

# Multi-Stage Programming with Explicit Annotations

Walid Taha & Tim Sheard

Oregon Graduate Institute of Science & Technology

{walid,t,heard}@cse.ogi.edu \*

## Abstract

We introduce MetaML, a statically-typed multi-stage programming language extending Nielson and Nielson's two stage notation to an arbitrary number of stages. MetaML extends previous work by introducing *four* distinct staging annotations which generalize those published previously [25, 12, 7, 6]

We give a static semantics in which type checking is done once and for all before the first stage, and a dynamic semantics which introduces a new concept of *cross-stage persistence*, which requires that variables available in any stage are also available in all future stages.

We illustrate that staging is a manual form of binding time analysis. We explain why, even in the presence of automatic binding time analysis, explicit annotations are useful, especially for programs with more than two stages.

A thesis of this paper is that multi-stage languages are useful as programming languages in their own right, and should support features that make it possible for programmers to write staged computations without significantly changing their normal programming style. To illustrate this we provide a simple three stage example, and an extended two-stage example elaborating a number of practical issues.

## 1 Introduction

Multi-stage languages have recently been proposed as intermediate representations for partial evaluation [12, 9, 10] and runtime code generation [7]. These languages generalize the well-known two-level notation of Nielson & Nielson [25] to an arbitrary number of levels.

A major thesis of this paper is that *multi-stage languages are useful not only as intermediate representations, but also as programming languages in their own right*. Multi-stage programming is important whenever performance is important. But there is very little language support for writing multi-stage programs. This paper extends previous work on multi-stage programming with features that are of *practical use to real programmers*.

\*The research reported in this paper was supported by the USAF Air Materiel Command, contract # F19628-93-C-0069, and NSF Grant IRI-9625462

Permission to make digital/hard copy of part or all this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage, the copyright notice, the title of the publication and its date appear, and notice is given that copying is by permission of ACM, Inc. To copy otherwise, to republish, to post on servers, or to redistribute to lists, requires prior specific permission and/or a fee. PEPM '97 Amsterdam, ND

© 1997 ACM 0-89791-917-3/97/0006...\$3.50

We introduce MetaML, a statically-typed multi-stage programming language extending Nielson and Nielson's two-level notation to an arbitrary number of stages (similar to their *B*-level language). MetaML is an extension of a Hindley-Milner polymorphically-typed [22] call-by-value  $\lambda$ -calculus [13] with support for sums, products, recursion, polymorphism, primitive datatypes and static type-inference. It provides the following extensions not found in previous work on multi-stage systems:

- Four distinct staging annotations, which we believe are necessary and sufficient for all multi-stage programming. (Section 5) These annotations generalize and safely combine those published previously [25, 12, 7, 6]. The formal proof of safety is ongoing work.
- A type system ensuring the well-formedness of acceptable multi-stage programs. Type checking is done once and for all before the first stage (Section 10.1).
- Variables of any stage are available in all future stages. This feature, in a language which also contains *run*<sup>1</sup> makes MetaML's annotations strictly more expressive than the languages of Nielsen & Nielsen [25, 24], Davies & Pfenning [7], and Davies [6]. We also deal with the interesting technical problem of ensuring the hygienic binding of free variables (Section 10.2) in code expressions.
- A non-Hindley-Milner, second-order type judgement for the *run* annotation to ensure that no code is ever run in a context in which it is undefined.

As a consequence of the above properties, MetaML provides a programming language suitable for expressing staged computations explicitly. We believe that MetaML can have positive implications for understanding and communicating ideas about multi-stage programs, partial evaluation and the complex process of binding-time analysis in much the same way that the boxed / unboxed(##) distinction provides a language for understanding boxing optimizations as source-to-source transformations [16].

## 2 Why Multi-stage Programs?

The concept of a stage arises naturally in a wide variety of situations. For a compiled language, there are two distinct stages: compile-time, and run-time. But three distinct

<sup>1</sup>An *eval*-like operator.

stages appear in the context of program generation: generation, compilation, and execution. For example, the Yacc parser generator first reads a grammar and generates C code; second, this program is compiled; third, the user runs the object code.

Yet despite the numerous examples of multi-stage software systems, almost all these systems have realized staging in ad-hoc ways. Our goal is to provide a language with well-designed support for multi-staged programming by using explicit staging annotations. In particular, a multi-stage programming language supplies a basis for generation technology. Generators can provide dramatic improvements in the following areas:

**Efficiency.** Specializing a function on a fixed argument can lead to dramatic efficiency gains. Program generators can provide the same efficiency gains that partial evaluation does.

**Productivity and reuse.** When a programming task or activity becomes routine, programmers can use program generators to encapsulate their knowledge of the routine task. This capture of a problem family rather than a single problem increases programmer productivity. Program generators let experts capture their knowledge in a clear (and hence reusable) notation that can then be used for synthesizing the desired software component [21, 17, 18].

**Reliability and quality.** The greatest source of errors in code maintenance is human intervention. When less human intervention is needed to modify a software product, there are proportionately fewer opportunities for error insertion and less rework of code is necessary. Automatically generated components require little manual rework after a re-generation.

Our language, MetaML, was designed as basis for an integrated generator system. It provides an approach radically different from, and superior to, the classic “programs-as-strings” view that seems to predominate in many ad-hoc software generation systems. MetaML is tightly integrated in this sense.

### 3 Relationship to Partial Evaluation

Today, the most sophisticated automatic staging techniques are found in partial evaluation systems [15]. Partial evaluation optimizes a program using a priori information about some of that program’s inputs. The goal is to identify and perform as many computations as possible in a program before run-time.

*Offline* partial evaluation has two distinct steps, *binding-time analysis* (BTA) and *specialization*. BTA determines which computations can be performed in an earlier stage given the names of inputs available before run-time (static inputs).

In essence, BTA performs automatic staging of the input program. After BTA, the actual values of the inputs are made available to the specializer. Following the annotations, the specializer either performs a computation, or produces text for inclusion in the output (*residual*) program.

The relationship between partial-evaluation and multi-stage programming is that the intermediate data structure between the two steps is a *two-stage annotated program* [2], and that the specialization phase is (the first stage in) the

execution of the two-stage annotated program produced by BTA. Recently, Glück and Jørgensen proposed *multi-level BTA* and showed that it is an efficient alternative to multiple specialization [9, 10]. Their underlying annotated language is closely related to MetaML.

### 4 Why Explicit Annotations?

If BTA performs staging automatically, why should programmers stage programs manually? They shouldn’t *have* to, but there are several important reasons why they may *want* to:

**Pragmatic.** While there are advantages to discussing the semantics of annotated programs and the techniques of BTA at the same time, we feel that the complexity of the semantics of annotated programs warrants studying them in (relative) isolation of other partial evaluation issues.

**Pedagogical tool.** It has been observed that it is sometimes hard for users to understand the workings of partial evaluation systems [14]. New users often lack a good mental model of how partial evaluation systems work. Although BTA is an involved process, requiring special expertise, the annotations it produces are relatively simple and easy to understand. However, new users are often uncertain: What is the output of a binding-time analysis? What are the annotations? How are they expressed? What do they really mean? The answers to these questions are crucial to the effective use of partial evaluation. Our observation is that programmers can understand the annotated output of BTA, without actually knowing how BTA works. Having a programming language with explicit staging annotations would help users of partial evaluation understand more of the issues involved in staged computation, and, hopefully, reduce the steep learning curve currently associated with learning to use a partial evaluator effectively [15]. Nielson & Nielson’s two-stage notation is the only widely accepted notation for expressing staged computation. But Nielson & Nielson’s notation is not widely viewed as a programming language, perhaps because over-bars and under-bars do not appear on the standard keyboard and no implementation of it is in widespread use.

**Controlling Evaluation Order.** Whenever performance is an issue, control of evaluation order is essential. BTA optimizes the evaluation order, but sometimes it is just easier to say what you want than to force a BTA to discover it. Automatic analyses like BTA are necessarily incomplete, and can only approximate the knowledge of the programmer. By using explicit annotations the programmer can exploit his full knowledge of the program domain.

In addition, BTA for programs with more than two stages is still imprecise. Hand annotation may be the only feasible mechanism for staging multi-stage programs, and may be the only mechanism expressive enough for the degree of control needed in many circumstances.

**High-Level Program Generation.** As we will also illustrate in this paper, staging annotations also provide a powerful tool for high-level program generation. No explicit construction of parse trees is needed. As a consequence, generators can be simpler and more reliable than their hand constructed counterparts. It is also easier to verify the correctness of both the generators and the programs they generate, as the issues of representation are hidden away from

the programmer.

## 5 MetaML's Multi-Stage Programming Annotations

The two-level notation of Nielson & Nielson [25] features two annotations: *over-bars* to mark computations of the first stage, and *under-bars* to mark those of the second stage. Although quite powerful, this is only a subset of the annotations needed for pragmatic multi-stage programming. MetaML has four programming constructs:

- Meta-Brackets ( $\langle \_ \rangle$ ) are the primary means for delaying a computation. For example, whereas the expression  $40+2$  specifies a current (or first) stage computation,  $\langle 40+2 \rangle$  specifies one for the next (or second) stage. A binary type constructor  $\langle \_ , \_ \rangle$  is used to distinguish the type of the latter expression from the first one. For example,  $7$  has type  $\text{int}$ , but  $\langle 7 \rangle$  has type  $\langle \text{int}, 'a \rangle$ , where (as is in ML)  $'a$  is a free type variable. The expression  $\langle\langle 1, \langle 2+1 \rangle \rangle\rangle$  has type:  $\langle\langle \text{int} * \langle \text{int}, 'c \rangle \rangle, 'b \rangle, 'a \rangle$  and the addition will be performed in the fourth stage. The second type in the code type constructor represents the name of the context in which this code can execute. In the examples above the context is completely unconstrained hence the type variables. More about this in section 10.
- Escape ( $\sim$ ) can occur only inside enclosing meta-brackets. It is the mechanism used to insert smaller delayed computations into larger ones. Escape allows its argument to escape the “freeze” imposed by a surrounding meta-bracket and to “splice” its result into the delayed computation being built. For example:

```
let val a =  $\langle 1+4 \rangle$  in  $\langle 72+\sim a \rangle$  end
```

returns the expression  $\langle 72+(1+4) \rangle$ . The escaped computation must yield a piece of code with a type that can be inserted in the context where the escape appears. The type system ensures that this is the case. For example, if  $x$  has type  $\langle \text{int}, 'a \rangle$ , then  $\langle (x, 1) \rangle$  has type  $\langle\langle \text{int}, 'a \rangle * \text{int} \rangle, 'b \rangle$  and  $\langle (\sim x, 1) \rangle$  has type  $\langle (\text{int} * \text{int}), 'a \rangle$ . Objects of type code are first class citizens, and can even be  $\lambda$ -abstracted. For example:

```
val add.72.later = fn a =>  $\langle 72+\sim a \rangle$ 
```

declares a first class function with type  $\langle \text{int}, 'a \rangle \rightarrow \langle \text{int}, 'a \rangle$ , and the expression  $\text{add.72.later } \langle 8 \rangle$  returns  $\langle 72+8 \rangle$ . From the language designer's point of view, escape poses a very interesting technical problem, as not all uses of escape are reasonable. We discuss this issue in Section 10.1.

- Run ( $\text{run } \_$ ) takes a code-valued argument and runs it. It is the only way a computation “frozen” using meta-brackets can be computed (or “forced”) in the current stage. The argument to run must be of *code* type. Having run in the language implies introducing a kind of reflection [30], and allows a future-stage computation to be performed now. To illustrate, consider the expression:

```
let val a =  $\langle 50-10 \rangle$  in  $2+(\text{run } a)$  end
```

This expression has type  $\text{int}$  and returns the value 42 when computed. Although  $\text{run}$  is not an annotation used in the result of BTA, it is an essential feature for a programmer who wants to use multi-stage programming to control evaluation order.

- Lift ( $\text{lift } \_$ ) allows the user to convert any *ground* value (not containing a function) into code. Contrast this with meta-brackets which converts any *syntactic expression* into a piece of code. Lift evaluates its arguments, and meta-brackets do not. Lift is most often used in conjunction with escape, because only pieces of code can be “spliced-in”. For example, in the expression  $\langle 1+\sim(\text{lift } 2+3) \rangle$ , the escape forces lift  $(2+3)$  to be computed in the first stage. The addition evaluates to the value 5, and lift converts this result into the piece of code  $\langle 5 \rangle$ , which is spliced (because of the escape) back into the original expression to return  $\langle 1+5 \rangle$ . Lift can be used on structured values such as tuples and lists as long as they do not contain functions. For example  $\text{lift } [(2,3), (2*1,4)]$  evaluates to  $\langle [(2,3), (2,4)] \rangle$ . Function values cannot be lifted using lift, as we cannot derive an intensional representation for them in general (This does not mean that function values cannot be delayed using meta-brackets. See Section 7.)

**Precedence Issues.** The escape operator ( $\sim$ ) has the highest precedence; even higher than function application. This allows us to write:  $\langle f \sim x y \rangle$  rather than  $\langle f (\sim x) y \rangle$ . The lift ( $\text{lift } \_$ ) and run ( $\text{run } \_$ ) operators have the lowest precedence. The scope of these operators extends to the right as far as possible. This makes it possible to write  $\langle f \sim(\text{lift } g y) \rangle z$  rather than  $\langle f \sim(\text{lift } (g y)) z \rangle$ .

## 6 Hand-Staging: A Short Example

Using MetaML, the programmer can stage programs by inserting the proper annotations at the right places in the program. The programmer uses these annotations to modify the default (strict) evaluation order of the program.

In our experience, starting with the type of the function to be hand-staged makes the number of different ways in which it can be annotated quite tractable. This leads us to believe that the location of the annotations in a staged version of a program is significantly constrained by its type. For example, consider the function `member` defined as:

```
(* member : int -> int list -> bool *)
fun member v l =
  if (null l)
  then false
  else if v=(hd l)
  then true
  else member v (tl l);
```

The function `member` has type  $\text{int} \rightarrow \text{List int} \rightarrow \text{bool}$ .<sup>2</sup> A good strategy for hand annotating a program is to first determine the target type of the desired annotated program. In the `member` example, the list parameter `l` is available in the first stage, and the element searched for will be available later. So, one target type for the hand-staged function is  $\langle \text{int}, 'a \rangle \rightarrow \text{List int} \rightarrow \langle \text{bool}, 'a \rangle$ .

<sup>2</sup>Function “=” has type  $(\text{int} * \text{int}) \rightarrow \text{bool}$  which forces `v` and `l` to have types  $\text{int}$  and  $\text{List int}$ , respectively.

Now we can begin annotating, starting with the whole expression, and working inwards until all sub-expressions are covered. At each step, we consider what annotations will “fix” the type of the expression so that the whole function has a type closer to the target type.

The following function realizes this type:

```
(* member : <int,'a> -> int list -> <bool,'a> *)
fun member v l =
  if (null l)
  then <false>
  else <if ~v~(lift hd l)
        then true
        else ~(member v (tl l))>;
```

The annotation `~(lift hd l)` is used rather than `hd l` in order to ensure that `hd` is performed during the first stage. Otherwise, all selections of the head element of the list would have been delayed until the code constructed was run in a later stage.

The meta-brackets around the branches of the outermost if-expression ensure that the return value of `member` will be a code type (`<?,?>`). The first branch (`false`) needs no further annotations, and makes the return value precisely a `<bool,'a>`. Moving inwards in the else branch, the condition of the inner if-expression (in particular `~v`) forces the type of the `v` parameter to have type `<int,'a>` as planned.

Just like the first branch of the outer if-statement, the whole of the inner if-statement must return `bool`. So, the first branch (`true`) is fine. But because the recursive call to `member` has type `<bool,'a>`, it must be escaped. Inserting this escape also implies that the recursion will be performed in the first stage, which is exactly the desired behavior. Thus, the result of the staged `member` function is a recursively-constructed piece of code with type `bool`.

Evaluating `<fn x => ~(member <x> [1,2,3])>` yields:

```
<fn d1 =>
  if d1 %= 1
  then true
  else if d1 %= 2
        then true
        else if d1 %= 3
              then true
              else false>
```

The percentage sign (%) at the beginning of an identifier indicates that it was bound to a value in the environment in which the code was constructed. Its precise meaning will be explained in Sections 7 and 10.2.

## 7 The Design of MetaML

MetaML was designed as a statically-typed programming language, and not as an internal representation for a multi-stage system. Our primary goals were: first, it should be easy to write multi-staged programs, second it should be as flexible as possible, and third it should ensure that only “reasonable things” can be done using the annotations. Therefore, our design choices were different from those of other multi-stage systems. In particular, we consider the following quality crucial to MetaML:

**Cross-stage Persistence:** *A variable  $i$  bound in stage  $n$ , will be available in stages  $n$ ,  $n + 1$  and all future stages.*

To the user, this means the ability to stage non-closed expressions. Non-closed expressions, like  $\lambda$ -abstractions with free variables, must resolve their free variable occurrences in the static environment where the meta-bracketed expression occurs. One can think of a code value as containing an environment which binds its free variables. For example the expression,

```
let val a=1+4 in <72+a> end
```

returns a value `<72+%a>`. The % sign indicates that the free variable `a` is bound in the value’s local environment. The % sign is printed by the display mechanism. The variable `a` has been bound during the first stage to the constant 5. In fact, in MetaML `%a` is not a variable, but rather, a new *constant*, and the name “a” is only hint to the user about where this constant originated. When `%a` is evaluated in a later stage, it will return 5 independent of the binding for the variable `a` in the new context since it is bound in the value’s local environment. Arbitrary values (including functions) can be delayed using this hygienic binding mechanism.

Specifying this behavior turns out to be non-trivial. In an interpreter for a multi-stage language, this requirement manifests itself as complex variable-binding rules, the use of closures, or capture-free substitutions. Our semantics addresses this in a rather unique way (See Section 10.2).

Cross-Stage Persistence poses a problem when staging is used for program generation. If the first stage is performed on one computer, and the second on another, we must “port” the local environments from the first machine to the second. Since arbitrary objects, such as functions and closures, can be bound in these local environments this can become a problem. Currently, MetaML assumes that the computing environment does not change between stages. This is part of what we mean by having an integrated system.

Cross-Stage Persistence can be relaxed by allowing variables to be available at exactly one stage. This seems to have been the case in all multi-stage languages known to us to date [25, 12, 9, 10, 7, 6]. The primary difficulty in implementing persistence is the proper hygienic treatment of free variables. We will show how this problem can be solved, thus allowing the user to stage significantly more expressions than was previously possible.

But even in MetaML, it will not be possible to stage every expression in the language. In particular, we must ensure that the user can only specify computations that respect the following condition:

**Cross-Stage Safety:** *An input first available at stage  $m$  cannot be used at a stage  $n$  if  $m > n$ .*

The problem arises with the use of the escape annotation. In particular, consider the expression

```
fn a => <fn b => ~(a+b)>
```

which is an (incorrectly) staged version of the function  $\lambda a. \lambda b. a + b$ . Operationally, the annotations require computing `a+b` in the first stage, while the value of `b` will be available only in the second stage! Therefore, MetaML’s type system was designed to ensure that “well-typed programs won’t go wrong”, where going wrong now includes the violation of

the cross-stage safety condition, as well as the standard notions of “going wrong” [22] in statically-typed languages. A formal proof of this property is ongoing work.

In our experience with the language, having a type system to screen-out programs containing this kind of error is a significant aid in hand-staging programs.

## 8 Isomorphism for Code Types

Recall the types of the staged `member` function: `<int, 'a> -> List int -> <bool, 'a>`, and the type of the term `<fn x => (member <x> [1,2,3])>` which is: `<int -> bool, 'a>`

This suggests that a function from code to code can be turned into the code of a function. This is important to users because `<alpha, 'a> -> <beta, 'a>` is a function and cannot be printed or observed, while `<alpha -> beta, 'a>` is a representation of a function, and can be printed and observed. We can define two functions to convert between these two types:

```
(* back: <'A, 'c> -> <'B, 'c> -> <('A -> 'B), 'c> *)
fun back f = <fn x => ~(f <x>>);

(* forth: <('A -> 'B), 'c> -> (<'A, 'c> -> <'B, 'c>) *)
fun forth f x = <f ~x>;
```

Here we use capitalized type variables to distinguish the type in the code from the context the code must evaluate in.

The conversion is not between syntactic forms, but semantic values. For example, the code produced by an application of `back` is in a language extended with a new construct that allows us to embed any value into syntax, without needing to know about its intentional representation. Thus, we are really not converting functions into source code, but rather, returning syntax that denotes this function under our semantics. Under this proviso (and disregarding termination issues) the composition of these two functions is identity under MetaML’s semantics (see Section 10). They define an isomorphism between values of type `<A, 'c> -> <B, 'c>` and `<A -> B, 'c>`. [3].

The `back` and `forth` functions are similar to two-level  $\eta$ -expansion [5]. In MetaML, however, `back` and `forth` are not only meta-level concepts or optimizations, but rather, first class functions in the language, and the user can apply them directly to values of the appropriate type.

This isomorphism can also be viewed as a formalization of the intuitive equivalence of a symbolic evaluator [23] `<A, 'c> -> <B, 'c>` and the syntactic representation of a function `<A -> B, 'c>`. It seems that this isomorphism, which MetaML has allowed us to make concrete, is at the heart of concise reduction systems, such as Danvy’s type-directed partial evaluator [4] and its extensions [27]. Under MetaML’s semantics, we can switch between the two types without needing to worry about substitution or variable capture.

This has profound implications for the writing of staged functions. In our experience annotating a function to have type `<A, 'c> -> <B, 'c>` requires less annotations than annotating it to have type `<A -> B, 'c>` and is often easier to think about. Because we are more used to reasoning about functions, this leads us to avoid creating functions of the latter kind except when we need to inspect the code.

The type of `back` is one of the axioms of the logic system motivating the type system of Davies [6]. MetaML’s type system was motivated purely by operational reasons. At

the same time, it is important for the programmer to have both coercions, thereby being able to switch back and forth between the two isomorphic types as the need arises.

This becomes even more important when writing programs with more than two stages. Consider the function:

```
fun back2 f = <fn x => <fn y => ~(f <x> <<y>>>>);
back2 : (<A, 'd> -> <<B, 'e>, 'd> -> <<C, 'e>, 'd>)
       -> <A -> <B -> C, 'e>, 'd>
```

This allows us to write a program which takes a `<a, 'd>` and a `<<b, 'e>, 'd>` as arguments and which produces a `<<c, 'e>, 'd>` and stage it into a three-stage function. Our experience is that such functions have considerably fewer annotations, and are easier to think about. We illustrate this in the next section.

## 9 A Multi-Stage Example

When information arrives in multiple phases it is possible to take advantage of this fact to get better performance. Consider a generic function for computing the inner product of two vectors. In the first stage the arrival of the size of the vectors offers an opportunity to specialize the inner product function on that size, removing the overhead of looping over the body of the computation  $n$  times. The arrival of the first vector affords a second opportunity for specialization. If the inner product of that vector is to be taken many times with other vectors it can be specialized by removing the overhead of looking up the elements of the first vector each time. This is exactly the case when computing the multiplication of 2 matrixes. For each row in the first matrix, the dot product of that row will be taken for each column of the second. This example has appeared in several other works [9, 20] and we give our version below. We give three versions of the inner product function. One (`iproduct`) with no staging annotations, the second (`iproduct2`) with two levels of annotations, and the third (`iproduct3`) with two levels of annotations but constructed with the `back2` function. In MetaML we quote relational operators involving `<` and `>` because of the possible confusion with meta-brackets.

```
(* iproduct : int -> Vector -> Vector -> int *)
fun iproduct n v w =
  if n == 0
  then ((nth v n) * (nth w n)) + (iproduct (n-1) v w)
  else 0;

(* iproduct2 : int -> <Vector -> <Vector -> int> *)
fun iproduct2 n = <fn v => <fn w =>
  ~(if n == 0
    then << (^(lift nth v n) * (nth w n)) +
      (^(~(iproduct2 (n-1)) v) w)
    >>
    else <<0>>) >>;

(* p3 : int -> <Vector> -> <<Vector>> -> <<int>> *)
fun p3 n v w =
  if n == 0
  then << (^(lift nth ~v n) * (nth ~w n)) +
    ~(p3 (n-1) v w) >>
  else <<0>>;

fun iproduct3 n = back2 (p3 n);
```

Notice that the staged versions are remarkably similar to the unstaged version, and that the version written with `back2` has fewer annotations. The type inference mechanism was a great help in placing the annotations correctly.

An important feature of MetaML is the visualization help that the system affords. By “testing” `iprod2` on some inputs we can “see” what the results are immediately.

```
val f1 = iprod3 3;
f1 : <Vector -> <Vector -> int>> =
<fn d1 =>
  <fn d5 =>
    (~(lift %nth d1 3) * (%nth d5 3)) +
    (~(lift %nth d1 2) * (%nth d5 2)) +
    (~(lift %nth d1 1) * (%nth d5 1)) +
    0 >>
```

When this piece of code is run it will return a function, which when applied to a vector builds another piece of code. This building process includes looking up each element in the first vector and splicing in the actual value using the `lift` operator. Using `lift` is especially valuable if we wish to inspect the result of the next phase. To do that we evaluate the code by running it, and apply the result to a vector.

```
val f2 = (run f1) [1,0,4];
f2: <Vector -> int> =
<fn d1 => (4 * (%nth d1 3)) +
          (0 * (%nth d1 2)) +
          (1 * (%nth d1 1)) + 0 >
```

Note how the actual values of the first array appear in the code, and how the access function `nth` appears as a constant expression applied to the second vector `d1`.

While this code is good, it does not take full advantage of all the information known in the second stage. In particular, note that we generate code for the third stage which may contain multiplication by 0 or 1. These multiplications can be optimized. To do this we write a second stage function `add` which given an index into a vector `i`, an actual value from the first vector `x`, and a piece of code which names the second vector `y`, constructs a piece of code which adds the result of the `x` and `y` multiplication to the code valued fourth argument `e`. When `x` is 0 or 1 special cases are possible.

```
(* add : int -> int -> <Vector> -> <int> *)
fun add i x y e =
  if x=0
  then e
  else if x=1
  then <(nth y ~(lift i)) + ~e>
  else <~(lift x) * (nth y ~(lift i))) + ~e>;
```

This specialized function is now used to build the second stage computation:

```
(* p3 : int -> <Vector> -> <<Vector>> -> <<int>> *)
fun p3 n v w =
  if n = 1
  then << ~ (add n (nth ~v n) ~w <0>) >>
  else << ~ (add n (nth ~v n) ~w
    < ~ (p3 (n-1) v w) >) >>;

fun iprod3 n = back2 (p3 n);
```

Now let us observe the result of the first stage computation.

```
val f3 = iprod3 3;
f3: <Vector -> <Vector -> int>> =
<fn d1 =>
  <fn d5 =>
    ~(%add 3 (%nth d1 3) <d5>
      < ~(%add 2 (%nth d1 2) <d5>
        < ~(%add 1 (%nth d1 1) <d5>
          <0>>>>) >>
```

This code is linear in the size of the vector; if we had actually inlined the calls to `add` it would be exponential. This is why being able to have free variables (such as `add`) in code is indispensable. Now let us observe the result of the second stage computation:

```
val f4 = (eval f3) [1,0,4];
f4: <Vector -> int> =
<fn d1 => (4 * (%nth d1 3)) + (%nth d1 1) + 0>
```

Note that now only the multiplications that contribute to the answer are evident in the third stage program. If the vector is sparse then this sort of optimization can have dramatic effects.

## 10 Semantics of $\lambda^M$

Figure 1 presents the static and dynamic semantics of the meta-lambda calculus,  $\lambda^M$ . This calculus is a mini-MetaML, which illustrates the relevant features of the staging annotations on the semantics of MetaML.

$\lambda^M$  is a call by value lambda calculus which supports integers, functions, and code ( $\text{int} \mid t \rightarrow t \mid \langle t \rangle^t \mid \alpha$ ). The syntax of terms includes integer constants, variables, applications, abstractions ( $i \mid x \mid ee \mid \lambda x^t.e$ ) and the four staging annotations: meta brackets, escape and run ( $\langle e \rangle \mid \sim e \mid \text{run } e$ ). In addition, the constant operator ( $\uparrow v$ ) allows us to injects a value into a term, and is crucial to the conciseness of our implementation of Cross-Stage Persistence. It is these constants that we print out as a % followed by a name. Note that users do not write programs with the constant operator; it is only introduced during reduction. Every shift in stage from a lower stage to a higher stage enriches the syntax passed to the higher stage with a new set of constants; the values of the previous stage that could still be referenced in the future.

Finally, we note that the `lift` annotation of MetaML is does not appear explicitly in the semantics of  $\lambda^M$ . Instead, the `lift` can be defined within  $\lambda^M$  for each type.

### 10.1 Static Semantics

The static semantics is expressed as a set of inference rules that determine if a term is well-formed, and determine its type. The judgement  $\Sigma \Delta \vdash^n x : \tau_1, \tau_2$  is read *under the context stack  $\Sigma$ , the type environment  $\Delta$ , the term  $x$  has type  $\tau_1$  at level  $n$  and may execute in the context with name  $\tau_2$ .*

The intuition behind contexts, is that any expression can only execute in a context which contains bindings for its free variables. The type inference algorithm assigns the same context name to expressions that must execute in the same context.

The type assignment  $\Delta$  maps variable to types and levels and context names. Every variable is bound at some particular level, namely, the level of the abstraction in which

it is bound (**Abs** rule). The role of  $n$  in the judgement  $\Sigma \Delta \vdash^n x : \sigma$  is to keep track of the level of the expression being typed. The *level* of a subexpression is the number of uncanceled surrounding brackets. One surrounding escape cancels one surrounding bracket. Hence,  $n$  is incremented for an expression inside meta-brackets (**Bracket**), and decremented for one inside an escape (**Escape**). Note that the rule **Escape** does not allow escape to appear at level 0. In other words, escape must appear inside uncanceled meta-brackets.

There are three main kinds of errors related to staging annotations that can occur at runtime:

- A variable is used in a stage before it is available, or
- Run or escape are passed values having a non-code type, or
- Run is passed a code-type value with free variables. This manifests itself in the type, where the name of the context is constrained.

The first kind of error is checked by the **Var** rules. Having no rule for  $m > n$  enforces Cross-Stage Persistence: Variables available in the current stage ( $m$ ) can be used in all future stages ( $n$ ). The second kind of error is checked by the **Run**  $n$  and **Esc**  $n+1$  rules. Detecting the third kind of error is an important contribution of this paper, and is accomplished by the free variable check in the rule **Run**  $n$ : *Only code whose context is completely unconstrained may be run.*

For the standard part of the language, *code* (now denoted by  $\langle \cdot \rangle$  for conciseness) is a normal type constructor that needs no special treatment and the level  $n$  is never changed. Similar type systems have been identified and used by Gomard and Jones [11], Davies & Pfenning [6] and Davies [7].

An important difference between these type systems and the one in Figure 1 is that in all previous statically-typed multi-stage languages [25, 7, 6], only the following *monolithic* type rule is used for variables:

$$\text{Var (Monolithic): } \frac{(\Delta x) = \tau^m}{\Delta \vdash^n x : \tau} \quad \text{when } m = n$$

Whereas we allow the more general condition  $m \leq n$ . This means any generated expressions may as well be evaluated in the empty environment since all well-typed terms are closed terms and cannot reference any free variables. For example the expression:

```
val lift_like = fn x => <x>
```

is accepted, because inside the meta-brackets,  $n = 1$ , and  $(\Delta x) = \alpha^0$ . This expression is not accepted by the monolithic variable rule. Note that while the whole function has type  $\alpha \rightarrow \langle \alpha \rangle$  it does not provide us with the functionality of lift, because the result of applying `lift_like` to any value always returns `<%x>`, and not a literal expression denoting the value. But this example demonstrates that meta-brackets can be used to “lift” any value, including functions. This is explained in the dynamic semantics.

The type system rejects the expression

```
fn a => <fn b => ~(a+b)>
```

because, inside the escape,  $n = 0$ , and  $(\Gamma b) = \alpha^1$ , but  $1 > 0$ .

## 10.2 Dynamic Semantics

The dynamic semantics provides meaning to terms. Values are a subset of terms, and we denote them with a small diamond superscript  $(i^\diamond \mid \{\lambda x^t . e\}^\diamond \mid \langle e \rangle^\diamond)$ . The semantics given in Figure 1, when applied to well-typed terms, maintains the invariant that no free variables ever occur in code values which will later be run.

The dynamic semantics is broken into two sets of rules, reduction and rebuilding. Reduction  $(\Gamma \vdash e \hookrightarrow v)$  maps terms to values and rebuilding  $(\Gamma \vdash e \xrightarrow{n+1} v)$  maps terms to terms and is indexed by a level  $n+1$ . Rebuilding “reconstructs” terms under the environment  $\Gamma$ .

The environment  $\Gamma$  binds a variable to a value. Bindings in environments come in two flavors: real ( $\text{Real}(v)$ ) and symbolic ( $\text{Sym}(x)^t$ ). The extension of the environment with real values occurs only in the rule **App** 0. Such values are returned under reduction (**Var** 0), or injected into constant terms (**RVar**  $n+1$ ) under rebuilding.

Several things about rebuilding should be noted:

1. Rebuilding replaces all known variables with a constant expression  $(\uparrow v)$  where the  $v$  comes from  $\text{Real}(v)$  bindings in  $\Gamma$  (**RVar**  $n+1$ ).
2. Rebuilding renames all bound variables. Symbolic  $\text{Sym}(x')^t$  bindings occur in rules **Abs** 0 and **Abs**  $n+1$  where a term is rebuilt, and new names must be introduced to avoid potential variable capture. These new names are projected from the environment in rule **SVar**  $n+1$ .
3. Rebuilding executes escaped expressions to obtain code to “splice” into the context where the escaped term occurs (**Escape** 1).

Without the staging annotations, rebuilding is simply capture-free substitution of the symbolic variables bound in  $\Gamma$ . Rebuilding is initiated in two places, in rule **Abs** 0 where it is used for capture-free substitution, and in rule **Bracket** 0 where it is applied to terms inside dynamic brackets and it describes how the delayed computations inside a dynamic value are constructed.

The type system ensures that in rule **Abs** 0, there are no embedded escapes at level 1 that will be encountered by the rebuilding process, so rebuilding actually implements capture-free substitution as advertised.

The rules **Escape** 1, **Run** 0, and **Bracket** 0 are at the heart of the dynamic semantics.

In the rebuilding rule **Escape** 1, an escaped expression at level 1 indicates a computation must produce a code valued result  $(\langle e_2 \rangle^\diamond)$ , and rebuilding returns the term  $e_2$ .

The reduction rule **Bracket** 0 describes how a code value is constructed from a meta-bracketed term  $\langle e_1 \rangle$ . The embedded expression is rebuilt at level 1, and the returned term is injected into the domain of values.

The reduction rule **Run** 0 describes how a code valued term is executed. The term is reduced to a code valued term, and the embedded term is then reduced in the empty environment to produce the answer. The empty environment is sufficient because all free variables in the original code valued term have been replaced by constant expressions  $(\uparrow v)$ .

## Domains and Relations

levels	$n \rightarrow 0 \mid 1 \mid n+1 \mid n+2 \mid \dots$			
integers	$i \rightarrow \dots \mid -2 \mid -1 \mid 0 \mid 1 \mid 2 \mid \dots$			
types	$\tau \rightarrow \text{int} \mid \tau \rightarrow \tau \mid \langle \tau \rangle^\tau \mid \alpha$			
terms	$e \rightarrow i \mid x \mid ee \mid \lambda x^\tau.e \mid \langle e \rangle \mid \sim e \mid \text{run } e \mid \uparrow v$			
values	$v \rightarrow i^\circ \mid \{\lambda x^t.e\}^\circ \mid \langle e \rangle^\circ$			
bindings	$b \rightarrow \text{Real}(v) \mid \text{Sym}(x)^t$			
environments	$\Gamma \rightarrow \bullet \mid \Gamma, x \mapsto b$	<b>where</b> $(\Gamma, x \mapsto b)y \equiv \text{if } x = y \text{ then } b \text{ else } \Gamma y$		
type environments	$\Delta \rightarrow \circ \mid \Delta, x \mapsto (\tau, \alpha)^n$	<b>where</b> $(\Delta, x \mapsto (\tau, \alpha)^n)y \equiv \text{if } x = y \text{ then } (\tau, \alpha)^n \text{ else } \Delta y$		
context stacks	$\Sigma \rightarrow [] \mid \alpha; \Sigma$			
reduction	$\Gamma \vdash e \hookrightarrow v$	rebuilding at level $n$	$\Gamma \vdash e \xrightarrow{n} e$	term typing at level $n$ $\Sigma \Delta \vdash^n e : \tau, \tau$

## Static Semantics

<b>Int <math>n</math>:</b> $\Sigma \Delta \vdash^n i : \text{int}, \alpha$	<b>Var =0:</b> $\frac{\Delta x = (\tau, \alpha_1)^0}{\Sigma \Delta \vdash^n x : \tau, \alpha_2}$	<b>Var <math>\leq n</math>:</b> $\frac{\Delta x = (\tau, \alpha)^i \quad i \neq 0 \wedge i \leq n}{\Sigma \Delta \vdash^n x : \tau, \alpha}$
<b>Br <math>n</math>:</b> $\frac{(\alpha_2; \Sigma) \Delta \vdash^{n+1} e : \tau, \alpha_1}{\Sigma \Delta \vdash^n \langle e \rangle : \langle \tau \rangle^{\alpha_1}, \alpha_2}$	<b>Abs <math>n</math>:</b> $\frac{\Sigma \Delta, x \mapsto (\tau_1, \alpha)^n \vdash^n e : \tau_2, \alpha}{\Sigma \Delta \vdash^n \lambda x^{\tau_1}.e : \tau_1 \rightarrow \tau_2, \alpha}$	<b>Esc <math>n+1</math>:</b> $\frac{\Sigma \Delta \vdash^n e : \langle \tau \rangle^{\alpha_1}, \alpha_2}{(\alpha_2; \Sigma) \Delta \vdash^{n+1} \sim e : \tau, \alpha_1}$
<b>Run <math>n</math>:</b> $\frac{\alpha_2 \notin FV(\Delta, \tau) \quad \Sigma \Delta \vdash^n e : \langle \tau \rangle^{\alpha_2}, \alpha_1}{\Sigma \Delta \vdash^n \text{run } e : \tau, \alpha_1}$	<b>App <math>n</math>:</b> $\frac{\Sigma \Delta \vdash^n e_1 : \tau_1 \rightarrow \tau, \alpha \quad \Sigma \Delta \vdash^n e_2 : \tau_1, \alpha}{\Sigma \Delta \vdash^n e_1 e_2 : \tau, \alpha}$	<b>Con <math>n</math>:</b> $\frac{\text{*Never appears in source terms*}}{\Sigma \Delta \vdash^n \uparrow v : \star}$

## Dynamic Semantics

<b>Int 0:</b> $\Gamma \vdash i \hookrightarrow i^\circ$	<b>Int <math>n+1</math>:</b> $\Gamma \vdash i \xrightarrow{n+1} i$
<b>Abs 0:</b> $\frac{\Gamma, x \mapsto \text{Sym}(x')^\tau \vdash e \xrightarrow{1} e_1}{\Gamma \vdash \lambda x^\tau.e \hookrightarrow \{\lambda x'^\tau.e_1\}^\circ}$	<b>Abs <math>n+1</math>:</b> $\frac{\Gamma, x \mapsto \text{Sym}(x')^\tau \vdash e_1 \xrightarrow{n+1} e_2}{\Gamma \vdash \lambda x^\tau.e_1 \xrightarrow{n+1} \lambda x'^\tau.e_2}$
<b>App 0:</b> $\frac{\Gamma \vdash e_1 \hookrightarrow \{\lambda x^\tau.e\}^\circ \quad \Gamma \vdash e_2 \hookrightarrow v_2 \quad \bullet, x \mapsto \text{Real}(v_2) \vdash e \hookrightarrow v}{\bullet, x \mapsto \text{Real}(v_2) \vdash e_1 e_2 \hookrightarrow v}$	<b>App <math>n+1</math>:</b> $\frac{\Gamma \vdash e_1 \xrightarrow{n+1} e_3 \quad \Gamma \vdash e_2 \xrightarrow{n+1} e_4}{\Gamma \vdash e_1 e_2 \xrightarrow{n+1} e_3 e_4}$
<b>Var 0:</b> $\frac{\Gamma x = \text{Real}(v)}{\Gamma \vdash x \hookrightarrow v}$	<b>SVar <math>n+1</math>:</b> $\frac{\Gamma x = \text{Sym}(x')^\tau}{\Gamma \vdash x \xrightarrow{n+1} x'}$ <b>RVar <math>n+1</math>:</b> $\frac{\Gamma x = \text{Real}(v)}{\Gamma \vdash x \xrightarrow{n+1} \uparrow v}$
	<b>EVar <math>n+1</math>:</b> $\frac{x \notin \Gamma}{\Gamma \vdash x \xrightarrow{n+1} x}$
<b>Bracket 0:</b> $\frac{\Gamma \vdash e_1 \xrightarrow{1} e_2}{\Gamma \vdash \langle e_1 \rangle \hookrightarrow \langle e_2 \rangle^\circ}$	<b>Bracket <math>n+1</math>:</b> $\frac{\Gamma \vdash e_1 \xrightarrow{n+1} e_2}{\Gamma \vdash \langle e_1 \rangle \xrightarrow{n} \langle e_2 \rangle}$
<b>Escape 1:</b> $\frac{\Gamma \vdash e_1 \hookrightarrow \langle e_2 \rangle^\circ}{\Gamma \vdash \sim e_1 \xrightarrow{1} e_2}$	<b>Escape <math>n+2</math>:</b> $\frac{\Gamma \vdash e_1 \xrightarrow{n+1} e_2}{\Gamma \vdash \sim e_1 \xrightarrow{n+2} \sim e_2}$
<b>Run 0:</b> $\frac{\Gamma \vdash e \hookrightarrow \langle e_1 \rangle^\circ \quad \bullet \vdash e_1 \hookrightarrow v_1}{\Gamma \vdash \text{run } e \hookrightarrow v_1}$	<b>Run <math>n+1</math>:</b> $\frac{\Gamma \vdash e_1 \xrightarrow{n+1} e_2}{\Gamma \vdash \text{run } e_1 \xrightarrow{n+1} \text{run } e_2}$
<b>Constant 0:</b> $\Gamma \vdash \uparrow v \hookrightarrow v$	<b>Constant <math>n+1</math>:</b> $\Gamma \vdash \uparrow v \xrightarrow{n+1} \uparrow v$

Figure 1: The Semantics of  $\lambda^M$



## 11 Optimizations

### 11.1 Safe Beta Reduction

To write multi-stage programs effectively, one needs to observe the programs produced, and these programs should be as simple as possible. For this reason, our implementation performs automatic *safe-beta reduction* on constants and variables. A beta reduction is safe if it does not change evaluation order, or effect termination properties. There is one safe case which is particularly easy to recognize, namely, Plotkin's  $\beta_v$  rule [26]. Whenever an application is constructed where the function part is an explicit lambda abstraction, and the argument part is a value, then that application can be symbolically beta reduced. In order to avoid duplicating code we restrict our optimizations to constants or variables (while Plotkin's  $\beta_v$  rule also allows the values to be lambda expressions). For example in:

```
val g = <fn x => x * 5>;
val h = <fn x => (~g x) - 2>;
```

The variable  $h$  evaluates to:  $\langle \text{fn } d1 \Rightarrow (d1 * 5) - 2 \rangle$  rather than  $\langle \text{fn } d1 \Rightarrow ((\text{fn } d2 \Rightarrow d2 * 5) d1) - 2 \rangle$ .

We realize of course that this might make it hard to understand *why* a particular program was generated. In our experience, the resulting smaller, simpler programs, are easier to understand and make this tradeoff worthwhile.

### 11.2 Nested Escapes

When we first wrote programs with more than two levels we observed that our programs took a long time to run. We traced this to rule **Escape  $n+2$**  of our semantics. Consider the case where a deeply bracketed term  $e$  at level  $n$  is escaped all the way to level 0. In order to execute this term (which escapes to level 0) it must be rebuilt  $n$  times. Consider the reduction sequence sketched below for the term  $\text{run } (\text{run } \langle \langle \sim e \rangle \rangle)$ , where  $e$  is bound in  $\Gamma$  to  $\langle 5 \rangle$ , of which we show only the innermost run.

$$\begin{array}{c}
 \frac{e \hookrightarrow \langle 5 \rangle^\circ}{\sim e \xrightarrow{1} \langle 5 \rangle} \\
 \frac{\sim e \xrightarrow{2} \sim \langle 5 \rangle}{\langle \sim e \rangle \xrightarrow{1} \langle \sim \langle 5 \rangle \rangle} \\
 \frac{\langle \sim e \rangle \xrightarrow{1} \langle \sim \langle 5 \rangle \rangle}{\langle \langle \sim e \rangle \rangle \hookrightarrow \langle \langle \sim \langle 5 \rangle \rangle \rangle^\circ} \\
 \frac{\langle \langle \sim e \rangle \rangle \hookrightarrow \langle \langle \sim \langle 5 \rangle \rangle \rangle^\circ}{\text{run } \langle \langle \sim e \rangle \rangle \hookrightarrow \langle 5 \rangle^\circ}
 \end{array}$$

For two levels the term is rebuilt 2 times. For three levels the term is rebuilt 3 times. A simple refinement can prevent this from happening. We change the rebuilding of escaped expressions at levels greater than 1 by adding the rule **Escape Opt  $n+2$**  in addition to the rule **Escape  $n+2$** .

$$\text{Escape Opt } n+2: \frac{\Gamma \vdash e_1 \xrightarrow{n+1} \langle e_2 \rangle}{\Gamma \vdash \sim e_1 \xrightarrow{n+2} e_2}$$

$$\text{Escape } n+2: \frac{\Gamma \vdash e_1 \xrightarrow{n+1} e_2}{\Gamma \vdash \sim e_1 \xrightarrow{n+2} \sim e_2}$$

Thus a long sequence of escapes surrounded by an equal number of brackets gets rebuilt exactly one. This optimization is safe since there are no variables in a rebuilt term. So rebuilding it more than once performs no useful work. This correctness of this optimization follows from the fact that under our semantics inside nesting brackets  $\sim \langle e \rangle$  is always equal to  $e$ .

## 12 Discussion and Related Works

A summary of the distinguishing characteristics of other work on multi-stage languages is shown in figure 2. In the figure we list whether or not that work contains the properties we believe to be important. Those properties are:

- **Staging:** Staging can be expressed in all these languages, in the sense that code representing expressions can be expressed, and that applications can be delayed to future stages. In this row, “2” means two stages are supported, and “+” means arbitrary number of stages is supported.
- **Strong Typing:** In a single stage language, strong typing ensures that a well-typed program “cannot go wrong” by applying functions to arguments of the wrong type. In addition, multi-stage programs can go wrong if a computation attempts to use a variable before it is available. A strongly typed multi-stage language protects against both kinds of errors.
- **Monolithic Variables:** Whether a variable from one stage can be used in all reasonable ways at that same stage. In the example,  $f$  is applied to  $\langle x \rangle$ , and hence  $f$  is acting as a “code transformer”.
- **run or eval function:** The languages also come with different set of multi-stage programming constructs. In particular, not all the language allow reflection in the form of a run or eval function. Most notably, even in the most recent work of Davies, it was not known how eval could be included in the language [6].
- **Cross-Stage Persistence:** This is the most distinguishing feature of MetaML, and has been discussed in detail in this paper. To our knowledge, this feature has not been proposed or incorporated into any multi-stage programming language. In essence, cross-stage persistence means that the programmer can use a variable bound in any stage in an expression executed in any future stage. Alternatively, it allows us to stage expressions containing free variables, as long as the type system is satisfied that these variables are available before this expression is evaluated.
- **Cross-Platform Portability:** The ability to print generated code as text. This allows code generated on one machine to be compiled and run on other machines. If code embeds its own “local environment” this becomes considerably more difficult. The loss of this ability is the price paid for cross stage persistence.

In what follows is an historical perspective of the work highlighted in the table:

- **NN:** Nielson and Nielson pioneered the investigation of staged languages with their two-level functional language [25, 24]. They presented rules for the well-formedness of the binding-times of expressions in the

Facility	Example	NN [25]	GA [11]	GB [9]	Th [31]	HG [12]	$\lambda^\square$ [7]	$\lambda^\circ$ [6]	$\lambda^M$
Staging	$\langle \lambda x. x \rangle$	2	2	+	2	+	+	+	+
Strong Typing		Y	1	N	N	N	Y	Y	Y
Monolithic Variables	$\langle \lambda x. \sim(f \langle x \rangle) \rangle$	Y	Y	Y	Y	Y	N	Y	Y
Reflection	<i>run or eval</i>	N	N	N	Y	N	Y	N	Y
X-Stage Persistence	$\lambda f. \langle \lambda x. f \ x \rangle$	N	N	N	N	N	N	N	Y
X-Platform Portability		Y	Y	Y	Y	Y	Y	Y	N

Figure 2: Comparative feature set

language, from which MetaML’s type rules are derived. They also sketched guidelines for a multi-stage (“B-level”) language. The two-level language is widely used to describe binding-time annotations in the partial evaluation literature.

- GA: Gomard and Jones use a statically-typed two-stage language for partial evaluation of the untyped  $\lambda$ -calculus [11]. The language allows the treatment of expressions containing monolithic free variables. They use a “const” construct only for constants of ground type. Our treatment of variables in the formal semantics is inspired by their work.
- GB: Glück and Jørgensen [9] present the novel idea of multi-level BTA, as a very efficient and effective alternate to multiple self-application. An untyped multi-level language based on Scheme is used for the presentation. Our study of MetaML is at a more basic level: MetaML is an abstract calculus. It is also notable that all intermediate results in GJ are printable, i.e., have an intensional representation. In MetaML, cross-stage persistence allows us to have intermediate results (between stages) that contain constants for which no intensional representation is available. While this is very convenient for run-time code generation, it makes the proper specification of MetaML more difficult. For example, we can’t use [9]’s “Generic Code Generation functions” to define the language. A second paper by Glück and Jørgensen [10] demonstrates the impressive efficiency of MBTA, and the use of constraints-solving methods to perform the analysis. The MBTA is type-based, but underlying language is dynamically typed.
- Th: Thiemann [31] studies a two-level language with *eval*, *apply*, and *call/cc*, in the context of studying the partial evaluation of a greater subset of scheme than was done previously. A BTA based on constraint-solving is presented. Although the problems with *eval* and *call/cc* are highlighted, and unlike with MetaML, there is no explicit notion of partially static types, and so, the complexity of introducing *eval* into a multi-stage language does not manifest itself. Thiemann also deals with the issue of variable-arity functions, which is a practical problem when dealing with *eval* in Scheme.
- HG: Hatcliff & Glück studied a multi-stage flow-chart language called S-Graph-n, and thoroughly investigated the issues involved in the implementation of such a language [12]. The syntax of S-Graph-n explicitly captures all the information necessary for specifying the staging of a computation: each construct is annotated with a number indicating the stage during which it is to be executed, and all variables are annotated with

a number indicating the stage of their availability. S-Graph-n is not statically typed, and the syntax and formal semantics of the language are quite sizable. Programming in S-Graph-n requires the user to annotate every construct and variable with stage annotations, and ensuring the consistency of the annotations is the user’s responsibility. In their work, Hatcliff & Glück identified *language-independence* of the internal representation of “code” as an important characteristic of any multi-stage language.

- $\lambda^\square$ : Davies & Pfenning presented the first statically-typed multi-stage language Mini-ML $^\square$  [6]. The type system is motivated by constructive modal logic, and a formal proof is presented for the equivalence of binding-time correctness and modal correctness. In contrast, the MetaML type-system was motivated primarily by operational considerations. Despite the different origins, the languages have a lot in common: meta-brackets, escape, and run roughly correspond to box, unbox, and eval respectively. Mini-ML $^\square$ ’s  $\square$  type constructor is also similar to *code*. Interestingly, in Mini-ML $^\square$ , eval is defined in terms of unbox, whereas in MetaML, neither run or escape can be defined in terms of each other. Also, while Mini-ML $^\square$  can simulate persistence for code values, a stage-zero function, for example, cannot be made persistent. Finally, Mini-ML $^\square$  allows expressing functions like *back* are not expressible in the language.
- $\lambda^\circ$ : The multi-stage language Mini-ML $^\circ$  [6] is motivated by a linear-time constructive modal logic. The language allows staged expressions to contain monolithic free variables. The two constructs of Mini-ML $^\circ$ ; *next* and *prev*, correspond quite closely to MetaML’s meta-brackets and escape. The type constructor  $\circ$  also corresponds (roughly) to *code*. Unfortunately, *eval* is no longer expressible in the language.

Sheard has also investigated richer type systems for multi-staged programming. Sheard and Nelson investigated a two-stage language for the purpose of program generation [28]. The base language was statically typed, and dependent types were used to generate a wider class of programs than is possible by MetaML restricted to two stages. Sheard and Shields [29] investigate a dynamic type systems for multi-staged programs where some type obligations of staged computations can be put off till run-time.

The type rule for run presented in this paper is motivated by the type system for runST [19].

### 13 Conclusion

We have described an  $n$ -stage multi-stage programming language which we call MetaML. MetaML was designed as a programming language. Our primary purpose was to support the writing of multi-stage programs. Because of this our design choices were different from those of other multi-stage systems. We find the following features essential when writing multi-stage programs.

- **Cross stage persistence.** The ability to use variables from any past stage is crucial to writing staged programs in the manner to which programmers are accustomed. Cross stage persistence provides a solution to hygienic macros in a typed language, i.e. macros which bind identifiers in the environment of definition, which are “captured” in the environment of use.
- **Multi-stage aware type system.** The type checker reports phase errors as well as type errors is crucial when debugging multi-stage programs, thus ensuring Cross-Stage Safety.
- **Display of code.** When debugging, it is important for users to observe the code produced by their programs. This implies a display mechanism (pretty-printer) for values of type code.
- **Display of Constants.** Constants originating from persistent variables are hard to identify. The % tags provide an approximation of where these constants came from. While potentially misleading they are often quite useful.
- **The Isomorphism between  $\langle A, 'c \rangle \rightarrow \langle B, 'c \rangle$  and  $\langle A \rightarrow B, 'c \rangle$ .** The isomorphism (which can only be expressed because of cross-stage persistence), reduces, sometimes drastically, the number of annotations needed to stage multi-stage programs.
- **Lift.** The lift annotation makes it possible to force computation in a early stage and lift this value into a program to be incorporated at a later stage. While never necessary (because of cross-stage persistence) it helps produce code which is easier to understand, because constants become explicit.
- **Safe beta reduction.** Safe beta reduction of the application of explicit abstractions to explicit variables reduces the clutter in generated code.

To further illustrate this we provide an extended example in Appendix A which stages a term rewriting system in which the the rewriting rules become known in the first stage and the terms to be rewritten become known only in later stages.

### 14 Future Work

We have built an implementation which was used to program the examples in this paper. Currently, the implementation supports polymorphic type-inference so that type information on bound variables is not necessary. We are currently extending this implementation to include all the features in core-ML.

We are also actively pursuing a subject reduction theorem for  $\lambda^M$ . The multi-level syntax makes the syntactic

approaches to type soundness [32] difficult to apply, because reduction contexts may appear inside lambda expressions at levels greater than zero. We have also found that the non-Hindley-Milner type judgement for the run annotation complicates matters considerably.

### References

- [1] C. Consel and O. Danvy. For a better support of static data flow. In J. Hughes, editor, *Functional Programming Languages and Computer Architecture, Cambridge, Massachusetts, August 1991 (Lecture Notes in Computer Science, vol. 523)*, pages 496–519. New York: ACM, Berlin: Springer-Verlag, 1991.
- [2] Charles Consel and Olivier Danvy. Tutorial notes on partial evaluation. In *ACM Symposium on Principles of Programming Languages*, pages 493–501, January 1993.
- [3] Roberto Di Cosmo. *Isomorphisms of Types: from  $\lambda$ -calculus to information retrieval and language design*. Progress in Theoretical Computer Science. Birkhäuser, 1995.
- [4] Olivier Danvy. Type-directed partial evaluation. In *ACM Symposium on Principles of Programming Languages*, pages 242–257, Florida, January 1996. New York: ACM.
- [5] Olivier Danvy, Karoline Malmkjaer, and Jens Palsberg. The essence of eta-expansion in partial evaluation. *LISP and Symbolic Computation*, 1(19), 1995.
- [6] Rowan Davies. A temporal-logic approach to binding-time analysis. In *Proceedings, 11<sup>th</sup> Annual IEEE Symposium on Logic in Computer Science*, pages 184–195, New Brunswick, New Jersey, 27–30 July 1996. IEEE Computer Society Press.
- [7] Rowan Davies and Frank Pfenning. A modal analysis of staged computation. In *23rd Annual ACM Symposium on Principles of Programming Languages (POPL'96)*, St.Petersburg Beach, Florida, January 1996.
- [8] Nachum Dershowitz. Computing with rewrite systems. *Information and Control*, 65:122–157, 1985.
- [9] Robert Glück and Jesper Jørgensen. Efficient multi-level generating extensions for program specialization. In *Programming Languages, Implementations, Logics and Programs (PLILP'95)*, volume 982 of *Lecture Notes in Computer Science*. Springer-Verlag, 1995.
- [10] Robert Glück and Jesper Jørgensen. Fasting binding-time analysis for multi-level specialization. In *PSI-96: Andrei Ershov Second International Memorial Conference, Perspectives of System Informatics*, volume 1181 of *Lecture Notes in Computer Science*. Springer-Verlag, 1996.
- [11] Carsten K. Gomard and Neil D. Jones. A partial evaluator for untyped lambda calculus. *Journal of Functional Programming*, 1(1):21–69, 1991.
- [12] John Hatcliff and Robert Glück. Reasoning about hierarchies of online specialization systems. In Olivier Danvy, Robert Glück, and Peter Thiemann, editors, *Partial Evaluation*, volume 1110 of *Lecture Notes in*

*Computer Science*, pages 161–182. Springer-Verlag, 1996.

- [13] J. Roger Hindley and Jonathan P. Seldin. *Introduction to Combinators and  $\lambda$ -Calculus*. Number 1 in London Mathematical Society Student Texts. Cambridge University Press, 1986.
- [14] Neil D. Jones. Mix ten years later. In *Partial Evaluation and Semantics-Based Program Manipulation, New Haven, Connecticut (Sigplan Notices, vol. 26, no. 9, September 1991)*, pages 24–38. New York: ACM, New York: ACM, June 1995.
- [15] Neil D. Jones, Carsten K Gomard, and Peter Sestoft. *Partial Evaluation and Automatic Program Generation*. Prentice-Hall, 1993.
- [16] Simon L. Peyton Jones and John Launchbury. Unboxed values as first class citizens in a non-strict functional language. In *Functional Programming and Computer Architecture*, September 91.
- [17] Richard B. Kieburtz, Francoise Bellegarde, Jef Bell, James Hook, Jeffrey Lewis, Dino Oliva, Tim Sheard, Lisa Walton, and Tong Zhou. Calculating software generators from solution specifications. In *TAPSOFT'95*, volume 915 of *LNCS*, pages 546–560. Springer-Verlag, 1995.
- [18] Richard B. Kieburtz, Laura McKinney, Jeffrey Bell, James Hook, Alex Kotov, Jeffrey Lewis, Dino Oliva, Tim Sheard, Ira Smith, and Lisa Walton. A software engineering experiment in software component generation. In *18th International Conference in Software Engineering*, March 1996.
- [19] John Launchbury and Amr Sabry. Monadic state: Axiomatization and type safety. In *Proceedings of the International Conference on Functional Programming*, Amsterdam, June 1997.
- [20] Mark Leone and Peter Lee. Deferred compilation: The automation of run-time code generation. Technical Report CMU-CS-93-225, Carnegie Mellon University, Dec. 1993.
- [21] Michael Lowry, Andrew Philpot, Thomas Pressburger, and Ian Underwood. Amphion: Automatic programming for scientific subroutine libraries. *NASA Science Information Systems Newsletter*, 31:22–25, 1994.
- [22] Robin Milner. A theory of type polymorphism in programming. *Journal of Computer and System Sciences*, 17:348–375, 1978.
- [23] Torben A. Mogensen. Efficient self-interpretation in lambda calculus. *Functional Programming*, 2(3):345–364, July 1992.
- [24] F. Nielson. Correctness of code generation from a two-level meta-language. In B. Robinet and R. Wilhelm, editors, *Proceedings of the European Symposium on Programming (ESOP 86)*, volume 213 of *LNCS*, pages 30–40, Saarbrücken, FRG, March 1986. Springer.
- [25] F. Nielson and H. R. Nielson. *Two-Level Functional Languages*. Number 34 in Cambridge Tracts in Theoretical Computer Science. Cambridge University Press, 1992.
- [26] G. D. Plotkin. Call-by-name, call-by-value- and the lambda-calculus. *Theoretical Computer Science*, 1:125–159, 1975.
- [27] Tim Sheard. A type-directed, on-line partial evaluator for a polymorphic language. In *Proceedings of the Symposium on Partial Evaluation and Semantics-Based Program Manipulation*, Amsterdam, June 1997.
- [28] Tim Sheard and Neal Nelson. Type safe abstractions using program generators. Technical Report OGI-TR-95-013, Oregon Graduate Institute of Science and Technology, 1995.
- [29] Tim Sheard and Mark Shields. Dynamic typing through staged type inference. In *Proceedings of the Second ACM International Conference on Functional Programming*, jun 1997. (submitted).
- [30] Brian C. Smith. *Reflection and Semantics in a Procedural Language*. PhD thesis, Massachusetts Institute of Technology, January 1982.
- [31] Peter Thiemann. Towards partial evaluation of full Scheme. In Gregor Kiczales, editor, *Reflection 96*, pages 95–106, San Francisco, CA, USA, April 1996.
- [32] Wright and Felleisen. A syntactic approach to type soundness. *Information and Computation (formerly Information and Control)*, 115, 1994.

## A Term Rewriting: An Extended Example

In this section we provide an extended example which illustrates multi-stage programming.

Dershowitz [8] defines a *term-rewriting system* as a set of directed rules. Each rule is made of a left-hand side and a right-hand side. A rule may be applied to a term  $t$  if a subterm  $s$  of  $t$  matches the left-hand under some substitution  $\sigma$ . A rule is applied by replacing  $s$  with  $t'$ , where  $t'$  is the result of applying the substitution  $\sigma$  to the right-hand side. We say “ $t$  rewrites (in one step) to  $t'$ ”, and write  $t \Rightarrow t'$ . The choice of which rule to apply is made non-deterministically. As an example, here are the rules for a Monoid [8]:

$$\begin{aligned} r_1 : \quad & x + 0 \rightarrow x \\ r_2 : \quad & 0 + x \rightarrow x \\ r_3 : \quad & x + (y + z) \rightarrow (x + y) + z \end{aligned}$$

Variables  $x$ ,  $y$ , and  $z$  in the rules can each match any term. If a variable occurs more than once on the left-hand side of a rule, all occurrences must match identical terms. We call this the *compatibility condition*. These rules allow us to have derivations such as:

$$\begin{aligned} & (a + b) + (0 + (d + e)) \\ \Rightarrow & \text{by } r_2, \sigma = [(d + e)/x] \\ & (a + b) + (d + e) \\ \Rightarrow & \text{by } r_3, \sigma = [(a + b)/x, d/y, e/z] \\ & ((a + b) + d) + e \end{aligned}$$

where the subterm being rewritten ( $s$ ) has been underlined.

Generally, the rules do not change over the life of the system. At the same time, the basic form of the matching function is a simultaneous traversal of a subject term and the left-hand side of the rule it is being matched against. This offers an opportunity for staging: We can “specialize”

matching over the rules in a first stage, and eliminate the overhead of traversing the left-hand side of the rules. Not only that, but as we will see, we can also remove a significant amount of administrative computations involved in constructing and applying the substitution  $\sigma$ . One would expect that this would significantly speed up the rewriting system.

Terms can be implemented using the following MetaML type:

```
datatype 'a T = Const of int
              | Var of string
              | Op of ('a * string * 'a);
```

The declaration introduces a type constructor  $T$  parameterized by another type,  $'a$ . The constructors of this type are `Const`, `Var` and `Op` representing one level of an integer, variable, or infix binary term, respectively. The following definitions will also be used:

```
datatype Term = In of (term T);
datatype 'a Maybe = Nothing of unit | Just of 'a;
type Sub = ((string * term) List) Maybe;
type Rule = (Term * Term);
```

```
fun const x = In (Const x);
fun var x = In (Var x);
fun op x = In (Op x);
fun out (In x) = x;
```

where `Term` ties the recursive knot in  $T$  to represent terms, and `Sub` is the type of substitutions, respectively. The `Maybe` type provides a mechanism to write functions that can fail to return a meaningful value by returning `Nothing ()` instead. Since pattern matching may fail the substitution type is a `Maybe` type. We represent rules as an ordered pair of terms. Thus,  $r_1$  would be `(op(var "x", "+", int "0"), var "x")`. The function `out : Term -> term T` removes a term from the recursive type to the 1 level type.

To simplify, we focus on `match`, a function that matches a rule with the whole subject term, that is, the special case when  $s = t$ . We assume that a separate helper function applies `match` to all subterms of the subject term.

Now, let us consider what the result of specializing `match` over a rule should look like. If we take  $r_1 : x + 0 \rightarrow x$  as an example, then an good result should have a form close to:

```
fun rewriteR1 term =
  case (out term) of
    Op (t1,s,t2) =>
      if streq(s, "+")
      then (case (out t2) of
              Const n => if n=0 then t1 else term
              | _ => term)
            else term
    | _ => term;
```

Here there is no interpretive overhead of traversing the rule, and the substitution operation has also been computed in the first stage. In the rest of this section, we will see how far towards this goal we can go using MetaML.

### A.1 An Implementation of match

Figure 3 presents the complete code for `match`, together with a staged version. Note that if we erase the annotations from the staged version, we get the source version back. The

source match function takes a pattern term `pat`, a substitution `msigma`, a term `term`, and returns a substitution. Recall that a substitution is a maybe type and may be `Nothing()`.

Hence, the type of the match function is  $\Theta \rightarrow \Sigma \rightarrow \Theta \rightarrow \Sigma$ , where  $\Theta \equiv \text{Term}$ , and  $\Sigma \equiv \text{Sub}$ .

If `match` is passed an invalid substitution `Nothing ()` the outer case-statement propagates the failure. Otherwise, if the pattern term is a variable, the substitution is extended appropriately after checking the compatibility condition. Similarly, if the pattern term is an operator, it must be checked that the subject term is also an operator with the same operation, and then the right- and left-hand-sides of the pattern and the subject term are recursively matched, extending the substitution. If the pattern term is an integer, then the subject term must also be an integer with the same value. In this case the original substitution is returned. In all other cases, `match` returns `Nothing ()`, indicating that `match` has failed.

The staged match function has type  $\Theta \rightarrow \langle \Sigma, 'a \rangle \rightarrow [4] \langle \Theta, 'a \rangle \rightarrow \langle \Sigma, 'a \rangle$ . The type indicates that the pattern term is inspected only in the first stage, and the result is a specialized function that can be run in a future stage. We can also define and annotate `rewrite` using the helper function `wrapper` as follows:

```
(* wrapper : Rule -> <Term> -> <Term> *)
fun wrapper (lhs,rhs) term =
  let val ms = match lhs <Just []> term
  in
    < case (~ms) of
      Nothing () => ~term
      | Just (sigma) =>
        subst sigma rhs > end ;

    (* rewrite : Rule -> <Term -> Term> *)
    fun rewrite rule = <fn x => ~(wrapper rule <x>)>;
```

In Figure 4 the code generated for  $r_1 : x + 0 \rightarrow x$  appears. The traversal of the pattern term has been performed. Yet, compared with the code we derived by hand, there are too many nested case-statements, and the calls to `subst` and `compatible` have not been performed.

Careful inspection shows that it should, in fact, be possible to reduce the nested case-statements by meaning-preserving transformations. If the outer case-statements could be “pushed” through the inner ones, then we would be able to simplify all the values at the leaves of the inner if-statements. In particular, in every case where we return `Nothing ()`, the unchanged term `t1` would be returned, and where `Just ( $\sigma$ )` is returned then the substitution `subst  $\sigma$`  could be performed on `rhs`.

Given that we were able to write the function `rewriteR1` by hand, it should be clear that values of `Maybe`-type need not appear in the generated code. Unfortunately, this cannot be achieved by using only the staging annotations on the current program.

In particular, the test of the compatibility condition in `match` cannot be performed at generation time since the subject term is unknown. At the same time, this call determines whether `Nothing` or `Just` is returned. This implies `match` must return a value of type  $\langle \text{Maybe}, 'a \rangle$ . This means the generated code will contain values of `Maybe` type.

Similarly, threading `msigma` through the recursive calls to `match` stops us from being able to reduce the calls to `compatible` and `subst` at generation time: The result of both of these functions depends on the substitution list, but because we are forced to annotate `msigma` to have type:

<pre>(* Source match: Term -&gt; Sub -&gt; Term -&gt; Sub *)  fun match pat msigma term = case (msigma) of   Nothing () =&gt; Nothing ()   Just (sigma) =&gt; (case (out pat) of   Var u =&gt;     if compatible u sigma term     then Just (cons((u,term),sigma))     else Nothing ()   Op (t11,s1,t12) =&gt;   (case (out term) of     Op (t21,s2,t22) =&gt;       (if streq(s2,(s1))        then (match t11               (match t12 msigma t22)               t21)        else Nothing ())     _ =&gt; Nothing ())   Const n =&gt;   (case (out term) of     Const u =&gt; if u=n                 then msigma                 else Nothing ()     _ =&gt; Nothing ());</pre>	<pre>(* Annotated match: Term -&gt; &lt;Sub&gt; -&gt; &lt;Term&gt; -&gt; &lt;Sub&gt; *)  fun match pat msigma term = &lt;case `msigma of   Nothing () =&gt; Nothing ()   Just (sigma) =&gt;   ~(case (out pat) of     Var u =&gt;       &lt;if compatible u sigma `term       then Just (cons((u,`term),sigma))       else Nothing ()&gt;     Op (t11,s1,t12) =&gt;     &lt;case (out `term) of       Op (t21,s2,t22) =&gt;         (if streq(s2, `(lift s1))          then ~(match t11                   (match t12 msigma &lt;t22&gt;)                   &lt;t21&gt;)          else Nothing ())       _ =&gt; Nothing ()&gt;     Const n =&gt;     &lt;case (out `term) of       Const u =&gt; if u= `(lift n)                   then `msigma                   else Nothing ()     _ =&gt; Nothing ()&gt;&gt;);</pre>
--	--

Figure 3: Normal and Annotated versions of the function match.

```
<<(fn t1 =>
  (case (case %out t1 of
    Op(t1,op1,tr) =>
      if %streq (op1,"+")
      then (case (case %out tr of
        Const n =>
          if %=(n,0) then Just(nil) else Nothing ()
        | Var s => Nothing ()
        | Op s => Nothing ()) of
          Nothing () => Nothing ()
          | Just s => if %compatible "x" s t1 then Just (%cons (("x",t1),s)) else Nothing ()
        else Nothing ()
        | Const n => Nothing ()
        | Var s => Nothing ()) of
          Nothing () => t1
          | Just s => %subst s In Var "x"))>>
```

Figure 4: Code generated from rule  $r_1 = x + 0 \rightarrow x$  by function match.

`<((string * Term)List)Maybe, 'a >` and not `((string * <Term, 'a >)List)Maybe`, these two functions cannot “access” the list at generation time.

That our first attempt at staging our example was only partly successful, should not be too surprising. Users of partial evaluation systems restructure their programs to help BTA succeed all the time. Just because binding-time annotations are placed manually shouldn't exempt us from this requirement. Because we can touch, see, and experiment with both the explicit annotations and the code returned, it helps us understand and reason about what is going on. Using the type system to filter out obviously incorrectly phased programs was also extremely useful. The strength of MetaML is the mental model it provides to reason about what is going on.

## A.2 Continuation-Passing Style

Our solution needs to propagate the context (performing a substitution over the right-hand side of the rule) into the leaves of the nested case expressions. This suggests rewriting the source program in continuation-passing style (CPS). This has been found to be quite useful in partial evaluation systems [1].

Figure 5 shows both a source version of `match` using an explicit continuation and an annotated two stage version of the same function. The function takes a pattern `pat`, a continuation `k`, a substitution `msigma`, and a term `term`. This function actually has a rather polymorphic type since nothing constrains the value returned by the continuation `k`. In the figure we have constrained the continuation to have type substitution to term. The continuation should apply the substitution to the right-hand side of the rule, or return the term unchanged if the substitution is a failure. Thus the wrapper function could be written:

```
fun wrapper (lhs,rhs) term =
```

<pre>(* Source *) type Cont = Sub -&gt; Term  (* match : Term -&gt; Cont -&gt; Sub -&gt; Term -&gt; Term *) fun match pat k msigma term = case (msigma) of   Nothing () =&gt; k (Nothing())   Just (sigma) =&gt; (case (out pat) of   Var u =&gt;     if compatible u sigma term     then k (Just (cons((u,term),sigma)))     else k (Nothing ())   Op (t11,s1,t12) =&gt;   (case (out term) of     Op (t21,s2,t22) =&gt;       (if streq(s2,s1)        then (match t11                 (fn s =&gt; match t12 k s t22)                 msigma                 t21)        else k (Nothing ()))     _ =&gt; k(Nothing ()))   Const n =&gt;   (case (out term) of     Const u =&gt; if u=n                 then k msigma                 else k (Nothing ())     _ =&gt; k(Nothing ())));</pre>	<pre>(* Annotated *) type Sub' = ((string * &lt;Term&gt;) List) Maybe; type Cont' = Sub' -&gt; &lt;Term&gt;  (* match : Term -&gt; Cont' -&gt; Sub' -&gt; &lt;Term&gt; -&gt; &lt;Term&gt; *) fun match pat k msigma term = case (msigma) of   Nothing () =&gt; k (Nothing())   Just (sigma) =&gt; (case (out pat) of   Var u =&gt;     &lt;if ~(compatible' u sigma term)     then ~(k (Just (cons((u,term),sigma))))     else ~(k (Nothing ()))&gt;   Op (t11,s1,t12) =&gt;   &lt;case (out ~term) of     Op (t21,s2,t22) =&gt;       (if streq(~(lift s1),s2)        then ~(match t11                 (fn s =&gt; match t12 k s &lt;t22&gt;)                 msigma                 &lt;t21&gt;)        else ~(k (Nothing ())))     _ =&gt; ~(k(Nothing ()))&gt;   Const n =&gt;   &lt;case (out ~term) of     Const u =&gt; if u= ~(lift n)                 then ~(k msigma)                 else ~(k (Nothing ()))     _ =&gt; ~(k(Nothing ()))&gt;;</pre>
---	--

Figure 5: Normal and Annotated versions of the CPS style match.

```
match lhs
  (fn Nothing () => term | Just s => subst rhs s)
  (Just [])
  term;
```

The annotated version is again remarkably similar to the original except for the annotations. The only difference is that we needed to write a staged version of the compatible function, since in the staged version a substitution (Sub') maps a string to a piece of code with type term, rather than a term. We call this function compatible'. It returns a piece of code with type bool.

```
fun compatible' u sigma term =
case find u sigma of
  Nothing() => <true>
| Just w => <termeq ~w ~term>;
```

Finally, the code generator can be constructed by supplying a suitable substitution continuation. This function needs an annotated substitution function which returns a piece of code with type term rather than a term, and is given below:

```
fun subst' (In t) sig =
case (t) of
  Var v =>
    (case find v sig of
      Nothing _ => <In (Var ~(lift v))>
    | Just w => w)
| Op (t1,s,t2) =>
  <In (Op (~(subst' t1 sig),
           ~(lift s),
           ~(subst' t2 sig))) >
| Const i => <In (Const ~(lift i))>;

fun wrapper (lhs,rhs) term =
match lhs
  (fn Nothing () => term | Just s => subst' rhs s)
```

```
(Just [])
term;
```

When this function is used to generate code for the rule  $r_1 = x + 0 \rightarrow x$  the following is generated:

```
<fn t =>
  (case %out t of
    Op(t1,s,tr) =>
      if %streq ("+",s)
      then (case %out tr of
        Const n => if %= (n,0) then t1 else t
        | Var s => t
        | Op z => t)
      else t
    | Const n => t
    | Var z => t)>
```

In addition to being compact and free of reducible nested case-statements or calls to subst or compatible, this output is virtually identical to the idealized code we presented for rewriteR1. We have observed that the code generated for a variety of other rules is equally as compact.