

# Typed Compilation of Recursive Datatypes

Joseph C. Vanderwaart    Derek R. Dreyer    Leaf Petersen  
Karl Crary                Robert Harper  
December 2001

School of Computer Science  
Carnegie Mellon University  
Pittsburgh, PA 15213

## Abstract

Standard ML employs an opaque (or generative) semantics of datatype definitions, in which every datatype definition produces a new type that is different from any other type, including other identically defined datatypes. A natural way of accounting for this is to consider the types defined in datatype definitions to be abstract types. When this interpretation is applied to type-preserving compilation, however, it has the unfortunate consequence that datatype constructors cannot be inlined, substantially increasing the run-time cost of constructor invocation compared to a traditional compiler. In this paper we examine two approaches to eliminating function call overhead from datatype constructors. First, we consider a transparent interpretation of datatypes that does away with generativity; and second, we alter the opaque interpretation by replacing datatype constructors with coercions that have no run-time effect or cost.

This research was sponsored by the Advanced Research Projects Agency CSTO under the title “The Fox Project: Advanced Languages for Systems Software”, ARPA Order No. C533, issued by ESC/ENS under Contract No. F19628-95-C-0050. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing official policies, either expressed or implied, of the Defense Advanced Research Projects Agency or the U.S. Government.

**Keywords:** typed compilation, datatypes, type theory, coercions

# 1 Introduction

The programming language Standard ML (SML) [6] provides a distinctive mechanism for defining recursive types, known as a *datatype declaration*. For example, the following declaration defines the type of lists of integers:

```
datatype intlist = Nil | Cons of int * intlist
```

This datatype declaration introduces the type `intlist` and two *constructors*: `Nil` represents the empty list, and `Cons` combines an integer and a list to produce a new list. For instance, the expression `Cons (1, Cons (2, Cons (3, Nil)))` has type `intlist` and corresponds to the list `[1, 2, 3]`. Values of this datatype are deconstructed by a case analysis that examines a list and determines whether it was constructed with `Nil` or with `Cons`, and in the latter case, extracting the original integer and list.

An important aspect of SML datatypes is that they are *generative*. That is, every datatype declaration defines a type that is distinct from any other type, including those produced by other, possibly identical, datatype declarations. The formal Definition of SML [6] makes this precise by stating that a datatype declaration produces a new type name, but does not associate that name with a definition; in this sense, datatypes are similar to abstract types. Harper and Stone [5] have given a type-theoretic interpretation of SML by exhibiting a translation from SML into a simpler typed *internal language*. This translation is faithful to the Definition of SML in the sense that, with a few well-known exceptions, it translates an SML program into a well-typed IL program if and only if the SML program is well-formed according to the Definition; consequently it is legitimate to consider the Harper-Stone interpretation as an alternative formulation of the Definition. Harper and Stone capture datatype generativity by translating a datatype declaration as a module containing an abstract type and functions to construct and deconstruct values of that type; thus in the setting of the Harper-Stone interpretation, datatypes *are* abstract types.

The generativity of datatypes poses some challenges for type-directed compilation of SML. In particular, although the Harper-Stone interpretation is easy to understand and faithful to the Definition of SML, it is inefficient when implemented naïvely. The problem is that construction and deconstruction of datatype values require calls to functions exported by the module defining the datatype; this is unacceptable given the ubiquity of datatypes in SML code. Conventional compilers, which disregard type information after an initial type-checking phase, may dispense with this cost by *inlining* those functions; that is, they may replace the function calls with the actual code of the corresponding functions to eliminate the call overhead. A type-directed compiler, however, does not have this option since all optimizations, including inlining, must be type-preserving. Moving the implementation of a datatype constructor across the module boundary violates type abstraction and thus results in ill-typed intermediate code. This will be made more precise in Section 2.

In this paper, we will discuss two potential ways of handling this performance problem, both in the context of the TILT/ML compiler developed at CMU [8, 11]. The first is to do away with datatype generativity altogether, replacing the abstract types in the Harper-Stone interpretation with concrete ones. We call this approach the *transparent interpretation of datatypes*. Clearly, a compiler that does this is *not* an implementation of Standard ML, and we will show that, although the modified language does admit inlining of datatype constructors, it has some unexpected properties. The second approach, which we have adopted in the most recent version of the TILT compiler, takes advantage of the way values of recursive type are represented at run time in programs compiled by TILT. In particular, since a value of recursive type is represented the same as its unrolling, we can observe that the mediating functions produced by the Harper-Stone interpretation all behave like the identity function at run time. We replace these functions with special

values called *coercions* and argue that this allows a compilation strategy that generates code with a run-time efficiency comparable to what would be attained if datatype constructors were inlined. We call this the *coercion interpretation of datatypes*.

The paper is structured as follows. Section 2 gives the details of the Harper-Stone interpretation of datatypes (which we also refer to as the *opaque interpretation of datatypes*) and points out the problems with inlining. Section 3 discusses the transparent interpretation; section 4 gives the coercion interpretation and discusses its properties. Section 5 discusses related work and concludes.

## 2 The Opaque Interpretation

In this section, we will review the parts of Harper and Stone’s interpretation of SML that are relevant to our discussion of datatypes. In particular, after defining the notation we use for our internal language, we will give an example of the Harper-Stone elaboration of datatypes. We will refer to this example throughout the paper. We will also review the way Harper and Stone define the matching of structures against signatures, and discuss the implications this has for datatypes. This will be important in Section 3, where we show some differences between signature matching in SML and signature matching under our transparent interpretation of datatypes.

### 2.1 Notation

Harper and Stone give their interpretation of SML as a translation, called *elaboration*, from SML into a typed internal language (IL). We will not give a complete formal description of the internal language we use in this paper; instead, we will use ML-like syntax for examples and employ the standard notation for function, sum and product types. For a complete discussion of elaboration, including a thorough treatment of the internal language, we refer the reader to Harper and Stone [5]. Since we are focusing our attention on datatypes, *recursive types* will be of particular importance. We will therefore give a precise description of the semantics of the form of recursive types we use.

The syntax for recursive types is given in Figure 1. Recursive types are separated into their

<i>Types</i>	$\sigma, \tau$	$::=$	$\dots \mid \alpha \mid \delta$
<i>Recursive Types</i>	$\delta$	$::=$	$\mu_i(\alpha_1, \dots, \alpha_n).(\tau_1, \dots, \tau_n)$
<i>Terms</i>	$e$	$::=$	$\dots \mid x \mid \text{roll}_\delta(e) \mid \text{unroll}_\delta(e)$
<i>Typing Contexts</i>	$\Gamma$	$::=$	$\epsilon \mid \Gamma, x : \tau \mid \Gamma, \alpha$

Figure 1: Syntax of Iso-recursive Types

own syntactic subcategory, ranged over by  $\delta$ . This is mostly a matter of notational convenience, as there are many times when we wish to make it clear that a particular type is a recursive one. A recursive type has the form  $\mu_i(\alpha_1, \dots, \alpha_n).(\tau_1, \dots, \tau_n)$ , where  $1 \leq i \leq n$  and each  $\alpha_j$  is a type variable that may appear free in any or all of  $\tau_1, \dots, \tau_n$ . Intuitively, this type is the  $i$ th in a system of  $n$  mutually recursive types. As such, it is isomorphic to  $\tau_i$  with each  $\alpha_j$  replaced by the  $j$ th component of the recursive bundle. Formally, it is isomorphic to the following somewhat unwieldy type:

$$\tau_i[\mu_1(\alpha_1, \dots, \alpha_n).(\tau_1, \dots, \tau_n), \dots, \mu_n(\alpha_1, \dots, \alpha_n).(\tau_1, \dots, \tau_n) / \tau_1, \dots, \tau_n]$$

(where, as usual, we denote by  $\tau[\sigma_1, \dots, \sigma_n / \alpha_1, \dots, \alpha_n]$  the simultaneous capture-avoiding substitution of  $\sigma_1, \dots, \sigma_n$  for  $\alpha_1, \dots, \alpha_n$  in  $\tau$ ). Since we will be writing such types often, we use some

notational conventions to make things clearer; these are shown in Figure 2. Using these shorthands, the above type may be written as  $\text{expand}(\mu_i(\alpha_1, \dots, \alpha_n).(\tau_1, \dots, \tau_n))$ .

$$\begin{aligned}
\vec{X} &\stackrel{\text{def}}{=} X_1, \dots, X_n \text{ for some } n \geq 1, \\
&\text{where } X \text{ is a metavariable, like } \alpha \text{ or } \tau \\
\text{length}(\vec{X}) &\stackrel{\text{def}}{=} n, \text{ where } \vec{X} = X_1, \dots, X_n \\
\mu\alpha.\tau &\stackrel{\text{def}}{=} \mu_1(\alpha).(\tau) \\
\vec{\mu}(\vec{\alpha}).(\vec{\tau}) &\stackrel{\text{def}}{=} \mu_1(\vec{\alpha}).(\vec{\tau}), \dots, \mu_n(\vec{\alpha}).(\vec{\tau}), \\
&\text{where } \text{length}(\vec{\alpha}) = \text{length}(\vec{\tau}) = n \\
\text{expand}(\delta) &\stackrel{\text{def}}{=} \tau_i[\vec{\mu}(\vec{\alpha}).(\vec{\tau})/\vec{\alpha}], \text{ where } \delta = \mu_i(\vec{\alpha}).(\vec{\tau})
\end{aligned}$$

Figure 2: Shorthand Definitions

The judgment forms of the static semantics of our internal language are given in Figure 3, and the rules relevant to recursive types are given in Figure 4. Note that the only rule that can be used to judge two recursive types equal requires that the two types in question are the same (*i*th) projection from bundles of the same length whose respective components are all equal. In particular, there is no “unrolling” rule stating that  $\delta \equiv \text{expand}(\delta)$ ; type theories in which this equality holds are said to have *equi-recursive* types and are significantly more complex. The recursive types in our theory are *iso-recursive* types that are isomorphic, but not equal, to their expansions. The isomorphism is embodied by the **roll** and **unroll** operations at the term level; the former turns a value of type  $\text{expand}(\delta)$  into one of type  $\delta$ , and the latter is its inverse.

$\Gamma \vdash \text{ok}$	Well-formed context.
$\Gamma \vdash \tau \text{ type}$	Well-formed type.
$\Gamma \vdash \sigma \equiv \tau$	Equivalence of types.
$\Gamma \vdash e : \tau$	Well-formed term.

Figure 3: Relevant Typing Judgements

$$\begin{array}{c}
\frac{i \in 1..n \quad \forall j \in 1..n. \Gamma, \alpha_1, \dots, \alpha_n \vdash \tau_j \text{ type}}{\Gamma \vdash \mu_i(\alpha_1, \dots, \alpha_n).(\tau_1, \dots, \tau_n) \text{ type}} \\
\\
\frac{i \in 1..n \quad \forall j \in 1..n. \Gamma, \alpha_1, \dots, \alpha_n \vdash \sigma_j \equiv \tau_j}{\Gamma \vdash \mu_i(\alpha_1, \dots, \alpha_n).(\sigma_1, \dots, \sigma_n) \equiv \mu_i(\alpha_1, \dots, \alpha_n).(\tau_1, \dots, \tau_n)} \\
\\
\frac{\Gamma \vdash e : \text{expand}(\delta)}{\Gamma \vdash \text{roll}_\delta(e) : \delta} \qquad \frac{\Gamma \vdash e : \delta}{\Gamma \vdash \text{unroll}_\delta(e) : \text{expand}(\delta)}
\end{array}$$

Figure 4: Typing Rules for Iso-recursive Types

## 2.2 Elaborating Datatype Definitions

The Harper-Stone interpretation of SML includes a full account of datatypes, including generativity. The main idea is to encode datatypes as recursive sum types but hide this implementation behind an opaque signature. A datatype definition therefore elaborates as a structure that exports a number of abstract types and functions that construct and deconstruct values of those types. For example, consider the following pair of mutually recursive datatypes, representing expressions and declarations in the abstract syntax of a toy language:

```
datatype exp = VarExp of var | LetExp of dec * exp
and dec = ValDec of var * exp | SeqDec of dec * dec
```

The Harper-Stone elaboration of this definition is given in Figure 5, using ML-like syntax for readability. To construct a value of one of these datatypes, a program must use the corresponding

```
structure ExpDec :> sig
  type exp
  type dec
  val exp_in : var + (dec * exp) -> exp
  val exp_out : exp -> var + (dec * exp)
  val dec_in : (var * exp) + (dec * dec) -> dec
  val dec_out : dec -> (var * exp) + (dec * dec)
end = struct
  type exp =  $\mu_1(\alpha, \beta).(\text{var} + \beta * \alpha, \text{var} * \alpha + \beta * \beta)$ 
  type dec =  $\mu_2(\alpha, \beta).(\text{var} + \beta * \alpha, \text{var} * \alpha + \beta * \beta)$ 
  fun exp_in x = rollexp(x)
  fun exp_out x = unrollexp(x)
  fun dec_in x = rolldec(x)
  fun dec_out x = unrolldec(x)
end
```

Figure 5: Harper-Stone elaboration of `exp-dec` example.

`in` function; these functions each take an element of the sum type that is the “unrolling” of the datatype and produce a value of the datatype. More concretely, we implement the constructors for `exp` and `dec` as follows:

$$\begin{aligned} \text{VarExp}(x) &\stackrel{\text{def}}{=} \text{ExpDec.exp\_in}(\text{inj}_1(x)) \\ \text{LetExp}(d, e) &\stackrel{\text{def}}{=} \text{ExpDec.dec\_in}(\text{inj}_2(d, e)) \\ \text{ValDec}(x, e) &\stackrel{\text{def}}{=} \text{ExpDec.dec\_in}(\text{inj}_1(x, e)) \\ \text{SeqDec}(d1, d2) &\stackrel{\text{def}}{=} \text{ExpDec.dec\_in}(\text{inj}_2(d1, d2)) \end{aligned}$$

Notice that the types `exp` and `dec` are held abstract by the opaque signature ascription. This captures the generativity of datatypes, since the abstraction prevents `ExpDec.exp` and `ExpDec.dec` from being judged equal to any other types. However, as we mentioned in Section 1, this abstraction also prevents inlining of the `in` and `out` functions: for example, if we attempt to inline `exp_in` in the definition of `VarExp` above, we get

$$\text{VarExp}(x) \stackrel{\text{def}}{=} \text{roll}_{\text{ExpDec.exp}}(\text{inj}_1(x))$$

but this is ill-typed outside of the `ExpDec` module because the fact that `exp` is a recursive type is not visible. Thus performing inlining on well-typed code can lead to ill-typed code, so we say that inlining across abstraction boundaries is *not type-preserving* and therefore not an acceptable strategy for a typed compiler. The problem is that since we cannot inline `in` and `out` functions, our compiler must pay the run-time cost of a function call every time a value of a datatype is constructed or case-analyzed. Since these operations occur very frequently in SML code, this performance penalty is significant.

## 2.3 Datatypes and Signature Matching

Standard ML makes an important distinction between datatype *definitions*, which appear at the top level or in structures, and datatype *specifications*, which appear in signatures. As we have seen, the Harper-Stone interpretation elaborates datatype definitions as opaquely sealed structures; datatype specifications are translated into specifications of structures. For example, the signature

```
signature S = sig
  datatype intlist = Nil | Cons of int * intlist
end
```

contains a datatype specification, and elaborates as follows:

```
signature S = sig
  struct Intlist : sig
    type intlist
    val intlist_in : unit + int * intlist -> intlist
    val intlist_out : intlist -> unit + int * intlist
  end
end
```

A structure  $M$  will match  $S$  if  $M$  contains a structure `Intlist` of the appropriate signature.<sup>1</sup> In particular, it is clear that the structure definition produced by the Harper-Stone interpretation for the datatype `intlist` defined in Section 1 has this signature, so that datatype definition matches the specification above.

What is necessary in general for a datatype declaration to match a specification under this interpretation? Since datatype declarations are translated as structures, and datatype specifications are translated as structure declarations with signatures, matching a datatype declaration against a spec boils down to matching an opaquely sealed structure against a signature. Specifically, a specification of the form

```
datatype  $t_1 = \tau_1$  and ... and  $t_n = \tau_n$ 
```

(where the  $\tau_i$  are sum types) elaborates to a specification of a structure with the following signature:

---

<sup>1</sup>Standard ML allows only datatypes to match datatype specifications, so the actual Harper-Stone elaboration must use a name for the datatype that cannot be guessed by a programmer. We will not discuss this issue further.

```

sig
  type t1
  :
  type tn
  val t1_in : τ1 -> t1
  val t1_out : t1 -> τ1
  :
  val tn_in : τn -> tn
  val tn_out : tn -> τn
end

```

In order to match this signature, the structure corresponding to a datatype definition must define types named  $t_1, \dots, t_n$  and must contain `in` and `out` functions of the appropriate type for each. (Note that in any structure produced by elaborating a datatype definition under this interpretation, the  $t_i$ 's will be abstract types.) Thus, for example, if  $m \geq n$  then the datatype definition

```
datatype t1 = σ1 and ... and tm = σm
```

matches the above specification if and only if  $\sigma_i \equiv \tau_i$  for  $1 \leq i \leq n$ , since this is necessary and sufficient for the types of the `in` and `out` functions to match for the types mentioned in the specification.

### 3 A Transparent Interpretation of Datatypes

A natural approach to enabling the inlining of datatypes in a type-preserving compiler is to do away with the generative semantics of datatypes. In the context of the Harper-Stone interpretation, this corresponds to replacing the abstract type specification in the signature of a datatype module with a transparent type definition, so we call this modified interpretation the *transparent interpretation of datatypes* (TID).

#### 3.1 Making Datatypes Transparent

The idea of the transparent interpretation is to expose the implementation of datatypes as recursive sum types during elaboration, rather than hiding it. In our `expdec` example, this corresponds to changing the declaration shown in Figure 5 to the following (we continue to use ML-like syntax for readability):

```

structure ExpDec :> sig
  type exp = μ1(α, β).(var + β * α, var * α + β * β)
  type dec = μ2(α, β).(var + β * α, var * α + β * β)
  (* ... specifications for in and out functions same as before ... *)
end =
  (* ... same structure as before ... *)

```

Importantly, this change extends to datatype specifications as well as datatype definitions. Thus, a structure that exports a datatype must export its implementation transparently, using a signature similar to the one above—otherwise a datatype inside a structure would appear to be generative outside that structure, and there would be little point to the new interpretation.



As we have mentioned before, altering the interpretation of datatypes to expose their implementation as recursive types really creates a new language, which is neither a subset nor a superset of Standard ML. An example of the most obvious difference can be seen in Figure 6. In the figure, two datatypes are defined by seemingly identical definitions. In SML, because datatypes are generative, the two types `List1.t` and `List2.t` are distinct; since the variable `l` has type `List1.t` but is passed to `List2.Cons`, which expects `List2.t`, the function `switch` is ill-typed. Under the transparent interpretation, however, the implementations of both datatypes are exported transparently as  $\mu\alpha.\text{unit} + \text{int} * \alpha$ . Thus under this interpretation, `List1.t` and `List2.t` are equal and so `switch` is a well-typed function. It is clear that many programs like this one fail to type-check in SML but succeed under the transparent interpretation; what is less obvious is that there are some programs for which the opposite is true. We will discuss two main reasons for this.

```

structure List1 = struct
  datatype t = Nil | Cons of int * t
end
structure List2 = struct
  datatype t = Nil | Cons of int * t
end
fun switch List1.Nil = List2.Nil
  | switch (List1.Cons (n,l)) = List2.Cons (n,l)

```

Figure 6: Non-generativity of Transparent Datatypes

## 3.2 Problematic Datatype Matchings

Recall that according to the Harper-Stone interpretation, a datatype matches a datatype specification if the types of the datatype's `in` and `out` functions match the types of the `in` and `out` functions in the specification, after the actual type has been substituted for the abstract type in the spec. (Note: the types of the `out` functions match if and only if the types of the `in` functions match, so we will hereafter refer only to the `in` functions.) Under the transparent interpretation, however, it is also necessary that the recursive type implementing the datatype match the one given in the specification. This is not a trivial requirement; we will now give two examples of matchings that succeed in SML but fail under the transparent interpretation.

### 3.2.1 A Simple Example

A very simple example of a problematic matching is the following. Under the opaque interpretation, matching the structure

```

struct
  datatype u = A of u * u | B of int
  type v = u * u
end

```

against the signature

```

sig
  type v
  datatype u = A of v | B of int
end

```

amounts to checking that the type of the `in` function for  $u$  defined in the structure matches that expected by the signature once  $u * u$  has been substituted for  $v$  in the signature. (No definition is substituted for  $u$ , since it is abstract.) After substitution, the type required by the signature for the `in` function is  $u * u + \text{int} \rightarrow u$ , which is exactly the type of the function given by the structure, so the matching succeeds.

Under the transparent interpretation, however, the structure defines  $u$  to be  $\mathbf{u}_{imp} \stackrel{\text{def}}{=} \mu\alpha. \alpha * \alpha + \text{int}$  but the signature specifies  $u$  as  $\mu\alpha. v + \text{int}$ . In order for matching to succeed, these two types must be equivalent after we have substituted  $\mathbf{u}_{imp} * \mathbf{u}_{imp}$  for  $v$  in the specification. That is, it is required that

$$\mathbf{u}_{imp} \equiv \mu\alpha. \mathbf{u}_{imp} * \mathbf{u}_{imp} + \text{int}$$

Observe that the type on the right is none other than  $\mu\alpha. \text{expand}(\mathbf{u}_{imp})$ . (Notice also the bound variable  $\alpha$  does not appear free in the body of this  $\mu$ -type. Hereafter we will write such types with a wildcard `_` in place of the type variable to indicate that it is not used in the body of the  $\mu$ .) This equivalence does not hold for iso-recursive types, so the matching fails.

### 3.2.2 A More Complex Example

Another example of a datatype matching that is legal in SML but fails under the transparent interpretation can be found by reconsidering our running example of `exp` and `dec`. Under the opaque interpretation, a structure containing this pair of datatypes matches the following signature, which hides the fact that `exp` is a datatype:

```

sig
  type exp
  datatype dec = ValDec of var * exp | SeqDec of dec * dec
end

```

When this datatype specification is elaborated under the transparent interpretation, however, the resulting HIL signature exposes the implementation of `dec` as

$$\text{dec}_{spec} \stackrel{\text{def}}{=} \mu\alpha. \text{var} * \text{exp} + \alpha * \alpha$$

Elaboration of the definitions of `exp` and `dec`, on the other hand, produces an implementation of `dec` that is mutually recursive with `exp`:

$$\text{dec}_{imp} \stackrel{\text{def}}{=} \mu_2(\alpha, \beta). (\text{var} + \beta * \alpha, \text{var} * \alpha + \beta * \beta)$$

Matching the structure containing the datatypes against the signature can only succeed if  $\text{dec}_{spec} \equiv \text{dec}_{imp}$  (under the substitution of the implementation of `exp` for `exp` in  $\text{dec}_{spec}$ ). As we have already remarked, this equivalence does not hold because the two  $\mu$ -types have different numbers of components.

### 3.3 Problematic Signature Constraints

The module system of SML provides two ways to express sharing of type information between structures. The first, `where type`, modifies a signature by “patching in” a definition for a type the signature originally held abstract. The second, `sharing type`, asserts that two or more type names (possibly in different structures) refer to the same type. Both of these forms of constraints are restricted so that multiple inconsistent definitions are not given to a single type name; in the case of `sharing type`, for example, it is required that all the names be *flexible*, that is, they must either be abstract or defined as equal to another type that is abstract. Under the opaque interpretation, datatypes are abstract and therefore flexible, meaning they can be shared; under the transparent interpretation, datatypes are concretely defined and hence can never be shared. For example, the following signature is legal in SML:

```
signature S = sig
  structure M : sig datatype t = A | B end
  structure N : sig datatype t = A | B end
  sharing type M.t = N.t
end
```

We can write a similar signature using `where type`, which is also valid SML:

```
signature S = sig
  structure M : sig datatype t = A | B end
  structure N : sig datatype t = A | B end where type t = M.t
end
```

Neither of these signatures elaborates successfully under the transparent interpretation of datatypes, since under that interpretation the datatypes are transparent and therefore ineligible for either sharing or `where type`.

Another example is the following signature:

```
signature AB = sig
  structure A : sig
    type s
    val C : s
  end
  structure B : sig
    datatype t = C | D of A.s * t
  end
  sharing type A.s = B.t
end
```

(Again, we can construct an analogous example with `where type`.) Since the name `B.t` is flexible under the opaque interpretation but not the transparent, this code is legal SML but must be rejected under the transparent interpretation.

### 3.4 Relaxing Recursive Type Equivalence

We will now consider some modifications to the theory of type equivalence in our intermediate language that allow the problematic datatype matchings we have discussed to typecheck under the

transparent interpretation. Specifically, we will show how to weaken type equivalence (*i.e.*, make it equate more pairs of types) so that the problematic matchings described earlier in this section succeed. The ideas in this section are based upon the equivalence algorithm adopted by Shao [9] for the FLINT/ML compiler.

To begin, consider the simple  $u$ - $v$  example of Section 3.2.1. Recall that in that example, matching the datatype definition against the spec required proving the equivalence

$$\mathbf{u}_{imp} \equiv \mu\alpha. \mathbf{u}_{imp} * \mathbf{u}_{imp} + \mathbf{int}$$

where the type on the right-hand side is just  $\mu \_ . \text{expand}(\mathbf{u}_{imp})$ . By simple variations on this example, it is easy to show that in general, for the transparent interpretation to be as permissive as the opaque, the following recursive type equivalence must hold:

$$\delta \equiv \mu \_ . \text{expand}(\delta)$$

We refer to this as the *boxed-unroll* rule. It says that a  $\mu$ -projection is equal to its unrolling “boxed” by a  $\mu$ . Intuitively, this rule is needed because datatype matching succeeds under the opaque interpretation whenever the unrolled form of the datatype implementation equals the unrolled form of the datatype spec (because these are both supposed to describe the domain of the  $\mathbf{in}$  function). An alternative formulation of the boxed-unroll rule (equivalent to the first one by transitivity) makes two  $\mu$ -projections equal if their unrollings are equal, *i.e.*:

$$\frac{\text{expand}(\delta_1) \equiv \text{expand}(\delta_2)}{\delta_1 \equiv \delta_2}$$

Although the boxed-unroll equivalence is necessary for the transparent interpretation of datatypes to admit all matchings admitted by the opaque one, it is not sufficient; to see this, consider the problematic  $\mathbf{exp}$ - $\mathbf{dec}$  matching from Section 3.2.2. The problematic constraint in that example is

$$\mu\alpha. \mathbf{var} * \mathbf{exp} + \alpha * \alpha \equiv \mu_2(\alpha, \beta). (\mathbf{var} + \beta * \alpha, \mathbf{var} * \alpha + \beta * \beta)$$

The boxed-unroll rule is insufficient to prove this equivalence. In order to apply boxed-unroll to prove these two types equivalent, we must be able to prove that their unrollings are equivalent, in other words that

$$\mathbf{var} * \mathbf{exp} + \mathbf{dec}_{spec} * \mathbf{dec}_{spec} \equiv \mathbf{var} * \mathbf{exp} + \mathbf{dec}_{imp} * \mathbf{dec}_{imp}$$

But we cannot prove this without first proving  $\mathbf{dec}_{spec} \equiv \mathbf{dec}_{imp}$ , which is exactly what we set out to prove in the first place. The boxed-unroll rule is therefore unhelpful in this case.

The trouble is that proving the premise of the boxed-unroll rule (the equivalence of  $\text{expand}(\delta_1)$  and  $\text{expand}(\delta_2)$ ) may require proving the conclusion (the equivalence of  $\delta_1$  and  $\delta_2$ ). Similar problems have been addressed in the context of general equi-recursive types. In that setting, deciding type equivalence involves assuming the conclusions of equivalence rules when proving their premises [1, 2]. Applying this idea provides a natural solution to the problem discussed in the previous section. We can maintain a “trail” of type-equivalence assumptions; when deciding the equivalence of two  $\mu$ -projections, we add that equivalence to the trail before comparing their unrollings.

Formally, the equivalence judgement itself becomes  $\Gamma; A \vdash \sigma \equiv \tau$ , where  $A$  is a set of assumptions, each of the form  $\tau_1 \equiv \tau_2$ . All the equivalence rules in the static semantics must be modified

to account for the trail. In all the rules except those for  $\mu$ -projections, the trail is simply passed unchanged from the conclusions to the premises. There are two new rules for  $\mu$ -projections:

$$\frac{\tau_1 \equiv \tau_2 \in A}{\Gamma; A \vdash \tau_1 \equiv \tau_2} \quad \frac{\Gamma; A \cup \{\delta_1 \equiv \delta_2\} \vdash \text{expand}(\delta_1) \equiv \text{expand}(\delta_2)}{\Gamma; A \vdash \delta_1 \equiv \delta_2}$$

The first rule allows an assumption from the trail to be used; the second rule is an enhanced form of the boxed-unroll rule that adds the conclusion to the assumptions of the premise. It is clear that the trail is just what is necessary in order to resolve the **exp-dec** anomaly described above; before comparing the unrollings of  $\text{dec}_{\text{spec}}$  and  $\text{dec}_{\text{imp}}$ , we add the assumption  $\text{dec}_{\text{spec}} \equiv \text{dec}_{\text{imp}}$  to the trail; we then use this assumption to avoid the cyclic dependency we encountered before.

In fact, the trailing version of the boxed-unroll rule is sufficient to ensure that the transparent interpretation accepts all datatype matchings accepted by SML. To see why, consider a datatype specification

$$\text{datatype } \mathbf{t}_1 = \tau_1 \text{ and } \dots \text{ and } \mathbf{t}_n = \tau_n$$

(where the  $\tau_i$  are sum types in which the  $\mathbf{t}_i$  may occur). Suppose that some implementation matches this spec under the opaque interpretation; the implementation of each type  $\mathbf{t}_i$  must be a  $\mu$ -projection  $\delta_i$ . Furthermore, the type of the  $\mathbf{t}_i$ -in function given in the spec is  $\tau_i \rightarrow \mathbf{t}_i$ , and the type of its implementation is  $\text{expand}(\delta_i) \rightarrow \delta_i$ . Because the matching succeeds under the opaque interpretation, we know that these types are equal after  $\delta_i$  has been substituted for  $\mathbf{t}_i$ ; thus we know that  $\text{expand}(\delta_i) \equiv \tau_i[\vec{\delta}/\vec{\mathbf{t}}]$  for each  $i$ .

When the specification is elaborated under the transparent interpretation, however, the resulting signature declares that the implementation of each  $\mathbf{t}_i$  is the appropriate projection from a recursive bundle determined by the spec itself. That is, each  $\mathbf{t}_i$  is transparently specified as  $\mu_i(\vec{\mathbf{t}}).(\vec{\tau})$ . In order for the implementation to match this transparent specification, it is sufficient that for each  $i$ ,  $\Gamma; \emptyset \vdash \delta_i \equiv \mu_i(\vec{\mathbf{t}}).(\vec{\tau})$ . We will now show that this equivalence can be derived using the trail algorithm. For any set  $S \subseteq \{1, \dots, n\}$ , define  $A_S = \{\delta_i \equiv \mu_i(\vec{\mathbf{t}}).(\vec{\tau}) \mid i \in S\}$ . We can prove the following lemma (we omit the context  $\Gamma$  from the judgments, since it never changes):

**Lemma 1** *If  $\text{expand}(\delta_i) \equiv \tau_i[\vec{\delta}/\vec{\mathbf{t}}]$  for  $1 \leq i \leq n$ , then for any  $S \subseteq \{1, \dots, n\}$  and any  $j \in \{1, \dots, n\}$ ,  $A_S \vdash \delta_j \equiv \mu_j(\vec{\mathbf{t}}).(\vec{\tau})$ .*

**Proof Sketch:** The proof is by induction on  $n - |S|$ . If  $n - |S| = 0$ , then for any  $j$  the required equivalence is an assumption in  $A_S$  and can therefore be concluded using the assumption rule. If  $n - |S| > 0$ , then there are two cases:

**Case:**  $j \in S$ . Then the required equivalence is an assumption in  $A_S$ .

**Case:**  $j \notin S$ . Then let  $S' = S \cup \{j\}$ . Note that  $|S'| > |S|$  and so  $n - |S'| < n - |S|$ . By the induction hypothesis,  $A_{S'} \vdash \delta_k \equiv \mu_k(\vec{\mathbf{t}}).(\vec{\tau})$  for every  $k \in \{1, \dots, n\}$ . Because substituting equal types into equal types gives equal results,  $A_{S'} \vdash \tau_j[\vec{\delta}/\vec{\mathbf{t}}] \equiv \tau_j[\vec{\mu}(\vec{\mathbf{t}}).(\vec{\tau})/\vec{\mathbf{t}}]$ . By assumption,  $\text{expand}(\delta_j) \equiv \tau_j[\vec{\delta}/\vec{\mathbf{t}}]$ , so by transitivity  $A_{S'} \vdash \text{expand}(\delta_j) \equiv \tau_j[\vec{\mu}(\vec{\mathbf{t}}).(\vec{\tau})/\vec{\mathbf{t}}]$ . The type on the right side of this equivalence is just  $\text{expand}(\mu_j(\vec{\mathbf{t}}).(\vec{\tau}))$ , so by the trailing boxed-unroll rule we can conclude  $A_S \vdash \delta_j \equiv \mu_j(\vec{\mathbf{t}}).(\vec{\tau})$ , as required.  $\square$

In the case where  $S = \emptyset$ , this lemma says that for each  $j \in \{1, \dots, n\}$  we can derive  $\Gamma; \emptyset \vdash \delta_j \equiv \mu_j(\vec{\mathbf{t}}).(\vec{\tau})$ , which is just what we needed in order to show the matching that succeeded under the opaque interpretation would succeed under the transparent one.

## 4 A Coercion Interpretation of Datatypes

In this section, we will discuss a treatment of datatypes based on a variant of Cray’s Coercion Calculus [3]. This solution will closely resemble the Harper-Stone interpretation, and thus will not require the boxed-unroll rule or a trail algorithm, but will not incur the run-time cost of a function call at constructor application sites.

### 4.1 Representation of Datatype Values

The calculus we have discussed in this paper can be given the usual structured operational semantics, in which an expression of the form  $\mathbf{roll}_\delta(v)$  is itself a value if  $v$  is a value. (From here on we will implicitly assume that the metavariable  $v$  ranges only over values.) In fact, it can be shown without difficulty that *any* closed value of a datatype  $\delta$  must have the form  $\mathbf{roll}_\delta(v)$  where  $v$  is a closed value of type  $\mathbf{expand}(\delta)$ . Thus the  $\mathbf{roll}$  operator plays a similar role to that of the  $\mathbf{inj}_1$  operator for sum types, as far as the high-level language semantics is concerned.

Although we specify the behavior of programs in our language with a formal operational semantics, it is our intent that programs be compiled into machine code for execution, which forces us to take a slightly different view of data. Rather than working directly with high-level language values, compiled programs manipulate *representations* of those values. A compiler is free to choose the representation scheme it uses, provided that the basic operations of the language can be faithfully performed on representations. For example, most compilers construct the value  $\mathbf{inj}_1(v)$  by attaching a tag to the value  $v$  and storing this new object somewhere. This tagging is necessary in order to implement the **case** construct. In particular, the representation of any value of type  $\tau_1 + \tau_2$  must carry enough information to determine whether it was created with  $\mathbf{inj}_1$  or  $\mathbf{inj}_2$  and recover a representation of the injected value.

What are the requirements for representations of values of recursive type? It turns out that they are somewhat weaker than for sums. The elimination form for recursive types is  $\mathbf{unroll}$ , which (unlike **case**) does not need to extract any information from its argument other than the original rolled value. In fact, the only requirement is that a representation of  $v$  can be extracted from any representation of  $\mathbf{roll}_\delta(v)$ . Thus one reasonable representation strategy is to represent  $\mathbf{roll}_\delta(v)$  exactly the same as  $v$ . Appendix A gives a more precise argument as to why this is reasonable, making use of two key insights. First, it is an invariant of TILT that the representation of any value fits in a single machine register; anything larger than 32 bits is *always* stored in the heap and referred to by a pointer. This means that all possible complications having to do with the sizes of recursive values are avoided. Second, we define representations for values, not types; that is, we define the set of machine words that can represent the value  $v$  by structural induction on  $v$ , rather than defining the set of words that can represent values of type  $\tau$  by induction on  $\tau$  as might be expected.

The TILT compiler adopts this strategy of identifying the representations of  $\mathbf{roll} v$  and  $v$ , which has the pleasant consequence that the  $\mathbf{roll}$  and  $\mathbf{unroll}$  operations are “no-ops”. For instance, the untyped machine code generated by the compiler for the expression  $\mathbf{roll}_\delta(e)$  need not differ from the code for  $e$  alone, since if the latter evaluates to  $v$  then the former evaluates to  $\mathbf{roll}_\delta(v)$ , and the representations of these two values are the same. The reverse happens for  $\mathbf{unroll}$ .

This, in turn, has an important consequence for datatypes. Since the **in** and **out** functions produced by the Harper-Stone elaboration of datatypes do nothing but roll or unroll their arguments, the code generated for any **in** or **out** function will be the same as that of the identity function. Hence, the only run-time cost incurred by using an **in** function to construct a datatype value is the overhead of the function call itself. In the remainder of this section we will explain how to

eliminate this cost by allowing the types of the `in` and `out` functions to reflect the fact that their implementations are trivial.

## 4.2 The Coercion Interpretation

To mark `in` and `out` functions as run-time no-ops, we will use a simplified form of Crary’s calculus of *coercions* [3]. Coercions are similar to functions, except that they are known to be no-ops and therefore no code needs to be generated for coercion applications. Crary’s calculus has a separate syntactic class of coercion terms, but this level of sophistication is unnecessary for our purposes. Instead, we incorporate coercions into the term level of our language and introduce special coercion types to which they belong. Figure 7 gives the changes to the syntax of our calculus. Note that while we have so far confined our discussion to monomorphic datatypes, the general case of polymorphic datatypes will require polymorphic coercions. The syntax we give here is essentially that used in the TILT compiler; it does not handle non-uniform datatypes.

$$\begin{array}{l} \text{Types} \quad \sigma, \tau \quad ::= \quad \dots \mid (\vec{\alpha}; \tau_1) \Rightarrow \tau_2 \\ \text{Terms} \quad e \quad ::= \quad \dots \mid \Lambda \vec{\alpha}. \text{fold}_\delta \mid \Lambda \vec{\alpha}. \text{unfold}_\delta \mid v@(\vec{\tau}; e) \end{array}$$

Figure 7: Syntax of Coercions

We extend the type level of the language with a type for (possibly polymorphic) coercions,  $(\vec{\alpha}; \tau_1) \Rightarrow \tau_2$ ; a value of this type is a coercion that takes  $\text{length}(\vec{\alpha})$  type arguments and then can change a value of type  $\tau_1$  into one of type  $\tau_2$  (where, of course, variables from  $\vec{\alpha}$  can appear in either of these types). Similarly, we extend the term level with the (possibly polymorphic) coercion values  $\Lambda \vec{\alpha}. \text{fold}_\delta$  and  $\Lambda \vec{\alpha}. \text{unfold}_\delta$ ; these take the place of `roll` and `unroll` expressions. Coercions are applied to (type and value) arguments in an expression of the form  $v@(\vec{\tau}; e)$ ; here  $v$  is the coercion,  $\vec{\tau}$  are the type arguments, and  $e$  is the value to be coerced. Note that the coercion is syntactically restricted to be a value; this makes the calculus more amenable to a simple code generation strategy, as we will discuss in Section 4.3. The typing rules for coercions are essentially the same as if they were ordinary polymorphic functions, and are shown in Figure 8.

$$\begin{array}{c} \frac{\Gamma, \vec{\alpha} \vdash \tau_1 \text{ type} \quad \Gamma, \vec{\alpha} \vdash \tau_2 \text{ type}}{\Gamma \vdash (\vec{\alpha}; \tau_1) \Rightarrow \tau_2 \text{ type}} \\ \\ \frac{}{\Gamma \vdash \Lambda \vec{\alpha}. \text{fold}_\delta : (\vec{\alpha}; \text{expand}(\delta)) \Rightarrow \delta} \quad \frac{}{\Gamma \vdash \Lambda \vec{\alpha}. \text{unfold}_\delta : (\vec{\alpha}; \delta) \Rightarrow \text{expand}(\delta)} \\ \\ \frac{\Gamma \vdash v : (\vec{\alpha}; \tau_1) \Rightarrow \tau_2 \quad \Gamma \vdash e : \tau_1[\vec{\sigma}/\vec{\alpha}] \quad \forall i \in 1..n. \Gamma \vdash \sigma_i \text{ type}}{\Gamma \vdash v@(\vec{\sigma}; e) : \tau_2[\vec{\sigma}/\vec{\alpha}]} \end{array}$$

Figure 8: Typing Rules for Coercions

With these modifications to the language in place, we can elaborate the datatypes `exp` and `dec` using coercions instead of functions to implement the `in` and `out` operations. The result of elaborating this pair of datatypes is shown in Figure 9. Note that the interface is exactly the same as the Harper-Stone interface shown in Section 2 except that the function arrows ( $\rightarrow$ ) have been replaced by coercion arrows ( $\Rightarrow$ ). This interface is implemented by defining `exp` and `dec` in the same way as in the Harper-Stone interpretation, and implementing the `in` and `out` coercions as the

appropriate `fold` and `unfold` values. The elaboration of a constructor application is superficially similar to the opaque interpretation, but a coercion application is generated instead of a function call. For instance, `LetExp(d,e)` elaborates as `exp_in@(inj2(d,e))`.

```

structure ExpDec :> sig
  type exp
  type dec
  val exp_in : var + (dec * exp) ⇒ exp
  val exp_out : exp ⇒ var + (dec * exp)
  val dec_in : (var * exp) + (dec * dec) ⇒ dec
  val dec_out : dec ⇒ (var * exp) + (dec * dec)
end = struct
  type exp = μ1(α,β).(var + β * α, var * α + β * β)
  type dec = μ2(α,β).(var + β * α, var * α + β * β)
  val exp_in = foldexp
  val exp_out = unfoldexp
  val dec_in = folddec
  val dec_out = unfolddec
end

```

Figure 9: Elaboration of `exp` and `dec` Under the Coercion Interpretation

### 4.3 Coercion Erasure

We are now ready to formally justify our claim that coercions may be implemented by *erasure*, that is, that it is sound for a compiler to consider coercions only as “retyping operators” and ignore them when generating code. First, we will describe the operational semantics of the coercion constructs we have added to our internal language. Next, we will give a translation from our calculus into an untyped one in which coercion applications disappear. Finally, we will state a theorem guaranteeing that the translation is safe.

The operational semantics of our coercion constructs are shown in Figure 10. We extend the class of values with the `fold` and `unfold` coercions, as well as the application of a `fold` coercion to a value. These are the canonical forms of coercion types and recursive types respectively. The two inference rules shown in Figure 10 define the manner in which coercion applications are evaluated. The evaluation of a coercion application is similar to the evaluation of a normal function application where the applicand is already a value. The rule on the left specifies that the argument is reduced until it is a value. If the applicand is a `fold`, then the application itself is a value. If the applicand is an `unfold`, then the argument must have a recursive type and therefore (by canonical forms) consist of a `fold` applied to a value  $v$ . The rule on the right defines `unfold` to be the left inverse of `fold`, and hence this evaluates to  $v$ .

As we have already discussed, the data representation strategy of TILT is such that no code needs to be generated to compute `fold v` from  $v$ , nor to compute the result of cancelling a `fold` with an `unfold`. Thus it seems intuitive that to generate code for a coercion application  $v@(\vec{\tau}; e)$ , the compiler can simply generate code for  $e$ , with the result that datatype constructors and destructors under the coercion interpretation have the same run-time costs as Harper and Stone’s functions would if they were inlined. To make this more precise, we now define an *erasure* mapping to translate terms of our typed internal language into an untyped language with no coercion application.



$$v ::= \dots \mid \Lambda \vec{\alpha}.\mathbf{fold}_\tau \mid \Lambda \vec{\alpha}.\mathbf{unfold}_\tau \mid (\Lambda \vec{\alpha}.\mathbf{fold}_\tau)@(\vec{\sigma}; v)$$

$$\frac{e \mapsto e'}{v@(\vec{\tau}; e) \mapsto v@(\vec{\tau}; e')} \quad \frac{}{(\Lambda \vec{\alpha}.\mathbf{unfold}_{\tau_1})@(\vec{\sigma}; ((\Lambda \vec{\beta}.\mathbf{fold}_{\tau_2})@(\vec{\sigma}'; v))) \mapsto v}$$

Figure 10: Operational Semantics for Coercions.

The untyped nature of the target language (and of machine language) is important: treating  $v$  as  $\mathbf{fold} v$  would destroy the subject reduction property of a typed language.

$$M ::= \dots \mid \lambda x.M \mid \mathbf{fold} \mid \mathbf{unfold}$$

$$\begin{aligned} x^- &= x \\ (\lambda x:\tau.e)^- &= \lambda x.e^- \\ (\Lambda \vec{\alpha}.\mathbf{fold}_\delta)^- &= \mathbf{fold} \\ (\Lambda \vec{\alpha}.\mathbf{unfold}_\delta)^- &= \mathbf{unfold} \\ (v@(\vec{\tau}; e))^- &= e^- \\ &\vdots \end{aligned}$$

Figure 11: Target Language Syntax; Type and Coercion Erasure

Figure 11 gives the syntax of our untyped target language and the coercion-erasing translation. The target language is intended to be essentially the same as our typed internal language, except that all types and coercion applications have been removed. It contains untyped coercion values  $\mathbf{fold}$  and  $\mathbf{unfold}$ , but no coercion application form. The erasure translation turns expressions with type annotations into expressions without them ( $\lambda$ -abstraction and coercion values are shown in the figure), and removes coercion applications so that the erasure of  $v@(\vec{\tau}; e)$  is just the erasure of  $e$ . In particular, for any value  $v$ ,  $v$  and  $\mathbf{fold} v$  are identified by the translation, which is consistent with our intuition about the compiler. The operational semantics of the target language is analogous to that of the source.

The language with coercions has the important type-safety property that if a term is well-typed, its evaluation does not get stuck. An important theorem is that the coercion-erasing translation preserves the safety of well-typed programs.

**Theorem 1 (Erasure Preserves Safety)** *If  $\Gamma \vdash e : \tau$ , then  $e^-$  is safe. That is, if  $e^- \mapsto^* f$ , then  $f$  is not stuck.*

**Proof:** See Appendix B. □

Note that the value restriction on coercions is crucial to the soundness of this “coercion erasure” interpretation. Since a divergent expression can be given an arbitrary type, any semantics in which a coercion expression is not evaluated before it is applied fails to be type-safe. Thus if arbitrary expressions of coercion type could appear in application positions, the compiler would have to generate code for them. Since values cannot diverge or have effects, we are free to ignore coercion applications when we generate code.

## 5 Conclusion

### 5.1 Related Work

Our trail algorithm for weakened recursive type equivalence is based on the one we understand to have been implemented by Shao for the Standard ML of New Jersey compiler [10, 9]. SML/NJ handles datatypes using a strategy similar to our transparent interpretation, using the relaxed equivalence algorithm to recover the expressiveness of SML. We have discovered, however, that SML/NJ is not quite successful in avoiding the type sharing pitfalls of the transparent interpretation: the final example of Section 3.3, which is valid SML, is rejected by the SML/NJ compiler.

Crary [3] and Curien and Ghelli [4] have defined languages that use coercions to replace substitution rules in languages with subtyping. Crary’s calculus of coercions includes `roll` and `unroll` for recursive types, but since the focus of his paper is on subtyping he does not explore the potential uses of these coercions in detail. Nevertheless, our notion of coercion erasure, and the proof of our safety preservation theorem, are based on Crary’s. The implementation of Typed Assembly Language for the x86 architecture (TALx86) [7] allows operands to be annotated with coercions that change their types but not their representations; these coercions include `roll` and `unroll` as well as introduction of sums and elimination of universal quantifiers.

Our intermediate language differs from these in that we include coercions in the term level of the language rather than treating them specially in the syntax. This simplifies the presentation of the coercion interpretation of datatypes, and it simplified our implementation because it required a smaller incremental change from earlier pre-release versions of the TILT compiler. However, including coercions in the term level is a bit unnatural, and our planned extension of TILT with a type-preserving back-end will likely involve a full-blown coercion calculus.

### 5.2 Conclusion

The generative nature of SML datatypes poses a significant challenge for efficient type-preserving compilation. Generativity can be correctly understood by interpreting datatypes as structures that hold their type components abstract, exporting functions that construct and deconstruct datatype values. Under this interpretation, the inlining of datatype construction and deconstruction operations is not type-preserving and hence cannot be performed by a typed compiler such as TILT.

In this paper, we have discussed two approaches to eliminating the function call overhead in a type-preserving way. The first, doing away with generativity by making the type components of datatype structures transparent, results in a new language that is different from, but neither more nor less permissive than, Standard ML. Some of the lost expressiveness can be regained by relaxing the rules of type equivalence in the intermediate language, at the expense of complicating the type theory. The fact that the transparent interpretation forbids datatypes to appear in `sharing type` or `where type` signature constraints is unfortunate; it is possible that a revision of the semantics of these constructs could remove the restriction.

The second approach, replacing the construction and deconstruction functions of datatypes with coercions that may be erased during code generation, eliminates the function call overhead without changing the static semantics of the external language. However, the erasure of coercions only makes sense in a setting where a recursive-type value and its unrolling are represented the same at run time. The coercion interpretation of datatypes has been implemented in the TILT compiler, and a detailed performance analysis of the two approaches will appear in a future revision of this paper.

## References

- [1] Roberto Amadio and Luca Cardelli. Subtyping recursive types. *ACM Transactions on Programming Languages and Systems*, 15(4):575–631, 1993.
- [2] Michael Brandt and Fritz Henglein. Coinductive axiomatization of recursive type equality and subtyping. *Fundamenta Informaticae*, 33:309–338, 1998. Invited submission to special issue featuring a selection of contributions to the 3d Int’l Conf. on Typed Lambda Calculi and Applications (TLCA), 1997.
- [3] Karl Crary. Typed compilation of inclusive subtyping. In *2000 ACM International Conference on Functional Programming*, Montreal, September 2000.
- [4] Pierre-Louis Curien and Giorgio Ghelli. Coherence of subsumption, minimum typing and type-checking in  $F_{\leq}$ . *Mathematical Structures in Computer Science*, 2(1):55–91, 1992.
- [5] Robert Harper and Chris Stone. A type-theoretic interpretation of Standard ML. In Gordon Plotkin, Colin Stirling, and Mads Tofte, editors, *Robin Milner Festschrift*. MIT Press, 1998.
- [6] Robin Milner, Mads Tofte, Robert Harper, and David MacQueen. *The Definition of Standard ML (Revised)*. MIT Press, Cambridge, Massachusetts, 1997.
- [7] Greg Morrisett, Karl Crary, Neal Glew, Dan Grossman, Richard Samuels, Frederick Smith, David Walker, Stephanie Weirich, and Steve Zdancewic. TALx86: A realistic typed assembly language. In *Second Workshop on Compiler Support for System Software*, pages 25–35, Atlanta, Georgia, May 1999.
- [8] Leaf Petersen, Perry Cheng, Robert Harper, and Chris Stone. Implementing the TILT internal language. Technical Report CMU-CS-00-180, School of Computer Science, Carnegie Mellon University, December 2000.
- [9] Zhong Shao. Personal communication.
- [10] Zhong Shao. An overview of the FLINT/ML compiler. In *1997 Workshop on Types in Compilation*, Amsterdam, June 1997. ACM SIGPLAN. Published as Boston College Computer Science Department Technical Report BCCS-97-03.
- [11] David Tarditi, Greg Morrisett, Perry Cheng, Chris Stone, Robert Harper, and Peter Lee. TIL: A type-directed optimizing compiler for ML. In *ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 181–192, Philadelphia, PA, May 1996.

## A Data Representation

In this appendix we will give a formal description of a data representation strategy similar to the one used by the TILT compiler. This is intended to clarify the sense in which the `roll` and `unroll` operations on recursive types are “no-ops” at run time.

A key invariant of the TILT system is that every value (except floating-point numbers) manipulated by a program at run-time is 32 bits wide. All values of record and sum types that require more than 32 bits of storage are *boxed*—*i.e.*, stored in the heap and referred to by a 32-bit pointer value. Our formalization of data representation will therefore attempt to characterize when, in the context of a particular heap, a particular integer represents a particular source-language value. Once we

have done this, we will argue that the representation strategy we have defined is *reasonable*—that is, that all the operations a program must perform on values may be faithfully performed on their representations under this strategy.

To begin, we make the simplifying assumption that any integer may be a valid memory address, and that a memory location can hold any integer. This can easily be restricted so that only integers between 0 and, say,  $2^{32}$  are used, but allowing arbitrary integers makes the presentation easier. Under these assumptions about the world, it makes sense to define heaps as follows:

**Definition 1** *A heap is a finite partial function  $H : \mathbb{N} \rightarrow \mathbb{N}$ .*

Next, we define the general notion of a representation strategy.

**Definition 2** *A representation strategy is a four-place relation  $S$  on heaps, natural numbers, closed values and closed types such that if  $(H, n, v, \tau) \in S$  and  $H \subseteq H'$ , then  $(H', n, v, \tau) \in S$ .*

If  $S$  is a representation strategy, we will use the notation  $H \vdash_S n \triangleleft v : \tau$  to mean that  $(H, n, v, \tau) \in S$ , omitting the subscript  $S$  if it is clear from context. That statement may be read as “in the heap  $H$ , the number  $n$  represents the value  $v$  at type  $\tau$ .”

We will now proceed to define a particular representation strategy, similar to the one used by the TILT compiler. Figure 12 gives the syntax of the types and terms we will be representing. The types include the type of integers,  $k$ -ary tuple types  $\langle \tau_1, \dots, \tau_k \rangle$ ,  $k$ -ary sum types  $[\tau_1, \dots, \tau_k]$  and (single) recursive types  $\mu\alpha.\tau$ . (We are not considering arrow types or function values here, because they complicate the presentation in ways not relevant to recursive types.)

$$\begin{array}{ll} \text{Types} & \sigma, \tau ::= \text{nat} \mid \langle \tau_1, \dots, \tau_k \rangle \mid [\tau_1, \dots, \tau_k] \mid \mu\alpha.\tau \\ \text{Values} & v ::= \bar{n} \mid (v_0, \dots, v_{m-1}) \mid \text{inj}_i v \mid \text{roll}_{\mu\alpha.\tau}(v) \end{array}$$

Figure 12: Syntax of types and values in the first-order fragment of a specialized intermediate language.

Figure 13 gives the definition of one possible representation strategy for these values. Note

$$\begin{array}{ll} H \vdash n \triangleleft \bar{n} : \text{nat} & \\ H \vdash n \triangleleft (v_0, \dots, v_{m-1}) : \langle \tau_0, \dots, \tau_{m-1} \rangle & \text{if } H \vdash H(n+i) \triangleleft v_i : \tau_i \text{ for each } i, 0 \leq i < m \\ H \vdash n \triangleleft \text{inj}_i v : [\tau_1, \dots, \tau_k] & \text{if } 1 \leq i < k \text{ and } H(n) = i \text{ and} \\ & H \vdash H(n+1) \triangleleft v : \tau_i \\ H \vdash n \triangleleft \text{roll}_{\mu\alpha.\tau} v : \mu\alpha.\tau & \text{if } H \vdash n \triangleleft v : \tau[\mu\alpha.\tau/\alpha]. \end{array}$$

Figure 13: Representation strategy for our intermediate language fragment.

that this strategy is well-defined, because all the values on the right-hand side of each clause are syntactically smaller than the one on the left. An integer value  $\bar{n}$  is represented by the integer  $n$ . A tuple of  $k$  values is represented by a pointer ( $n$ ) giving the location of the first of  $k$  consecutive heap cells containing representations of the tuple’s components. An injection into a sum type is represented by a pointer to what is essentially a two-element tuple: the first element is a tag that identifies a branch of the sum type, and the second is a representation of the injected value. Finally, a value of recursive type is represented by a representation of the rolled value.

This data representation strategy is similar to, but considerably simpler than, the one used by the TILT/ML compiler. The main difference is in the handling of sum types; TILT uses a more

complex version of sum type representations that saves space and improves run-time efficiency. Fortunately, the treatment of sums and the treatment of recursive types are orthogonal, so the differences between the representation strategy in Figure 13 and that of TILT is unimportant.

We can now argue that this data representation strategy is reasonable. The following property may be proved using the definition.

**Construction Property:** *Let  $H$  be a heap. If  $\emptyset \vdash v : \tau$ , then there is a heap  $H' \supseteq H$  and an integer  $n$  such that  $H' \vdash n \triangleleft v : \tau$ .*

The construction property states that any well-typed closed value can be represented in some extension of any initial heap. This roughly means that all the introduction forms of the language can be faithfully implemented under our representation strategy. It remains to argue that the elimination forms may be implemented as well.

The elimination forms for integers are arithmetic operations and comparisons. Since the representations of integer values are the integers themselves, these operations can be implemented. The elimination form for tuples is projection; we need to show that if  $1 \leq i \leq k$  and  $H \vdash n \triangleleft v : \langle \tau_0, \dots, \tau_{k-1} \rangle$  then we can find some  $n'$  that represents the value of  $\pi_i v$ . By canonical forms, the value  $v$  must be  $(v_0, \dots, v_{k-1})$ , and so  $\pi_i v$  evaluates to  $v_i$ . But by the definition of our representation strategy,  $H \vdash H(n+i) \triangleleft v_i : \tau_i$ , so projection can be implemented.

The elimination form for sums, the **case** construct, is a little different in that it may take many more than one step to produce a value, if it produces one at all. Only the first of these steps, however—the one in which the **case** chooses one of its subexpressions to continue executing and passes the appropriate value to that branch—is relevant to the reasonableness of our representation for sums. In particular, if  $\emptyset \vdash v : (\tau_1, \dots, \tau_k)$  then the expression

$$\text{case } v \text{ of } \mathbf{inj}_1 x \Rightarrow e_1 \mid \dots \mid \mathbf{inj}_k x \Rightarrow e_k$$

will certainly take a step to  $e_j[v'/x]$  for some  $j$  and some value  $v'$  of the appropriate type. In order to show **case** is implementable, it suffices to show that given  $H \vdash n \triangleleft v : [\tau_1, \dots, \tau_k]$  we can compute the appropriate  $i$  and a representation of the appropriate value  $v'$ . But note that by canonical forms  $v$  must have the form  $\mathbf{inj}_i v_i$ , in which case the branch to select is the  $i^{\text{th}}$  one, and the value to pass is  $v_i$  (i.e.,  $j = i$  and  $v' = v_i$ ). According to our representation strategy,  $H(n) = i$  and  $H \vdash H(n+1) \triangleleft v_i : \tau_i$ , so we can implement **case**.

Finally, consider the elimination form for recursive types,  $\mathbf{unroll}_{\mu\alpha.\tau}$ . In order to implement this operation, it must be the case that if  $\emptyset \vdash v : \mu\alpha.\tau$  and  $H \vdash n \triangleleft v : \mu\alpha.\tau$  then we can construct some  $n'$  such that  $n'$  represents the value of  $\mathbf{unroll}_{\mu\alpha.\tau} v$ . By canonical forms,  $v = \mathbf{roll}_{\mu\alpha.\tau} v'$  and  $\mathbf{unroll}_{\mu\alpha.\tau} v$  evaluates to  $v'$ . But under our representation strategy,  $H \vdash n \triangleleft v' : \tau[\mu\alpha.\tau/\alpha]$  and so we can implement **unroll**.

Notice also that a representation of a value  $v$  of recursive type also represents the value produced by  $\mathbf{unroll} v$ , just as was the case for **roll**. This justifies our notion that **roll** and **unroll** may be implemented as no-ops in a compiler that uses this data representation strategy.

## B Proof of Coercion Erasure

In this appendix we will prove the coercion erasure safety theorem stated in Section 4 for a simple fragment of our intermediate language. The theorem states that if an expression is well-typed, then evaluation of its erasure does get stuck; that is, it never reaches an expression that is not a value but for which no transition rule applies. To prove this, we will establish a correspondence between

the transition relation of the typed calculus and that of the untyped one; we will then be able to show that the type-preservation and progress lemmas for the typed calculus guarantee safety of their erasures.

The first two lemmas we will use are easy to prove using the definition of erasure given in Section 4.

**Lemma 2 (Erasure of Values)** *If  $v$  is a value, then  $v^-$  is a value.*

**Proof:** By structural induction on  $v$ . □

**Lemma 3 (Substitution and Erasure Commute)**

*For any expression  $e$  and value  $v$ ,  $(e[v/x])^- = e^-[v^-/x]$ .*

**Proof:** By structural induction on  $e$ . □

$$\frac{}{(\lambda x:\tau . e_1) v_2 \mapsto_e e_1[v_2/x]} \quad \frac{}{(\Lambda \vec{\alpha} . \text{unfold}_{\tau_1})@(\vec{\sigma}; ((\Lambda \vec{\beta} . \text{fold}_{\tau_2})@(\vec{\sigma}'; v))) \mapsto_c v}$$

$$\frac{e_1 \mapsto_\varphi e'_1}{e_1 e_2 \mapsto_\varphi e'_1 e_2} \quad \frac{e_2 \mapsto_\varphi e'_2}{v_1 e_2 \mapsto_\varphi v_1 e'_2} \quad \frac{e \mapsto_\varphi e'}{v@(\vec{\tau}; e) \mapsto_\varphi v@(\vec{\tau}; e')}$$

Figure 14: Annotated Operational Semantics for Typed Language

In order to properly characterize the relationship between the typed and untyped transition relations we must distinguish, for the typed calculus, between evaluation steps like  $\beta$ -reduction that correspond to steps in the untyped semantics and those like coercion application that do not. To this end, we annotate transitions of the typed calculus as shown in Figure 14. The flag  $\varphi$  adorning each transition may be either  $e$ , indicating an “evaluation” step that is preserved by erasure, or  $c$ , indicating (to use Crary’s [CITE] terminology) a “canonicalization” step that is removed by erasure. We will continue to use the unannotated  $\mapsto$  to stand for the union of  $\mapsto_e$  and  $\mapsto_c$ , and as usual we will use  $\mapsto^*$ ,  $\mapsto_e^*$  and  $\mapsto_c^*$  to denote the reflexive, transitive closures of these relations. With these definitions in place, we can prove the following lemma, which states that evaluation steps are preserved by erasure, but terms that differ by a canonicalization step are identified by erasure. It follows that any sequence of steps in the typed language erases to some sequence of steps in the untyped language.

**Lemma 4 (Simulation)**

1. *If  $e_1 \mapsto_e e_2$ , then  $e_1^- \mapsto e_2^-$ .*
2. *If  $e_1 \mapsto_c e_2$ , then  $e_1^- = e_2^-$ .*

**Proof:** By induction on derivations, using the definition of erasure and Lemmas 2 and 3. □

**Lemma 5**

1. *If  $e_1 \mapsto_c^* e_2$ , then  $e_1^- = e_2^-$ .*
2. *If  $e_1 \mapsto_e^* e_2$ , then  $e_1^- \mapsto^* e_2^-$ .*

3. If  $e_1 \mapsto^* e_2$ , then  $e_1^- \mapsto^* e_2^-$ .

**Proof:** Each part is proved separately by induction on the length of the transition sequence.  $\square$

We are now ready to prove the equivalent of a Progress lemma for the untyped language. It states that a term whose erasure is a value canonicalizes to a value in some number of steps, and a term whose erasure is not a value will eventually take an evaluation step.

**Lemma 6 (Progress under Erasure)** *If  $\emptyset \vdash e : \tau$  then either*

1.  $e^-$  is a value and  $e \mapsto_c^* v$  for some value  $v$ , OR

2.  $e^-$  is not a value and  $e \mapsto_c^* e' \mapsto_e e''$  for some  $e'$  and  $e''$ .

**Proof:** By induction on the typing derivation for  $e$ . Note that if  $e$  is a value, then so is  $e^-$  (by Lemma 2) and  $e \mapsto_c^* e$  by definition. Thus we need only consider the typing rules in which the expression being typed may be a non-value.

**Case:**

$$\frac{\emptyset \vdash e_1 : \tau' \rightarrow \tau \quad \emptyset \vdash e_2 : \tau'}{\emptyset \vdash e_1 e_2 : \tau}$$

Note that  $(e_1 e_2)^- = e_1^- e_2^-$ , which is not a value. Thus we must show that  $e_1 e_2 \mapsto_c^* e' \mapsto_e e''$ . There are three sub-cases:

**Sub-case:**  $e_1^-$  is not a value. By the induction hypothesis,  $e_1 \mapsto_c^* e'_1 \mapsto_e e''_1$ . Thus,  $e_1 e_2 \mapsto_c^* e'_1 e_2 \mapsto_e e''_1 e_2$ .

**Sub-case:**  $e_1^-$  is a value,  $e_2^-$  is not. By the induction hypothesis,  $e_2 \mapsto_c^* e'_2 \mapsto_e e''_2$ . Thus,  $e_1 e_2 \mapsto_c^* e_1 e_2 \mapsto_c^* v_1 e_2 \mapsto_c^* v_1 e'_2 \mapsto_e v_1 e''_2$ .

**Sub-case:**  $e_1^-$  and  $e_2^-$  are both values. By the induction hypothesis,  $e_1 \mapsto_c^* v_1$  and  $e_2 \mapsto_c^* v_2$ . By type preservation,  $\emptyset \vdash v_1 : \tau' \rightarrow \tau$ . By canonical forms,  $v_1 = \lambda x:\tau'.e_3$ . Thus,  $e_1 e_2 \mapsto_c^* (\lambda x:\tau'.e_3) e_2 \mapsto_c^* (\lambda x:\tau'.e_3) v_2 \mapsto_e e_3[v_2/x]$ .

**Case:**

$$\frac{\Gamma \vdash v_c : (\vec{\alpha}; \tau_1) \Rightarrow \tau_2 \quad \Gamma \vdash e_1 : \tau_1[\vec{\sigma}/\vec{\alpha}] \quad \forall i \in 1..n. \Gamma \vdash \sigma_i \text{ type}}{\Gamma \vdash v_c @(\vec{\sigma}; e_1) : \tau_2[\vec{\sigma}/\vec{\alpha}]}$$

Note that  $(v_c @(\vec{\sigma}; e_1))^- = e_1^-$ . If  $e_1^-$  is not a value, then by the induction hypothesis we get  $e_1 \mapsto_c^* e'_1 \mapsto_e e''_1$ , and hence  $v_c @(\vec{\sigma}; e_1) \mapsto_c^* v_c @(\vec{\sigma}; e'_1) \mapsto_e v_c @(\vec{\sigma}; e''_1)$  as required.

If  $e_1^-$  is a value, then we must show that  $v_c @(\vec{\sigma}; e_1) \mapsto_c^* v$  for some value  $v$ . By the induction hypothesis we have  $e_1 \mapsto_c^* v_1$ . There are two sub-cases.

**Sub-case:** The coercion  $v_c$  is  $\Lambda \vec{\alpha}.\text{fold}_\delta$ . Then  $v_c @(\vec{\sigma}; v_1)$  is a value and  $v_c @(\vec{\sigma}; e_1) \mapsto_c^* v_c @(\vec{\sigma}; v_1)$  as required.

**Sub-case:** The coercion  $v_c$  is  $\Lambda \vec{\alpha}.\text{unfold}_\delta$ . By inversion,  $\tau_1 \equiv \delta$ . By type preservation,  $\emptyset \vdash v_1 : \delta[\vec{\sigma}/\vec{\alpha}]$ . By canonical forms,  $v_1 = (\Lambda \vec{\beta}.\text{fold}_{\delta'}) @(\vec{\sigma}'; v'_1)$ . Thus,  $v_c @(\vec{\sigma}; e_1) \mapsto_c^* v_c @(\vec{\sigma}; v_1) = (\Lambda \vec{\alpha}.\text{unfold}_\delta) @(\vec{\sigma}; (\Lambda \vec{\beta}.\text{fold}_{\delta'}) @(\vec{\sigma}'; v'_1)) \mapsto_c v'_1$ .  $\square$

Next, we would like to prove an analogue of type preservation for our target calculus. Clearly it is meaningless to prove type preservation for an untyped calculus, so we must prove instead that if the erasure of a well-typed term takes a step, the result is itself the erasure of a well-typed term. Because type preservation does hold for the typed calculus, the following lemma is sufficient to show this.

**Lemma 7 (Preservation under Erasure)** *If  $\emptyset \vdash e : \tau$  and  $e^- \mapsto f$ , then there is some  $e'$  such that  $e \mapsto^* e'$  and  $(e')^- = f$ .*

**Proof:** By induction on the typing derivation for  $e$ . Note that since  $e^- \mapsto f$ ,  $e^-$  cannot be a value and hence neither can  $e$ . Thus, as for the progress lemma, we only need to consider typing rules in which the expression being typed may be a non-value.

**Case:**

$$\frac{\emptyset \vdash v_c : (\vec{\alpha}; \tau_1) \Rightarrow \tau_2 \quad \emptyset \vdash e_1 : \tau_1[\vec{\sigma}/\vec{\alpha}] \quad \forall i \in 1..n. \emptyset \vdash \sigma_i \text{ type}}{\emptyset \vdash v_c @(\vec{\sigma}; e_1) : \tau_2[\vec{\sigma}/\vec{\alpha}]}$$

Note that  $(v_c @(\vec{\sigma}; e_1))^- = e_1^-$ , so in fact  $e_1^- \mapsto f$ . By the induction hypothesis,  $e_1 \mapsto^* e'_1$  where  $(e'_1)^- = f$ . Thus  $v_c @(\vec{\sigma}; e_1) \mapsto^* v_c @(\vec{\sigma}; e'_1)$ , and  $(v_c @(\vec{\sigma}; e'_1))^- = (e'_1)^- = f$ .

**Case:**

$$\frac{\emptyset \vdash e_1 : \tau' \rightarrow \tau \quad \emptyset \vdash e_2 : \tau'}{\emptyset \vdash e_1 e_2 : \tau}$$

Note that  $(e_1 e_2)^- = e_1^- e_2^-$ , and so  $e_1^- e_2^- \mapsto f$ . There are three possibilities for the last rule used in the derivation of this transition.

**Sub-case:**

$$\frac{e_1^- \mapsto f_1}{e_1^- e_2^- \mapsto f_1 e_2^-} \quad (\text{where } f = f_2 e_2^-.)$$

By the induction hypothesis,  $e_1 \mapsto^* e'_1$  such that  $(e'_1)^- = f_1$ . Thus  $e_1 e_2 \mapsto e'_1 e_2$  and  $(e'_1 e_2)^- = (e'_1)^- e_2^- = f_1 e_2^- = f$ .

**Sub-case:**

$$\frac{e_2^- \mapsto f_2}{w_1 e_2^- \mapsto w_1 f_2} \quad (\text{where } e_1^- = w_1 \text{ is a value, and } f = w_1 f_2)$$

By Lemma 6,  $e_1 \mapsto^* v_1$ . By Lemma 5,  $v_1^- = w_1$ . By the induction hypothesis,  $e_2 \mapsto^* e'_2$  such that  $(e'_2)^- = f_2$ . Thus,  $e_1 e_2 \mapsto^* v_1 e_2 \mapsto^* v_1 e'_2$ , and  $(v_1 e'_2)^- = v_1^- (e'_2)^- = w_1 f_2 = f$ .

**Sub-case:**

$$\overline{(\lambda x.f_3) w_2 \mapsto f_3[w_2/x]} \quad (\text{where } e_1^- = \lambda x.f_3, e_2^- = w_2 \text{ is a value, and } f = f_3[w_2/x].)$$

By Lemma 6,  $e_1 \mapsto^* v_1$  and  $e_2 \mapsto^* v_2$ . By Lemma 5,  $v_1^- = \lambda x.f_3$  and  $v_2^- = w_2$ . By type preservation,  $\emptyset \vdash v_1 : \tau' \rightarrow \tau$ , and so by canonical forms  $v_1 = \lambda x : \tau'. e_3$ . But this means that  $v_1^- = \lambda x.e_3^-$ , so it must be the case that  $e_3^- = f_3$ . Thus we have

$$e_1 e_2 \mapsto^* (\lambda x:\tau'. e_3) e_2 \mapsto^* (\lambda x:\tau'. e_3) v_2 \mapsto^* e_3[v_2/x]$$

and by Lemma 3,  $(e_3[v_2/x])^- = e_3^-[v_2^-/x] = f_3[w_2/x] = f$ . □



We can now extend this lemma to transition sequences of any length, and prove the safety theorem. This last theorem effectively states that the erasure of a well-typed expression can never get stuck.

**Lemma 8** *If  $\emptyset \vdash e : \tau$  and  $e^- \mapsto^* f$ , then there is some  $e'$  such that  $e \mapsto^* e'$  and  $(e')^- = f$ .*

**Proof:** By induction on the length of the transition sequence. □

**Theorem 2 (Erasure Preserves Safety)** *If  $\emptyset \vdash e : \tau$  and  $e^- \mapsto^* f$ , then either  $f$  is a value or there exists an  $f'$  such that  $f \mapsto f'$ .*

**Proof:** By Lemma 8, there is an  $e'$  such that  $e \mapsto^* e'$  and  $(e')^- = f$ . By type preservation,  $\emptyset \vdash e' : \tau$ . Suppose  $f$  is not a value; then by Lemma 6 there are  $e''$  and  $e'''$  such that  $e' \mapsto_c^* e'' \mapsto_e e'''$ . By Lemma 5,  $(e')^- = (e'')^-$  and  $(e'')^- \mapsto (e''')^-$ . Therefore,  $f \mapsto (e''')^-$ . □