold function can side-effect the store of the sharing machine because of the sequencing imposed by functional composition. Landin defined that "Imperatives are treated as null-list producing

functions" † . The assignandhold imperative is subtly different in that it enables Algol's compound statements to be handled. The function takes a store location and a value as its argument, and performs the assignment to the store of the sharing machine, returning the value assigned as a result of the function. Because Landin assumed applicative order reduction, the **K**-combinator ‡ was used to return (), and the imperative was evaluated as a side effect by the unused argument of the **K**-combinator. He therefore handled statements by wrapping them in a lambda expression that takes () as an argument. Two consecutive Algol-60 assignments would therefore be encoded in the lambda calculus as:

Algol 60	Lambda Calculus
x:= 2;	$((\lambda() \to \mathbf{K}()(assign and hold\mathbf{x}2))\tilde{\circ}$
x := -3;	$(\lambda() \to \mathbf{K}() (assign and hold \mathbf{x}(-3))))()$

By using a lambda with () as its parameter, () can be thought of as the "state of the world" that is threaded throughout a program by functional composition.

9.3. Making the leaky fold-left law work

In a parallel implementation of a loop, what the previous sections have taught us is that the fold-map fusion law and the first duality theorem are crucial in the transformation into a vectorizable fold-left.

Given the monoid (M', bind, unit), and the set of functions $\mathcal{C}(\mathcal{A}, \mathcal{M}(\{()\}))$ (in the running example, this set of functions is the singleton set containing print), then we require that there exists a function g of type $\alpha \to M'$ that models all the functions in the hom-set.

The composition of instances of the function g can be used to create new functions that can be used in the body of the loop—e.g., gv_1 'bind' gv_2 . To successfully vectorize a loop where the body is a computation created from the compositions of the function g, the amalgamated function must be an instance of the lambda expression $\lambda s i \to s \oplus f i$ required by the fold-map fusion law. We define g to be $gv = \lambda s \to s \oplus f v$ where the operator \oplus and f are functions specific to the definition of the hom-set $\mathcal{C}(\mathcal{A}, \mathcal{M}(\{()\}))$, and \oplus is associative. The result of unfolding the definition of bind into the composition of two instances of g produces:

```
\begin{array}{ll} g \, v_1 \text{ `bind' } g \, v_2 \\ \Rightarrow g \, v_2 \, \circ \, g \, v_1 & \text{unfolding bind} \\ \Rightarrow (\lambda s \to s \oplus f \, v_2) \, \circ \\ (\lambda s \to s \oplus f \, v_1) & \text{unfolding } g \\ \Rightarrow \lambda s \to (s \oplus f \, v_1) \oplus f \, v_2 & \text{definition of } \circ \\ \Rightarrow \lambda s \to s \oplus (f \, v_1 \oplus f \, v_2) & \text{associativity of } \oplus. \end{array}
```

As can be seen from the last transformation above, because of the associativity of \oplus , any combination of the compositions of g produces a function that is also an instance of the lambda expression required by the fold-map fusion law. Using the definition of g, the

[†] In Landin's paper, () is the syntactic representation of the empty list and not the unit.

 $^{^{\}ddagger} \mathbf{K} = \lambda \mathbf{x} \mathbf{y} \to \mathbf{x}.$

monoid (M', bind, unit) is converted into $(M'', \oplus, \oplus_{id})$, and the set of monad operations modelled by g is replaced by f.

This transformation is repeatedly applied to the monoid, until M'' is a non-functional type, or the process fails. The relationship between the monad used before and after this transformation is that of a $monoid\ homomorphism$ —i.e., if hom is the monoid homomorphism, then the following holds:

$$hom(unit) \equiv \bigoplus_{Id}$$

 $hom(x \text{'bind'} y) \equiv hom(x) \oplus hom(y)$

The monoid defined by step 1 in the running example only has one operation in the associated hom-set. This function print can be coerced into the form required by g as shown in eqn 9.1.

From the equation, the monoid (String, ++, []) can be used as the new definitions of *unit* and bind, and print becomes the identity function, completing vectorization.

```
type IO = String
unit = []

l 'bind' r = l ++ r
print str = str
```

Finished: the final monoid.

The function helloworlds of loop 4 can be left syntactically unchanged, and unlike the result of §5, because the monoid contains a non-functional type (i.e., *String*), the fold can be truly implemented in parallel.

assoc

10. Conclusions

What this paper has shown is that a small collection of laws can be used to transform imperative for-loops into a form that can be implemented in terms of fold and scan. Instead of a compiler transforming a for-loop using these laws, we have used monads to develop a restricted programming language in which only vectorizable programs can be expressed. All that is required of a compiler is that it performs a minor syntactic translation of a for-loop into the vectorization monad. If the translated program is well-typed, then the loop is vectorizable. The monad therefore guarantees that a loop can be vectorized, and a proof for the monad can be given once, independently of the compiler.

References

Aho, A.V., Sethi, R., Ullman, J.D. (1986). Compilers: Principles, techniques and tools. Addison-Wesley. Bird, R.S. (1986). An introduction of the theory of lists. Technical Report PRG-56, Oxford University Computing Laboratory.

Bird, R.S., Wadler, P. (1988). Introduction to Functional Programming. Prentice-Hall International.
 Blelloch, G.E. (1993). NESL: A nested data-parallel language. Technical Report CS-93-129, Carnegie Mellon University.

- Henderson, P. (1980). Functional Programming: Application and implementation. Prentice-Hall International.
- Hill, J.M.D. (1993a). The aim is laziness in a data-parallel language. In O'Donnell, J.T., Hammond, K., (eds), Glasgow functional programming workshop, Workshops in computing, pp. 83-99. Berlin, Springer-Verlag. Available by FTP from ftp.dcs.qmw.ac.uk in /pub/cpc/jon_hill/aimDpLaziness.ps.
- Hill, J.M.D. (1993b). Vectorizing a non-strict functional language for a data-parallel "Spineless (not so) Tagless G-Machine". In *Proc. 5th international workshop on the implementation of functional languages*, Nijmegen, Holland. Available by FTP.
- Hill, J.M.D. (1994). Data-parallel lazy functional programming. PhD thesis, Department of Computer Science, Queen Mary & Westfield College, University of London.
 Hill, J.M.D., Clarke, K.M. (1993). Parallel Haskell: The vectorization monad. Technical Report 658,
- Hill, J.M.D., Clarke, K.M. (1993). Parallel Haskell: The vectorization monad. Technical Report 658, Department of Computer Science, Queen Mary & Westfield College, University of London. Available by FTP.
- Hudak, P., Fasel, J. (1992). A gentle introduction to Haskell. SIGPLAN Notices, 27(5).
- Hudak, P., Peyton Jones, S.L., Wadler, P. (1992). Report on the Programming Language Haskell, A Non-strict Purely Functional Language (Version 1.2). SIGPLAN Notices, 27(5).
- Landin, P. (1966). The next 700 programming languages. Communications of the ACM, 9(3):157–166.
 Landin, P.J. (1965). A correspondence between ALGOL 60 and Church's lambda notation. Communications of the ACM, 8(2):89–101.
- Launchbury, J. (1993). Lazy imperative programming. In *Proc. ACM Sigplan workshop on State in Programming Languages*, pp. 46–56.
- Mac-Lane, S. (1971). Categories for the working mathematician. Berlin, Springer-Verlag.
- Moggi, E. (1989). Computational lambda-calculus and monads. In *IEEE symposium on Logic in computer science*.
- Nikhil, R.S., Arvind, Hicks, J., Aditya, S., Augustsson, L., Maessen, J., Zhou, Y. (1995). pH Language Reference Manual. MIT Laboratory for Computer Science, 1.0 edition.
- Peyton Jones, S.L. (1987). The Implementation of Functional Programming Languages. Prentice-Hall International.
- Peyton Jones, S.L., Wadler, P. (1993). Imperative functional programming. In ACM Principles of Programming Languages.
- Rice Univ. (1993). High Performance Fortran Language Specification. Houston, Texas. Version 1.0.
- Sparud, J. (1993). Fixing some spaces leaks without a garbage collector. In Functional Programming Languages and Computer Architecture.
- Talpin, J.-P. (1993). Aspects théoriques et praticques de l'inférence de type et d'effets. PhD thesis, L'Ecole nationale supérieure des mines de Paris. (In English).
- Wadler, P. (1990). Comprehending monads. In ACM Conference on Lisp and functional programming, pp. 61–78.