

Typed Compilation of Recursive Datatypes*

Joseph C. Vanderwaart

Karl Crary

Derek Dreyer

Robert Harper

Leaf Petersen

Perry Cheng[†]

School of Computer Science
Carnegie Mellon University
Pittsburgh, PA 15213

Abstract

Standard ML employs an opaque (or generative) semantics of datatypes, in which every datatype declaration 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 datatypes to be abstract. 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, altering the semantics of SML; and second, we propose an interpretation of datatype constructors as coercions, which have no run-time effect or cost and faithfully implement SML semantics.

Categories and Subject Descriptors

D.3.3 [Programming Languages]: Language Constructs and Features—*Abstract data types*; D.3.4 [Programming Languages]: Processors—*Compilers*; F.3.3 [Logics and Meanings of Programs]: Studies of Program Constructs—*Type structure*

General Terms

Languages, Theory, Performance

Keywords

Typed compilation, Standard ML, recursive types, coercions

*The ConCert Project is supported by the National Science Foundation under grant number 0121633: "ITR/SY+SI: Language Technology for Trustless Software Dissemination".

[†]IBM, TJ Watson, P.O. Box 704, Yorktown, NY 10598

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

TLDI'03, January 18, 2003, New Orleans, Louisiana, USA.

Copyright 2003 ACM 1-58113-649-8/03/0001 ...\$5.00

1 Introduction

The programming language Standard ML (SML) [9] 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 [9] 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 [7] (hereafter, HS) give 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, we consider HS to be a suitable foundation for type-directed compilation of SML. Furthermore, it seems likely that any other suitable type-theoretic interpretation (*i.e.*, one that is faithful to the Definition) will encounter the same issues we explore in our analysis.

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 HS interpretation, datatypes *are* abstract types. The generativity of datatypes poses some challenges for type-directed compilation of SML. In particular, although the HS interpretation is easy to understand and faithful to the Definition of SML, it is inefficient when implemented naively. 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. We will present these alternatives in the context of the TILT/ML compiler developed at CMU [11, 14]; they are relevant, however, not just to TILT, but to understanding the definition of the language and type-preserving compilation in general.

The first approach is to do away with datatype generativity altogether, replacing the abstract types in the HS 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. In particular, it is *not* the case that every well-formed SML program is allowed under the transparent interpretation.

In contrast, the second approach, which we have adopted in the most recent version of the TILT compiler, offers an efficient way of implementing datatypes in a typed setting that is consistent with the Definition. In particular, since a value of recursive type is typically represented at run time in the same way as its unrolling, we can observe that the mediating functions produced by the HS interpretation all behave like the identity function at run time. We replace these functions with special values that are distinguished from ordinary functions by the introduction of “coercion types”. We call this the *coercion interpretation of datatypes*, and argue that it allows a compilation strategy that generates code with a run-time efficiency comparable to what would be attained if datatype constructors were inlined.

The paper is structured as follows: Section 2 gives the details of the HS interpretation of datatypes (which we also refer to as the *opaque interpretation of datatypes*) and illustrates the problems with inlining. Section 3 discusses the transparent interpretation. Section 4 gives the coercion interpretation and discusses its properties. Section 5 gives a performance comparison of the three interpretations. Section 6 discusses related work and Section 7 concludes.

2 The Opaque Interpretation of Datatypes

In this section, we 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 HS 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.

<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>	$\Gamma ::= \varepsilon \mid \Gamma, x : \tau \mid \Gamma, \alpha$

Figure 1. Syntax of Iso-recursive Types

\bar{X}	$\stackrel{\text{def}}{=} X_1, \dots, X_n \text{ for some } n \geq 1,$ where X is a metavariable, such as α or τ
$\text{length}(\bar{X})$	$\stackrel{\text{def}}{=} n, \text{ where } \bar{X} = X_1, \dots, X_n$
$\mu\alpha. \tau$	$\stackrel{\text{def}}{=} \mu_1(\alpha).(\tau)$
$\bar{\mu}(\bar{\alpha}).(\bar{\tau})$	$\stackrel{\text{def}}{=} \mu_1(\bar{\alpha}).(\bar{\tau}), \dots, \mu_n(\bar{\alpha}).(\bar{\tau}),$ where $\text{length}(\bar{\alpha}) = \text{length}(\bar{\tau}) = n$
$\text{expand}(\delta)$	$\stackrel{\text{def}}{=} \tau_i[\bar{\mu}(\bar{\alpha}).(\bar{\tau})/\bar{\alpha}], \text{ where } \delta = \mu_i(\bar{\alpha}).(\bar{\tau})$

Figure 2. Shorthand Definitions

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 [7]. 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 own syntactic subcategory, ranged over by δ . 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 α_j is a type variable that may appear free in any or all of τ_1, \dots, τ_n . Intuitively, this type is the i th in a system of n mutually recursive types. As such, it is isomorphic to τ_i with each α_j replaced by the j th component of the recursive bundle. Formally, it is isomorphic to the following somewhat unwieldy type:

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

(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 τ). 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 $\text{expand}(\mu_i(\bar{\alpha}).(\bar{\tau}))$.

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 “un-

$\Gamma \vdash \text{ok}$	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 Judgments

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

Figure 4. Typing Rules for Iso-recursive Types

rolling” 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 [5]. 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 `expand`(δ) into one of type δ , and the latter is its inverse.

2.2 Elaborating Datatype Declarations

The HS 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 declaration therefore elaborates to 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 HS elaboration of this declaration 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 `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:

```

VarExp(x)  $\stackrel{\text{def}}{=}$  ExpDec.exp_in(inj_1(x))
LetExp(d,e)  $\stackrel{\text{def}}{=}$  ExpDec.exp_in(inj_2(d,e))
ValDec(x,e)  $\stackrel{\text{def}}{=}$  ExpDec.dec_in(inj_1(x,e))
SeqDec(d1,d2)  $\stackrel{\text{def}}{=}$  ExpDec.dec_in(inj_2(d1,d2))

```

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

```

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 = roll_exp(x)
  fun exp_out x = unroll_exp(x)
  fun dec_in x = roll_dec(x)
  fun dec_out x = unroll_dec(x)
end

```

Figure 5. Harper-Stone Elaboration of `exp-dec` Example

out functions: for example, if we attempt to inline `exp_in` in the definition of `VarExp` above, we get

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

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.

One strategy that can alleviate this somewhat is to hold the implementation of a datatype abstract *during elaboration*, but to expose its underlying implementation *after elaboration* to other code defined in the same compilation unit. Calls to the constructors of a locally-defined datatype can then be safely inlined. In the setting of whole-program compilation, this approach can potentially eliminate constructor call overhead for all datatypes except those appearing as arguments to functors. However, in the context of separate compilation, the clients of a datatype generally do not have access to its implementation, but rather only to the specifications of its constructors. As we shall see in Section 3, the specifications of a datatype’s constructors do not provide sufficient information to correctly predict how the datatype is actually implemented, so the above compilation strategy will have only limited success in a true separate compilation setting.

2.3 Datatypes and Signature Matching

Standard ML makes an important distinction between datatype *declarations*, which appear at the top level or in structures, and datatype *specifications*, which appear in signatures. As we have seen, the HS interpretation elaborates datatype declarations 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.¹ In particular, it is clear that the structure definition produced by the HS interpretation for the datatype `intlist` defined in Section 1 has this signature, so that datatype declaration 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 opaquely sealed structures, and datatype specifications are translated as structure specifications, matching a datatype declaration against a spec boils down to matching one signature—the one opaquely sealing the declaration structure—against another signature.

Suppose we wish to know whether the signature S matches the signature T ; that is, whether a structure with signature S may also be given the signature T . Intuitively, we must make sure that for every specification in T there is a specification in S that is compatible with it. For instance, if T contains a value specification of the form `val x : τ` , then S must also contain a specification `val x : τ'` , where $\tau' \equiv \tau$. For an abstract type specification of the form `type t` occurring in T , we must check that a specification of t also appears in S ; furthermore, if the specification in S is a transparent one, say `type t = τ_{imp}` , then when checking the remainder of the specifications in T we may assume in both signatures that $t = \tau_{imp}$. Transparent type specifications in T are similar, but there is the added requirement that if the specification in T is `type t = τ_{spec}` and the specification in S is `type t = τ_{imp}` , then τ_{spec} and τ_{imp} must be equivalent.

Returning to the specific question of datatype matching, a specification of the form

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

(where the τ_i may be sum types) elaborates to a specification of a structure with the following signature:

```
sig
  type t1
  :
  type tn
  val t1_in :  $\tau_1$  -> t1
  val t1_out : t1 ->  $\tau_1$ 
  :
  val tn_in :  $\tau_n$  -> tn
  val tn_out : tn ->  $\tau_n$ 
end
```

¹Standard ML allows only datatypes to match datatype specifications, so the actual HS elaboration must use a name for the datatype that cannot be guessed by a programmer.

```
structure ExpDec :> sig
  type exp =  $\mu_1(\alpha, \beta).(\text{var} + \beta * \alpha, \text{var} * \alpha + \beta * \beta)$ 
  type dec =  $\mu_2(\alpha, \beta).(\text{var} + \beta * \alpha, \text{var} * \alpha + \beta * \beta)$ 
  (* ... specifications for in and out functions
     same as before ... *)
end =
  (* ... same structure as before ... *)
```

Figure 6. The Transparent Elaboration of Exp and Dec

In order to match this signature, the structure corresponding to a datatype declaration 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 declaration under this interpretation, the t_i 's will be abstract types.) Thus, for example, if $m \geq n$ then the datatype declaration

datatype $t_1 = \sigma_1$ and ... and $t_m = \sigma_m$

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 HS 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 that shown in Figure 6 (we continue to use ML-like syntax for readability).

Importantly, this change must extend to datatype specifications as well as datatype declarations. Thus, a structure that exports a datatype must export its implementation transparently, using a signature similar to the one in the figure—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 7. In the figure, two datatypes are defined by seemingly identical declarations. 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

```

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 7. Non-generativity of Transparent Datatypes

obvious is that there are some programs for which the opposite is true. We will discuss two main reasons for this.

3.2 Problematic Datatype Matchings

Recall that according to the HS 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. (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 in the structure.) 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 $u_{\text{imp}} \stackrel{\text{def}}{=} \mu\alpha. \alpha * \alpha + \text{int}$ 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 $u_{\text{imp}} * u_{\text{imp}}$ for v in the specification. That is, it is required that

$$u_{\text{imp}} \equiv \mu\alpha. u_{\text{imp}} * u_{\text{imp}} + \text{int}$$

Observe that the type on the right is none other than $\mu\alpha. \text{expand}(u_{\text{imp}})$. (Notice also that the bound variable α does not

appear free in the body of this μ -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 μ .) 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 IL signature looks like:

```

sig
  type exp
  type dec = decspec
  :
end

```

where $\text{dec}_{\text{spec}} \stackrel{\text{def}}{=} \mu\alpha. \text{var} * \text{exp} + \alpha * \alpha$. Elaboration of the declarations of `exp` and `dec`, on the other hand, produces the structure in Figure 6, which has the signature:

```

sig
  type exp = expimp
  type dec = decimp
  :
end

```

where we define

$$\begin{aligned} \text{exp}_{\text{imp}} &\stackrel{\text{def}}{=} \mu_1(\alpha, \beta). (\text{var} + \beta * \alpha, \text{var} * \alpha + \beta * \beta) \\ \text{dec}_{\text{imp}} &\stackrel{\text{def}}{=} \mu_2(\alpha, \beta). (\text{var} + \beta * \alpha, \text{var} * \alpha + \beta * \beta) \end{aligned}$$

Matching the structure containing the datatypes against the signature can only succeed if $\text{dec}_{\text{spec}} \equiv \text{dec}_{\text{imp}}$ (under the assumption that $\text{exp} \equiv \text{exp}_{\text{imp}}$). This equivalence does not hold because the two μ -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
    type s
    datatype t = A | B of s
  end
  structure N : sig
    type s
    datatype t = A | B of s
  end
  sharing type M.t = N.t
end

```

We can write an equivalent signature by replacing the `sharing type` line with `where type t = M.t`, which is also valid SML. 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 describe a way of weakening type equivalence (*i.e.*, making it equate more types) so that the problematic datatype matchings described in Section 3.2 succeed under the transparent interpretation. (This will not help with the problematic sharing constraints of Section 3.3.) The ideas in this section are based upon the equivalence algorithm adopted by Shao [8] 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 declaration against the spec required proving the equivalence

$$u_{imp} \equiv \mu\alpha. u_{imp} * u_{imp} + \text{int}$$

where the type on the right-hand side is just $\mu_. \text{expand}(u_{imp})$. 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 recursive type is equal to its unrolling “boxed” by a μ . An alternative formulation, equivalent to the first one by transitivity, makes two recursive types equal if their unrollings are equal, *i.e.*:

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

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 `in` function).

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 `exp-dec` matching from Section 3.2.2. The problematic constraint in that example is:

$$\text{dec}'_{spec} \equiv \text{dec}_{imp}$$

where $\text{dec}'_{spec} = \text{dec}_{spec}[\text{exp}_{imp}/\text{exp}]$ (substituting exp_{imp} for exp in dec_{imp} has no effect, since the variable does not appear free). Expanding the definitions of these types, we get the constraint:

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

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

$$\begin{aligned} \text{var} * \text{exp}_{imp} + \text{dec}'_{spec} * \text{dec}'_{spec} &\equiv \\ \text{var} * \text{exp}_{imp} + \text{dec}_{imp} * \text{dec}_{imp} &\equiv \end{aligned}$$

But we cannot prove this without first proving $\text{dec}'_{spec} \equiv \text{dec}_{imp}$, 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 δ_1 and δ_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 recursive types, 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 recursive types, the trail is simply passed unchanged from the conclusions to the premises. There are two new rules that handle recursive types:

$$\frac{\tau_1 \equiv \tau_2 \in A}{\Gamma; A \vdash \tau_1 \equiv \tau_2}$$

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

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 dec_{spec} and dec_{imp} , we add the assumption $\text{dec}_{spec} \equiv \text{dec}_{imp}$ 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

$$\text{datatype } \tau_1 = \tau_1 \text{ and } \dots \text{ and } \tau_n = \tau_n$$

(where the τ_i may be sum types in which the τ_i may occur). Suppose that some implementation matches this spec under the opaque interpretation; the implementation of each type τ_i must be a recursive type δ_i . Furthermore, the type of the τ_i -in function given in the spec is $\tau_i \rightarrow \tau_i$, 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 each δ_i has been substituted for τ_i ; thus we know that $\text{expand}(\delta_i) \equiv \tau_i[\tilde{\delta}/\tilde{\tau}]$ for each i .

When the specification is elaborated under the transparent interpretation, however, the resulting signature declares that the implementation of each τ_i is the appropriate projection from a recursive bundle determined by the spec itself. That is, each τ_i is transparently specified as $\mu_i(\tilde{\tau}).(\tilde{\tau})$. In order for the implementation to match this transparent specification, it is thus sufficient to prove the following theorem:

Theorem 1 *If $\forall i \in 1..n, \Gamma; \emptyset \vdash \text{expand}(\delta_i) \equiv \tau_i[\tilde{\delta}/\tilde{\tau}]$, then $\forall i \in 1..n, \Gamma; \emptyset \vdash \delta_i \equiv \mu_i(\tilde{\tau}).(\tilde{\tau})$.*

Proof: See Appendix A. \square

3.5 Discussion

While we have given a formal argument why the trailing version of the boxed-unroll rule is flexible enough to allow the datatype matchings of SML to typecheck under the transparent interpretation, we have not been precise about how maintaining a trail relates to the rest of type equivalence. In fact, the only work regarding trails we are aware of is the seminal work of Amadio and Cardelli [1] on subtyping equi-recursive types, and its later coinductive axiomatization by Brandt and Henglein [2], both of which are conducted in the context of the simply-typed λ -calculus. Our trailing boxed-unroll rule can be viewed as a restriction of the corresponding rule in Amadio and Cardelli's trailing algorithm so that it is only applicable when both types being compared are recursive types.

It is not clear, though, how trails affect more complex type systems that contain type constructors of higher kind, such as Girard's F^ω [6]. In addition to higher kinds, the MIL (Middle Intermediate Language) of TILT employs singleton kinds to model SML's type sharing [13], and the proof that MIL typechecking is decidable is rather delicate and involved. While we have implemented the above trailing algorithm in TILT for experimental purposes (see Section 5), the interaction of trails and singletons is not well-understood.

As for the remaining conflict between the transparent interpretation and type sharing, one might argue that the solution is to broaden SML's semantics for sharing constraints to permit sharing of *rigid* (non-abstract) type components. The problem is that the kind of sharing that would be necessary to make the examples of Section 3.3 typecheck under the transparent interpretation would require some form of type unification. It is difficult to determine where to draw the line between SML's sharing semantics and full higher-order unification, which is undecidable. Moreover, unifica-

tion would constitute a significant change to SML's semantics, disproportionate to the original problem of efficiently implementing datatypes.

4 A Coercion Interpretation of Datatypes

In this section, we will discuss a treatment of datatypes based on *coercions*. 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 either.

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 $\text{roll}_\delta(v)$ is itself a value if v is a value. (From here on we will assume that the metavariable v ranges only over values.) In fact, it can be shown without difficulty that *any* closed value of a datatype δ must have the form $\text{roll}_\delta(v)$ where v is a closed value of type $\text{expand}(\delta)$. Thus the roll operator plays a similar role to that of the 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 $\text{inj}_1(v)$ 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 inj_1 or 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 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 $\text{roll}_\delta(v)$. Thus one reasonable representation strategy is to represent $\text{roll}_\delta(v)$ exactly the same as v . In the companion technical report [15], we give a more precise argument as to why this is reasonable, making use of two key insights. First, it is an invariant of the TILT compiler that the representation of any value fits in a single machine register; anything larger than 32 bits is *always* stored in the heap. 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 τ by induction on τ as might be expected.

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

<i>Types</i>	$\sigma, \tau ::= \dots \mid (\vec{\alpha}; \tau_1) \Rightarrow \tau_2$
<i>Terms</i>	$e ::= \dots \mid \Lambda \vec{\alpha}. \text{fold}_\delta \mid \Lambda \vec{\alpha}. \text{unfold}_\delta$ $\mid v@(\vec{\tau}; e)$

Figure 8. Syntax of Coercions

the representations of these two values are the same. The reverse happens for `unroll`.

This, in turn, has an important consequence for datatypes. Since the `in` and `out` functions produced by the HS 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 use *coercions*, which are similar to functions, except that they are known to be no-ops and therefore no code needs to be generated for coercion applications. We incorporate coercions into the term level of our language and introduce special coercion types to which they belong. Figure 8 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 address non-uniform datatypes.

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 τ_1 into one of type τ_2 (where, of course, variables from $\vec{\alpha}$ can appear in either of these types). When $\vec{\alpha}$ is empty, we will write $(\vec{\alpha}; \tau_1) \Rightarrow \tau_2$ as $\tau_1 \Rightarrow \tau_2$.

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 9.

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 10. Note that the interface is exactly the same as the HS 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 HS 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))`.

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 11. 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 11 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. The untyped nature of the target language (and of machine language) is important: treating v as `fold v` would destroy the subject reduction property of a typed language.

Figure 12 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 `fold` and `unfold`, but no coercion application form. The erasure translation turns expressions with type annotations into expressions without them (λ -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 `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 2 (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 the companion technical report [15]. □

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

Figure 9. Typing Rules for Coercions

```

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 =  $\mu_2(\alpha, \beta).(\text{var} + \beta * \alpha, \text{var} * \alpha + \beta * \beta)$ 
  val exp_in = foldexp
  val exp_out = unfoldexp
  val dec_in = folddec
  val dec_out = unfolddec
end

```

Figure 10. Elaboration of exp and dec Under the Coercion Interpretation

$$\begin{array}{l}
\text{Values } v ::= \dots \mid \Lambda \vec{\alpha}. \text{fold}_{\tau} \mid \Lambda \vec{\alpha}. \text{unfold}_{\tau} \mid (\Lambda \vec{\alpha}. \text{fold}_{\tau}) @ (\vec{\sigma}; v) \\
\\
\frac{e \mapsto e'}{v @ (\vec{\tau}; e) \mapsto v @ (\vec{\tau}; e')} \quad \frac{}{(\Lambda \vec{\alpha}. \text{unfold}_{\tau_1}) @ (\vec{\sigma}; ((\Lambda \vec{\beta}. \text{fold}_{\tau_2}) @ (\vec{\sigma}'; v))) \mapsto v}
\end{array}$$

Figure 11. Operational Semantics for Coercions

$$\begin{array}{l}
M ::= \dots \mid \lambda x. M \mid \text{fold} \mid \text{unfold} \\
\\
\begin{array}{l}
x^- = x \\
(\lambda x : \tau. e)^- = \lambda x. e^- \\
(\Lambda \vec{\alpha}. \text{fold}_{\delta})^- = \text{fold} \\
(\Lambda \vec{\alpha}. \text{unfold}_{\delta})^- = \text{unfold} \\
(v @ (\vec{\tau}; e))^- = e^- \\
\vdots
\end{array}
\end{array}$$

Figure 12. Target Language Syntax; Type and Coercion Erasure

Test	HS	CID	TID
life	8.233	2.161	2.380
leroy	5.497	4.069	3.986
fft	22.167	17.619	16.509
boyer	2.031	1.559	1.364
simple	1.506	1.003	0.908
tyan	16.239	8.477	9.512
msort	1.685	0.860	1.012
pia	1.758	1.494	1.417
lexgen	11.052	5.599	5.239
frank	37.449	25.355	26.017
TOTAL	107.617	68.199	68.344

Figure 13. Performance Comparison

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, including a coercion 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 Performance

To evaluate the relative performance of the different interpretations of datatypes we have discussed, we performed experiments using three different versions of the TILT compiler: one that implements a naïve Harper-Stone interpretation in which the construction of a non-locally-defined datatype requires a function call²; one that implements the coercion interpretation of datatypes; and one that implements the transparent interpretation. We compiled ten different benchmarks using each version of the compiler; the running times for the resulting executables (averaged over three trials) are shown in Figure 13. All tests were run on an Ultra-SPARC Enterprise server; the times reported are CPU time in seconds.

The measurements clearly indicate that the overhead due to datatype constructor function calls under the naïve HS interpretation is significant. The optimizations afforded by the coercion and transparent interpretations provide comparable speedups over the opaque interpretation, both on the order of 37% (comparing the total running times). Given that, of the two optimized approaches, only the coercion interpretation is entirely faithful to the semantics of SML, and since the theory of coercion types is a simpler and more orthogonal extension to the HS type theory than the trailing algorithm of Section 3.4, we believe the coercion interpretation is the more robust choice.

6 Related Work

Our trail algorithm for weakened recursive type equivalence is based on the one implemented by Shao in the FLINT intermediate language of the Standard ML of New Jersey compiler [12]. The typing rules in Section 3.4 are based on the formal semantics for FLINT given by League and Shao [8], although we are the first to give a formal argument that their trailing algorithm actually works. It is important to note that SML/NJ only implements the transpar-

ent interpretation *internally*: the opaque interpretation is employed during elaboration, and datatype specifications are made transparent only afterward. As the examples of Section 3.3 illustrate, there are programs that typecheck according to SML but not under the transparent interpretation even with trailing equivalence, so it is unclear what SML/NJ does (after elaboration) in these cases. As it happens, the final example of Section 3.3, which is valid SML, is rejected by the SML/NJ compiler.

Curien and Ghelli [4] and Cray [3] have defined languages that use coercions to replace subsumption rules in languages with subtyping. Cray’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 Cray’s. The implementation of Typed Assembly Language for the x86 architecture (TALx86) [10] 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 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 coercion calculus.

7 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.

²In particular, we implement the strategy described at the end of Section 2.2.

Although we have presented our analysis of SML datatypes in the context of Harper-Stone and the TILT compiler, the idea of “coercion types” is one that we think is generally useful. Terms that serve only as retyping operations are pervasive in typed intermediate languages, and are usually described as “coercions” that can be eliminated before running the code. However, applications of these informal coercions cannot in general be erased if there is no way to distinguish coercions from ordinary functions by their types; this is a problem especially in the presence of true separate compilation. Our contribution is to provide a simple mechanism which permits coercive terms to be recognized as such and their applications to be safely eliminated without requiring significant syntactic and meta-theoretic overhead.

8 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 Cray. 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] Vladimir Gapeyev, Michael Levin, and Benjamin Pierce. Recursive subtyping revealed. In *2000 ACM International Conference on Functional Programming*, 2000. To appear in *Journal of Functional Programming*.
- [6] Jean-Yves Girard. *Interprétation fonctionnelle et élimination des coupures de l’arithmétique d’ordre supérieur*. PhD thesis, Université Paris VII, 1972.
- [7] 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.
- [8] Christopher League and Zhong Shao. Formal semantics of the FLINT intermediate language. Technical Report Yale-CS-TR-1171, Yale University, 1998.
- [9] Robin Milner, Mads Tofte, Robert Harper, and David MacQueen. *The Definition of Standard ML (Revised)*. MIT Press, Cambridge, Massachusetts, 1997.
- [10] Greg Morrisett, Karl Cray, 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.
- [11] 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.
- [12] 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.
- [13] Christopher A. Stone and Robert Harper. Deciding type equivalence in a language with singleton kinds. In *Twenty-Seventh ACM Symposium on Principles of Programming Languages*, pages 214–227, Boston, January 2000.
- [14] 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.
- [15] Joseph C. Vanderwaart, Derek Dreyer, Leaf Petersen, Karl Cray, Robert Harper, and Perry Cheng. Typed compilation of recursive datatypes. Technical Report CMU-CS-02-200, School of Computer Science, Carnegie Mellon University, December 2002.

A Proof of Theorem 1

Suppose that $\forall i \in 1..n, \Gamma; \emptyset \vdash \text{expand}(\delta_i) \equiv \tau_i[\vec{\delta}/\vec{\tau}]$. Then we can prove the following lemma:

Lemma 1 *For any set $S \subseteq \{1, \dots, n\}$, define $A_S = \{\delta_i \equiv \mu_i(\vec{\tau}).(\vec{\tau}) \mid i \in S\}$. Then for any $S \subseteq \{1, \dots, n\}$ and any $j \in \{1, \dots, n\}, \Gamma; A_S \vdash \delta_j \equiv \mu_j(\vec{\tau}).(\vec{\tau})$.*

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, $\Gamma; A_{S'} \vdash \delta_k \equiv \mu_k(\vec{\tau}).(\vec{\tau})$ for every $k \in \{1, \dots, n\}$. Because substituting equal types into equal types gives equal results, $\Gamma; A_{S'} \vdash \tau_j[\vec{\delta}/\vec{\tau}] \equiv \tau_j[\vec{\mu}(\vec{\tau}).(\vec{\tau})/\vec{\tau}]$. By assumption, $\text{expand}(\delta_j) \equiv \tau_j[\vec{\delta}/\vec{\tau}]$, so by transitivity $\Gamma; A_{S'} \vdash \text{expand}(\delta_j) \equiv \tau_j[\vec{\mu}(\vec{\tau}).(\vec{\tau})/\vec{\tau}]$. The type on the right side of this equivalence is just $\text{expand}(\mu_j(\vec{\tau}).(\vec{\tau}))$, so by the trailing boxed-unroll rule we can conclude $\Gamma; A_S \vdash \delta_j \equiv \mu_j(\vec{\tau}).(\vec{\tau})$, as required. \square

The desired result then follows as a corollary:

Corollary 1 *For $j \in \{1, \dots, n\}, \Gamma; \emptyset \vdash \delta_j \equiv \mu_j(\vec{\tau}).(\vec{\tau})$.*

Proof: Choose $S = \emptyset$. By the Lemma, $\Gamma; A_\emptyset \vdash \delta_j \equiv \mu_j(\vec{\tau}).(\vec{\tau})$. But $A_\emptyset = \emptyset$, so we are done. \square