

# PHANTOM TYPES AND SUBTYPING

Matthew Fluet and Riccardo Pucella

*Department of Computer Science*

*Cornell University*

{fluet,riccardo}@cs.cornell.edu

**Abstract** We investigate a technique from the literature, called the phantom types technique, that uses parametric polymorphism, type constraints, and unification of polymorphic types to model a subtyping hierarchy. Hindley-Milner type systems, such as the one found in ML, can be used to enforce the subtyping relation. We show that this technique can be used to encode any finite subtyping hierarchy (including hierarchies arising from multiple interface inheritance). We then formally demonstrate the suitability of the phantom types technique for capturing subtyping by exhibiting a type-preserving translation from a simple calculus with bounded polymorphism to a calculus embodying the type system of ML.

## 1. Introduction

It is well known that traditional type systems, such as the one found in Standard ML [10], with parametric polymorphism and type constructors can be used to capture program properties beyond those naturally associated with a Hindley-Milner type system [9]. For concreteness, let us review a simple example, due to Leijen and Meijer [8]. Consider a type of atoms, either booleans or integers, that can be easily represented as an algebraic datatype:

```
datatype atom = I of int | B of bool
```

There are a number of operations that we may perform on such atoms (see Figure 1(a)). When the domain of an operation is restricted to only one kind of atom, as with `conj` and `double`, a run-time check must be made and an error or exception reported if the check fails.

One aim of static type checking is to reduce the number of run-time checks by catching type errors at compile time. Of course, in the example above, the ML type system does not consider `conj (mkI 3, mkB true)` to be ill-typed; evaluating this expression will simply raise a run-time exception.

If we were working in a language with subtyping, we would like to consider integer atoms and boolean atoms as distinct subtypes of the general type of atoms and use these subtypes to refine the types of the operations. Then the type system would report a type error in the expression `double (mkB false)` at compile time. Fortunately, we can write the operations in a way that utilizes the ML type system to do just this. We

---

<pre> fun mkI (i:int):atom = I (i) fun mkB (b:bool):atom = B (b)  fun toString (v:atom):string =   (case v    of I (i) =&gt; Int.toString (i)       B (b) =&gt; Bool.toString (b)) fun double (v:atom):atom =   (case v    of I (i) =&gt; I (i * 2)       _ =&gt; raise Fail "type mismatch") fun conj (v1:atom,          v2:atom):atom =   (case (v1,v2)    of (B (b1), B (b2)) =&gt; B (b1 andalso b2)       _ =&gt; raise Fail "type mismatch") </pre>	<pre> fun mkI (i:int):int atom = I (i) fun mkB (b:bool):bool atom = B (b)  fun toString (v:'a atom):string =   (case v    of I (i) =&gt; Int.toString (i)       B (b) =&gt; Bool.toString (b)) fun double (v:int atom):int atom =   (case v    of I (i) =&gt; I (i * 2)       _ =&gt; raise Fail "type mismatch") fun conj (v1:bool atom,          v2:bool atom):bool atom =   (case (v1,v2)    of (B (b1), B (b2)) =&gt; B (b1 andalso b2)       _ =&gt; raise Fail "type mismatch") </pre>
(a) Unsafe operations	(b) Safe operations

---

Figure 1

change the definition of the datatype to the following:

```
datatype 'a atom = I of int | B of bool
```

and constrain the types of the operations (see Figure 1(b)). We use the superfluous type variable in the datatype definition to encode information about the kind of atom. (Because instantiations of this type variable do not contribute to the run-time representation of atoms, it is called a *phantom type*.) The type *int atom* is used to represent integer atoms and *bool atom* is used to represent boolean atoms. Now, the expression `conj (mkI 3, mkB true)` results in a compile-time type error, because the types *int atom* and *bool atom* do not unify. (Observe that our use of `int` and `bool` as phantom types is arbitrary; we could have used any two types that do not unify to make the integer versus boolean distinction.) On the other hand, both `toString (mkI 3)` and `toString (mkB true)` are well-typed; `toString` can be used on any atom. This is the essence of the technique explored in this paper: using a free type variable to encode subtyping information and using an ML-like type system to enforce the subtyping. This “phantom types” technique, where user-defined restrictions are reflected in the constrained types of values and functions, underlies many interesting uses of type systems [14, 12, 2, 13, 6, 8, 5, 11, 1].

The main contributions of this paper are to exhibit a general encoding of subtyping hierarchies and to give one formalization of the use of the phantom types technique. We present a type-preserving translation from a calculus with subtyping to a calculus with let-bounded polymorphism. The kind of subtyping that can be captured turns out to be an interesting variant of bounded polymorphism [3], with a very restricted subsumption rule.

This paper is structured as follows. In the next section, we describe a simple recipe for deriving an interface enforcing a given subtyping hierarchy. The interface is parameterized by an encoding, via phantom types, of the subtyping hierarchy. In Section 3, we focus on a simple encoding for hierarchies. In Section 4, we extend the recipe to capture a limited form of bounded polymorphism. In Section 5, we formally define the

kind of subtyping captured by our encodings by giving a simple calculus with subtyping and showing that our encodings provide a type-preserving translation to a variant of the Damas-Milner calculus, embodying the essence of the ML type system. We conclude with some problems inherent to the approach and a consideration of future work. Due to space considerations, proofs of our results, a more involved discussion of the encodings in Section 3, as well as the full typing rules for the formalization in Section 5 have been deferred to the full paper.

## 2. From subtyping to polymorphism

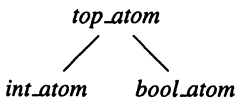


Figure 2

The example in the introduction has the following features: an underlying primitive type of values (the original type *atom*), a set of operations, and “implicit” subtypes that correspond to the sensible domains of the operations. The subtyping hierarchy corresponding to the example is given in Figure 2. The subtyping hierarchy

is modeled by assigning a type to every implicit subtype in the hierarchy. For instance, integer atoms with implicit subtype *int\_atom* are encoded by the ML type *int\_atom*. The appropriate use of polymorphic type variables in the type of an operation indicates the maximal type in the domain of the operation. For instance, the operation `toString` has the conceptual type  $top\_atom \rightarrow string$  which is encoded by the ML type  $'a\ atom \rightarrow string$ . The key observation is the use of type unification to enforce the subtyping hierarchy: an *int\_atom* can be passed to a function expecting an *'a\_atom*, because these types unify.

We consider the following problem. Given an abstract type  $\tau_p$ , a subtyping hierarchy, and an implementation of  $\tau_p$  and its operations, we wish to derive a “safe” ML signature which uses phantom types to encode the subtyping and a “safe” implementation from the “unsafe” implementation. We will call the elements of the subtyping hierarchy *implicit types* and talk about *implicit subtyping* in the hierarchy. All values share the same underlying representation and each operation has a single implementation that acts on this underlying representation. The imposed subtyping captures restrictions that arise because of some external knowledge about the semantics of the operations; intuitively, it captures a “real” subtyping relationship that is not exposed by the abstract type.

We first consider deriving the safe interface. The new interface defines a type  $\alpha\ \tau$  corresponding to the abstract type  $\tau_p$ . The type variable  $\alpha$  will be used to encode implicit subtype information. We require an encoding  $\langle \sigma \rangle$  of each implicit type  $\sigma$  in the hierarchy; this encoding should yield a type in the underlying ML type system, with the property that  $\langle \sigma_1 \rangle$  unifies with  $\langle \sigma_2 \rangle$  if and only if  $\sigma_1$  is an implicit subtype of  $\sigma_2$ . An obvious issue is that we want to use unification (a symmetric relation) to capture subtyping (an asymmetric relation). The simplest approach is to use two encodings  $\langle \cdot \rangle_C$  and  $\langle \cdot \rangle_A$  defined over all the implicit types in the hierarchy. A *value* of implicit type  $\sigma$  will be assigned a type  $\langle \sigma \rangle_C\ \tau$ . We call  $\langle \sigma \rangle_C$  the *concrete* subtype encoding of  $\sigma$ , and we assume that it uses only ground types (i.e., no type variables). In order to restrict the domain of an operation to the set of values in any implicit subtype of  $\sigma$ , we use  $\langle \sigma \rangle_A$ , the *abstract* subtype encoding of  $\sigma$ . In order for the underlying type system to enforce the subtype hierarchy, we require the encodings  $\langle \cdot \rangle_C$  and  $\langle \cdot \rangle_A$  to be

---

<pre>signature ATOM = sig   type atom   val int : int -&gt; atom   val bool : bool -&gt; atom   val toString : atom -&gt; string   val double : atom -&gt; atom   val conj : atom * atom -&gt; atom end</pre>	<pre>signature SAFE_ATOM = sig   type 'a atom   val int : int -&gt; (int)<sub>C</sub> atom   val bool : bool -&gt; (bool)<sub>C</sub> atom   val toString : (top)<sub>A</sub> atom -&gt; string   val double : (int)<sub>A</sub> atom -&gt; (int)<sub>C</sub> atom   val conj : (bool)<sub>A</sub> atom * (bool)<sub>A</sub> atom -&gt; (bool)<sub>C</sub> atom end</pre>
(a) Unsafe signature	(b) Safe signature

---

Figure 3

respectful by satisfying the following property:

for all  $\sigma_1$  and  $\sigma_2$ ,  $\langle \sigma_1 \rangle_C$  matches  $\langle \sigma_2 \rangle_A$  iff  $\sigma_1 \leq \sigma_2$ .

For example, the encodings used in the introduction are respectful:

$$\begin{array}{ll} \langle \text{top\_atom} \rangle_A = 'a \text{ atom} & \langle \text{top\_atom} \rangle_C = \text{unit atom} \\ \langle \text{int\_atom} \rangle_A = \text{int atom} & \langle \text{int\_atom} \rangle_C = \text{int atom} \\ \langle \text{bool\_atom} \rangle_A = \text{bool atom} & \langle \text{bool\_atom} \rangle_C = \text{bool atom} \end{array}$$

The utility of the phantom types technique relies on being able to find respectful encodings for subtyping hierarchies of interest.

To allow for matching, the abstract subtype encoding will introduce free type variables. Since in a Hindley-Milner type system, a type cannot contain free type variables, the abstract encoding will be part of the larger type scheme of some polymorphic function operating on the value of implicit subtypes. This leads to some restrictions on when we should constrain values by concrete or abstract encodings. We will restrict ourselves to using concrete encodings in all covariant type positions, and using abstract encodings in most contravariant type positions. We will return to this issue in Section 5.

Consider again the example from the introduction. Assume we have encodings  $\langle \cdot \rangle_C$  and  $\langle \cdot \rangle_A$  for the hierarchy and a structure `Atom` implementing the “unsafe” operations, with the signature given in Figure 3(a). Deriving an interface using the recipe above, we get the safe signature given in Figure 3(b).

We must now derive a corresponding “safe” implementation. We need a type  $\alpha \tau$  isomorphic to  $\tau_p$  such that the type system considers  $\tau_1 \tau$  and  $\tau_2 \tau$  equivalent iff  $\tau_1$  and

---

<pre>structure SafeAtom1 :&gt; SAFE_ATOM = struct   type 'a atom = Atom.atom   val int = Atom.int   val bool = Atom.bool   val toString = Atom.toString   val double = Atom.double   val conj = Atom.conj end</pre>	<pre>structure SafeAtom2 : SAFE_ATOM = struct   datatype 'a atom = C of Atom.atom   fun int (i) = C (Atom.int (i))   fun bool (b) = C (Atom.bool (b))   fun toString (C v) = Atom.toString (v)   fun double (C v) = C (Atom.double (v))   fun conj (C b1, C b2) = C (Atom.conj (b1,b2)) end</pre>
(a) Opaque signature	(b) Datatype declaration

---

Figure 4

$\tau_2$  are equivalent. (Note that this requirement precludes the use of type abbreviations of the form  $\alpha \tau = \tau_p$ , which define constant type functions.) We can then constrain the types of values and operations using  $\langle \sigma \rangle_C \tau$  and  $\langle \sigma \rangle_A \tau$ . In ML, the easiest way to achieve this is to use an abstract type at the module system level, as shown in Figure 4(a). The use of an opaque signature is critical to get the required behavior in terms of type equivalence. The advantage of this method is that there is no overhead.

In a language without abstract types at the module level, another approach is to wrap the primitive type  $\tau_p$  using a datatype declaration

```
datatype 'a  $\tau = C$  of  $\tau_p$ 
```

The type  $\alpha \tau$  behaves as required, because the datatype declaration defines a generative type operator. However, we must explicitly convert primitive values to and from  $\alpha \tau$  to witness the isomorphism. This yields the implementation given in Figure 4(b).

We should stress that the “safe” interface must ensure that the type  $\alpha \tau$  is abstract—either through the use of opaque signature matching, or by hiding the value constructors of the type. Otherwise, it may be possible to create values that do not respect the subtyping invariants enforced by the encodings. Similarly, the use of an abstract subtype encoding in a covariant type position can lead to violations in the subtyping invariants.

We now have a way to derive a safe interface and implementation, by adding type information to a generic, unsafe implementation. In the next section, we show how to construct respectful encodings  $\langle \cdot \rangle_C$  and  $\langle \cdot \rangle_A$  by taking advantage of the structure of the subtyping hierarchy.

### 3. Encoding subtyping hierarchies

The framework presented in the previous section relies on having concrete and abstract encodings of the implicit subtypes in the subtyping hierarchy with the property that unification of the results of the encoding respects the subtype relation. In this section, we describe one general construction for such encodings.

We first consider a particular lattice that will be useful in our development. Recall that a lattice is a hierarchy where every set of elements has both a least upper bound and a greatest lower bound. Given a finite set  $S$ , we let the *powerset lattice* of  $S$  be the lattice of all subsets of  $S$ , ordered by inclusion, written  $(\wp(S), \subseteq)$ . We now exhibit an encoding of powerset lattices.

Let  $n$  be the cardinality of  $S$  and assume an ordering  $s_1, \dots, s_n$  on the elements of  $S$ . We encode subset  $X$  of  $S$  as an  $n$ -tuple type, where the  $i^{\text{th}}$  entry expresses that  $s_i \in X$  or  $s_i \notin X$ . First, we introduce a datatype definition:

```
datatype 'a  $z = Z$ 
```

(The name of the datatype constructor is irrelevant, because we will never construct values of this type.) The encoding of an arbitrary subset of  $S$  is given by:

$$\begin{aligned} \langle X \rangle_C &= (t_1, \dots, t_n) \text{ where } t_i = \begin{cases} \text{unit} & \text{if } s_i \in X \\ \text{unit } z & \text{otherwise} \end{cases} \\ \langle X \rangle_A &= (t_1, \dots, t_n) \text{ where } t_i = \begin{cases} \alpha_i & \text{if } s_i \in X \\ \alpha_i z & \text{otherwise} \end{cases} \end{aligned}$$

Note that  $\langle \cdot \rangle_A$  requires every type variable  $\alpha_i$  to be a fresh type variable, unique in its context. This ensures that we do not inadvertently refer to any type variable bound in the context where we are introducing the abstractly encoded type.

As an example, consider the powerset lattice of  $\{1, 2, 3, 4\}$ , which encodes into a four-tuple. We can verify, for example, that the concrete encoding for  $\{2\}$ , namely  $(\text{unit } z, \text{unit}, \text{unit } z, \text{unit } z)$ , unifies with the abstract encoding for  $\{1, 2\}$ , namely  $(\alpha_1, \alpha_2, \alpha_3 z, \alpha_4 z)$ . On the other hand, the concrete encoding of  $\{1, 2\}$  does not unify with the abstract encoding of  $\{2, 3\}$ .

The main reason we introduced powerset lattices is the fact that any finite hierarchy can be embedded in the powerset lattice of a set  $S$ . It is a simple matter, given a hierarchy  $H'$  embedded in a hierarchy  $H$ , to derive an encoding for  $H'$  given an encoding for  $H$ . Let  $\text{inj}(\cdot)$  be the injection from  $H'$  to  $H$  witnessing the embedding and let  $\langle \cdot \rangle_{C_H}$  and  $\langle \cdot \rangle_{A_H}$  be the encodings for the hierarchy  $H$ . Deriving an encoding for  $H'$  simply involves defining  $\langle \sigma \rangle_{C_{H'}} = \langle \text{inj}(\sigma) \rangle_{C_H}$  and  $\langle \sigma \rangle_{A_{H'}} = \langle \text{inj}(\sigma) \rangle_{A_H}$ . It is straightforward to verify that if  $\langle \cdot \rangle_{C_H}$  and  $\langle \cdot \rangle_{A_H}$  are respectful encodings, so are  $\langle \cdot \rangle_{C_{H'}}$  and  $\langle \cdot \rangle_{A_{H'}}$ . By the result above, this allows us to derive an encoding for an arbitrary finite hierarchy.

We have presented a strategy for obtaining respectful encodings, which is sufficient for the remainder of this paper. However, there are encodings for specific hierarchies that are in general more efficient than their embedding in a powerset lattice, for instance, the encoding for tree hierarchies found in [6]. We discuss such encodings and address the issue of encoding extensibility in the full paper.

## 4. Towards bounded polymorphism

As mentioned in Section 3, the handling of type variables is somewhat delicate. If we allow common type variables to be used across abstract encodings, then we can capture a form of *bounded polymorphism* as in  $F_{<}$ . [3]. Bounded polymorphism à la  $F_{<}$  is a typing discipline which extends both parametric polymorphism and subtyping. From parametric polymorphism, it borrows type variables and universal quantification; from subtyping, it allows one to set bounds on quantified type variables. For example, one can guarantee that the argument and return types of a function are the same and a subtype of  $\sigma$ , as in  $\forall \alpha \leq \sigma. \alpha \rightarrow \alpha$ . Similarly, one can guarantee that two arguments have the same type that is a subtype of  $\sigma$ , as in  $\forall \alpha \leq \sigma. (\alpha \times \alpha) \rightarrow \sigma$ . Notice that neither function can be written in a language that supports only subtyping.

Returning to the example from the introduction, consider adding natural numbers as a subtype of integers, so that  $\text{nat\_atom}$  is a subtype of  $\text{int\_atom}$ . Using bounded polymorphism, we can assign to `double` the reasonable type  $\forall \alpha \leq \text{int\_atom}. \alpha \rightarrow \alpha$ . However, bounded polymorphism has its limitations. One reasonable type for a plus operation is  $\forall \alpha \leq \text{int\_atom}. \alpha \times \alpha \rightarrow \alpha$  where the same kind of atom is required for both arguments. In order to add an integer and a natural number we need a function `toInt` (operationally, an identity function) to coerce the type of the natural number to that of an integer.

We can adapt our “recipe” from Section 2 to types of the form  $\forall \beta \leq \sigma_1. (\beta \times \sigma_2) \rightarrow \beta$ . Let the “safe” interface use types of the form  $\alpha \tau$ . Since  $\beta$  stands for a subtype of  $\sigma_1$ , we let  $\phi_\beta = \langle \sigma_1 \rangle_A$ , the abstract encoding of the bound. We then translate the type

as we did in Section 2, but replace occurrences of the type variable  $\beta$  by  $\phi_\beta$  instead of applying  $\langle \cdot \rangle_A$  repeatedly, thereby sharing the type variables introduced by  $\langle \sigma_1 \rangle_A$ . Hence, we get the type  $\phi_\beta \tau \times \langle \sigma_2 \rangle_A \tau \rightarrow \phi_\beta \tau$ . In fact, we can further simplify the process by noting that we can “pull out” all the subtyping into bounded polymorphism. If a function expects an argument of any implicit subtype of  $\sigma$ , we can introduce a fresh type variable for that argument and bound it by  $\sigma$ . For example, the type above can be rewritten as:  $\forall \beta \leq \sigma_1, \gamma \leq \sigma_2. (\beta \times \gamma) \rightarrow \beta$ .

Unfortunately, this technique does not generalize to full  $F_{<}$ . For example, we cannot encode bounded polymorphism where the bound on a type variable uses a type variable, such as a function  $\mathbf{f}$  with type  $\forall \alpha \leq \sigma, \beta \leq \alpha. \alpha \times \beta \rightarrow \alpha$ . Encoding this type as  $\phi_\alpha \tau \times \phi_\beta \tau \rightarrow \phi_\alpha \tau$  where  $\phi_\alpha = \langle \sigma \rangle_A$  and  $\phi_\beta = \langle \alpha \rangle_A$  fails, because we have no definition of  $\langle \alpha \rangle_A$ . Essentially, we need a different encoding of  $\beta$  for each instantiation of  $\alpha$  at each application of  $\mathbf{f}$ , something that cannot be accommodated by a single encoding of the type at the definition of  $\mathbf{f}$ .

Likewise, we cannot encode first-class polymorphism, such as a function  $\mathbf{g}$  with type  $\forall \alpha \leq \sigma_1. \alpha \rightarrow (\forall \beta \leq \sigma_2. \beta \rightarrow \beta)$ . Applying the technique yields a type  $\phi_\alpha \tau \rightarrow \phi_\beta \tau \rightarrow \phi_\beta \tau$  where  $\phi_\alpha$  and  $\phi_\beta$  contain free type variables. A Hindley-Milner style type system requires quantification over these variables in prenex position, which doesn’t match the intuition of the original type. In fact, because we are translating into a language with prenex polymorphism, we can only capture bounded polymorphism that is itself in prenex form.

In other words, we cannot account for the general substitution rule found in  $F_{<}$ . Instead, we require all subtyping to occur at type application. This is the real motivation for the simplification above which “pulls out” all subtyping into bounded polymorphism. By introducing type variables for each argument, we move the resolution of the subtyping to the point of type application (when we instantiate the type variables).

These two restrictions impose one final restriction on the kind of subtyping we can encode. Consider a higher-order function  $\mathbf{h}$  with type  $\alpha \leq (\sigma_1 \rightarrow \sigma_2). \alpha \rightarrow \sigma_2$ . What are the possible encodings of the bound  $\sigma_1 \rightarrow \sigma_2$  that allow subtyping? Clearly encoding the bound as  $\langle \sigma_1 \rangle_C \tau \rightarrow \langle \sigma_2 \rangle_C \tau$  does not allow any subtyping. Encoding the bound as  $\langle \sigma_1 \rangle_A \tau \rightarrow \langle \sigma_2 \rangle_A \tau$  or  $\langle \sigma_1 \rangle_A \tau \rightarrow \langle \sigma_2 \rangle_C \tau$  leads to an unsound system. (Consider applying the argument function to a value of type  $\sigma_0 \geq \sigma_1$ , which would type-check in the encoding, because  $\langle \sigma_0 \rangle_C$  unifies with  $\langle \sigma_1 \rangle_A$  by the definition of a respectful encoding.) However, we can soundly encode the bound as  $\langle \sigma_1 \rangle_C \tau \rightarrow \langle \sigma_2 \rangle_A \tau$ . This corresponds to a subtyping rule on functional types that asserts  $\tau_1 \rightarrow \tau_2 \leq \tau_1 \rightarrow \tau'_2$  iff  $\tau_2 \leq \tau'_2$ .

Despite these restrictions, the phantom types technique is still a viable method for encoding subtyping in a language like ML. All of the examples of phantom types found in the literature satisfy these restrictions. In practice, one rarely needs first-class polymorphism or complicated dependencies between the subtypes of function arguments, particularly when implementing a safe interface to existing library functions.

## 5. A formalization

There are subtle issues regarding the kind of subtyping that can be captured using phantom types. In this section, we clarify the picture by exhibiting a typed calculus

with a suitable notion of subtyping that can be faithfully translated into a language such as ML, via a phantom types encoding. The idea is simple: to see if an interface can be implemented using phantom types, first express the interface in this calculus in such a way that the program type-checks. If it is possible to do so, our results show that a translation using phantom types exists. The target of the translation is a calculus embodying the essence of ML, essentially the calculus of Damas and Milner [4], a predicative polymorphic  $\lambda$ -calculus.

Let us first introduce the source calculus,  $\lambda_{<}^{\text{DM}}$ , also a variant of the Damas-Milner calculus, but with a very restricted notion of subtyping, and allowing multiple types for constants. We assume a partially ordered set  $(T, \leq)$  of basic types. The types and prenex quantified type schemes of  $\lambda_{<}^{\text{DM}}$  are as follows:

$$\begin{aligned}\tau &::= t \mid \alpha \mid \tau_1 \rightarrow \tau_2 \\ \sigma &::= \forall \alpha_1 <: \tau_1, \dots, \alpha_n <: \tau_n. \tau\end{aligned}$$

(where  $t \in T$ ). Furthermore, we make a syntactic restriction that precludes the use of type variables in the bounds of quantified type variables.

An important aspect of our calculus, at least for our purposes, is the constants that we allow. We distinguish between two types of constants: basic constants and primitive operations. Basic constants, taken from a set  $C_b$ , are constants representing values of basic types  $t \in T$ . We suppose a function  $\pi_b : C_b \rightarrow T$  assigning a basic type to every basic constant. The primitive operations, taken from a set  $C_p$ , are operations acting on constants and returning constants.<sup>1</sup> Rather than giving primitive operations polymorphic types, we assume that the operations can have multiple types, which encode the allowed subtyping. The primitive operation `double` in our example would get the types  $\text{int\_value} \rightarrow \text{int\_value}$  and  $\text{nat\_value} \rightarrow \text{nat\_value}$ . We suppose a function  $\pi_p$  assigning to every constant  $c \in C_p$  a set of types  $\pi_p(c)$ , each type a functional type of the form  $t \rightarrow t'$  (for  $t, t' \in T$ ).

Our expression language is again typical:

$$\begin{aligned}e &::= c \mid \lambda x : \tau. e \mid e_1 e_2 \mid x \mid p[\tau_1, \dots, \tau_n] \mid \text{let } x = p \text{ in } e \\ p &::= x \mid \Lambda \alpha_1 <: \tau_1, \dots, \alpha_n <: \tau_n. e \\ v &::= c \mid \lambda x : \tau. e \\ E &::= [] \mid E e \mid v E \mid E[\tau_1, \dots, \tau_n] \mid \text{let } x = E \text{ in } e\end{aligned}$$

(where  $c \in C_b \cup C_p$ ). The operational semantics are given using a standard rewriting system. The basic reductions are

$$\begin{aligned}(\lambda x : \tau. e) v &\rightarrow_{<} e\{v/x\} \\ (\Lambda \alpha_1 <: \tau_1, \dots, \alpha_n <: \tau_n. e) [\tau'_1, \dots, \tau'_n] &\rightarrow_{<} e\{\tau'_1/\alpha_1, \dots, \tau'_n/\alpha_n\} \\ \text{let } x = v \text{ in } e &\rightarrow_{<} e\{v/x\} \\ c_1 c_2 &\rightarrow_{<} c_3 \quad \text{iff } \delta(c_1, c_2) = c_3\end{aligned}$$

where  $\delta : C_p \times C_b \rightarrow C_p$  is a partial function defining the result of applying a primitive operation to a basic constant. This reduction extends to contexts via the rule:

$$E[e_1] \rightarrow_{<} E[e_2] \quad \text{iff } e_1 \rightarrow_{<} e_2$$

<sup>1</sup>For simplicity, we will not deal with higher-order functions here—they would simply complicate the formalism without bringing any new insight. Likewise, allowing primitive operations to act on and return tuples of values is a simple extension of the formalism presented here.



As previously noted, we only allow primitive operations to be monotyped. However, we can easily use the fact that they can take on many types to write polymorphic wrappers. Returning to the `double` example, we can write a polymorphic wrapper  $\Lambda\alpha <: \text{int\_value}.\lambda x:\alpha.\text{double } x$  to capture the expected behavior. We will see shortly that this function is well-typed.

The typing rules for  $\lambda_{<}^{\text{DM}}$  are the standard Damas-Milner typing rules, modified to account for subtyping. Subtyping is given by a judgment  $\Delta \vdash_{<} \tau_1 <: \tau_2$ , and is derived from the subtyping on the basic types. The interesting rules are:

$$\frac{t_1 \leq t_2}{\Delta \vdash_{<} t_1 <: t_2} \quad \frac{\Delta \vdash_{<} \tau_2 <: \tau'_2}{\Delta \vdash_{<} \tau_1 \rightarrow \tau_2 <: \tau_1 \rightarrow \tau'_2}$$

Notice that subtyping at higher types only involves the result type. The typing rules are given by judgments  $\Delta; \Gamma \vdash_{<} e : \tau$  for monotypes and  $\Delta; \Gamma \vdash_{<} p : \sigma$  for type schemes. The rule for primitive operations is interesting:

$$\frac{\forall \tau'_1 <: \tau_1, \dots, \tau'_n <: \tau_n \quad \tau\{\tau'_1/\alpha_1, \dots, \tau'_n/\alpha_n\} \in \pi_p(c)}{\Delta, \alpha_1 <: \tau_1, \dots, \alpha_n <: \tau_n; \Gamma \vdash_{<} c : \tau} \quad \left( \begin{array}{l} c \in C_p, \\ FV(\tau) = \{\alpha_1, \dots, \alpha_n\} \end{array} \right)$$

The syntactic restriction on type variable bounds ensures that each  $\tau_i$  has no type variables, so each  $\tau'_i <: \tau_i$  is well-defined. The rule captures the notion that any subtyping on a primitive operation through the use of bounded polymorphism is in fact realized by the “many types” interpretation of the operation.

Subtyping occurs at type application:

$$\frac{\Delta; \Gamma \vdash_{<} p : \forall \alpha_1 <: \tau_1, \dots, \alpha_n <: \tau_n. \tau \quad \Delta \vdash_{<} \tau'_1 <: \tau_1 \quad \dots \quad \Delta \vdash_{<} \tau'_n <: \tau_n}{\Delta; \Gamma \vdash_{<} p[\tau'_1, \dots, \tau'_n] : \tau\{\tau'_1/\alpha_1, \dots, \tau'_n/\alpha_n\}}$$

As discussed in the previous section, there is no subsumption in the system: subtyping must be witnessed by type application. Hence, there is a difference between the type  $t_1 \rightarrow t_2$  (where  $t_1, t_2 \in T$ ) and  $\forall \alpha <: t_1. \alpha \rightarrow t_2$ ; namely, the former does not allow any subtyping. The restrictions of Section 4 are formalized by prenex quantification and the syntactic restriction on type variable bounds.

Clearly, type soundness of the above system depends on the definition of  $\delta$  over the constants. We say that  $\pi_p$  is sound with respect to  $\delta$  if for all  $c_1 \in C_p$  and  $c_2 \in C_b$ , we have  $\vdash_{<} c_1 c_2 : \tau$  implies that  $\delta(c_1, c_2)$  is defined and  $\pi_b(\delta(c_1, c_2)) = \tau$ . This definition ensures that any application of a primitive operation  $c_1$  to a basic constant  $c_2$  yields exactly one value  $\delta(c_1, c_2)$  at exactly one type  $\pi_b(\delta(c_1, c_2)) = \tau$ . This leads to the following conditional type soundness result for  $\lambda_{<}^{\text{DM}}$ :

**Theorem 1** *If  $\pi_p$  is sound with respect to  $\delta$ ,  $\vdash_{<} e : \tau$ , and  $e \rightarrow_{<} e'$ , then  $\vdash_{<} e' : \tau$  and either  $e'$  is a value or there exists  $e''$  such that  $e' \rightarrow_{<} e''$ .*

Our target calculus,  $\lambda_{\top}^{\text{DM}}$ , is meant to capture the appropriate aspects of ML that are relevant for the phantom types encoding of subtyping. Essentially, it is the Damas-Milner calculus [4] extended with a single type constructor  $\top$ . Formally,

$$\begin{aligned} \tau &::= \alpha \mid \tau_1 \rightarrow \tau_2 \mid \top \tau \mid 1 \mid \tau_1 \times \tau_2 \\ \sigma &::= \forall \alpha_1, \dots, \alpha_n. \tau \\ e &::= c \mid \lambda x:\tau. e \mid e_1 e_2 \mid p[\tau_1, \dots, \tau_n] \mid x \mid \text{let } x = p \text{ in } e \\ v &::= c \mid \lambda x:\tau. e \\ p &::= x \mid \Lambda \alpha_1 \dots, \alpha_n. e \\ E &::= [] \mid E e \mid v E \mid E [\tau_1, \dots, \tau_n] \mid \text{let } x = E \text{ in } e \end{aligned}$$

The operational semantics (via a reduction relation  $\longrightarrow_\tau$ ) and most typing rules (via a judgment  $\Delta; \Gamma \vdash_\tau e : \tau$ ) are standard. As before, we assume that we have constants  $C_b$  and  $C_p$  and a function  $\delta$  providing semantics for primitive applications. Likewise, we assume that  $\pi_b$  and  $\pi_p$  provide types for constants, with the same restrictions. The typing rule for primitive operations in  $\lambda_{\tau}^{\text{DM}}$  is similar to the corresponding rule in  $\lambda_{\tau}^{\text{M}}$ . Given two types  $\tau$  and  $\tau'$  in  $\lambda_{\tau}^{\text{DM}}$ , we define their unification  $\text{unify}(\tau, \tau')$  to be a sequence of bindings  $\langle (\alpha_1, \tau_1), (\alpha_2, \tau_2), \dots \rangle$  in depth-first, left-to-right order of appearance of  $\alpha_1, \dots, \alpha_n$  in  $\tau$ , or  $\emptyset$  if  $\tau'$  is not a substitution instance of  $\tau$ . Given a type  $\tau$  in  $\lambda_{\tau}^{\text{DM}}$ , we define  $FV(\tau)$  to be the sequence of free type variables appearing in  $\tau$ , in depth-first, left-to-right order.

$$\frac{\forall \tau' \in \pi_b(C_b) \text{ with } \text{unify}(\tau, \tau') = \langle (\alpha_1, \tau'_1), \dots, (\alpha_n, \tau'_n), \dots \rangle \quad \left( \begin{array}{l} c \in C_p, \\ FV(\tau_1 \rightarrow \tau_2) = \\ \langle \alpha_1, \dots, \alpha_n \rangle \end{array} \right)}{\Delta, \alpha_1, \dots, \alpha_n; \Gamma \vdash_\tau c : \tau_1 \rightarrow \tau_2}$$

Again, this rule captures our notion of “subtyping through unification” by ensuring that the operation is defined at every basic type that unifies with its argument type. Our notion of soundness of  $\pi_p$  with respect to  $\delta$  carries over and we can again establish a conditional type soundness result:

**Theorem 2** *If  $\pi_p$  is sound with respect to  $\delta$ ,  $\vdash_\tau e : \tau$ , and  $e \longrightarrow_\tau e'$ , then  $\vdash_\tau e' : \tau$  and either  $e'$  is a value or there exists  $e''$  such that  $e' \longrightarrow_\tau e''$ .*

Note that the types  $\top$ ,  $1$ , and  $\tau_1 \times \tau_2$  have no corresponding introduction and elimination expressions. We include these types for the exclusive purpose of constructing the phantom types used by the encodings. We could add other types to allow more encodings, but these suffice for the lattice encodings of Section 3.

Thus far, we have a calculus  $\lambda_{\tau}^{\text{DM}}$  embodying the notion of subtyping that interests us and a calculus  $\lambda_{\tau}^{\text{PM}}$  capturing the essence of the ML type system. We now establish a translation from the first calculus into the second using phantom types to encode the subtyping, showing that we can indeed capture that particular notion of subtyping in ML. Moreover, we show that the translation preserves the soundness of the types assigned to constants, thereby guaranteeing that if the original system was sound, the system obtained by translation is sound as well.

We first describe how to translate types in  $\lambda_{\tau}^{\text{DM}}$ . Since subtyping is only witnessed at type abstraction, the type translation realizes the subtyping using the phantom types encoding of abstract and concrete subtypes. The translation is parameterized by an environment  $\rho$  associating every (free) type variable with a type in  $\lambda_{\tau}^{\text{PM}}$  representing the abstract encoding of the bound.

$$\begin{aligned} \mathcal{T}[\alpha]\rho &= \rho(\alpha) \\ \mathcal{T}[t]\rho &= \top \langle t \rangle_C \\ \mathcal{T}[\tau_1 \rightarrow \tau_2]\rho &= \mathcal{T}[\tau_1]\rho \rightarrow \mathcal{T}[\tau_2]\rho \\ \mathcal{T}[\forall \alpha_1 <: \tau_1, \dots, \alpha_n <: \tau_n. \tau]\rho &= \\ &\forall \alpha_{11}, \dots, \alpha_{1k_1}, \dots, \alpha_{n1}, \dots, \alpha_{nk_n}. \mathcal{T}[\tau]\rho[\alpha_i \mapsto \tau_i^A] \\ &\text{where } \tau_i^A = \mathcal{A}[\tau_i] \text{ and } FV(\tau_i^A) = \langle \alpha_{i1}, \dots, \alpha_{ik_i} \rangle \end{aligned}$$

If  $\rho$  is empty, we will simply write  $\mathcal{T}[\tau]$ . To compute the abstract and concrete encodings of a type, we define:

$$\begin{aligned} \mathcal{A}[t] &= \top \langle t \rangle_A & \mathcal{C}[t] &= \top \langle t \rangle_C \\ \mathcal{A}[\tau_1 \rightarrow \tau_2] &= \mathcal{C}[\tau_1] \rightarrow \mathcal{A}[\tau_2] & \mathcal{C}[\tau_1 \rightarrow \tau_2] &= \mathcal{C}[\tau_1] \rightarrow \mathcal{C}[\tau_2] \end{aligned}$$

Note that the syntactic restriction on type variable bounds ensures that  $\mathcal{A}$  and  $\mathcal{C}$  are always well-defined, as they will never be applied to type variables. Furthermore, observe that the above translation depends on the fact that the type encodings  $\langle t \rangle_C$  and  $\langle t \rangle_A$  are expressible in the  $\lambda_{\top}^{\text{DM}}$  type system using  $\top$ ,  $1$ , and  $\times$ .

We extend the type transformation  $\mathcal{T}$  to type contexts  $\Gamma$  in the obvious way:

$$\mathcal{T}[x_1 : \tau_1, \dots, x_n : \tau_n]\rho = x_1 : \mathcal{T}[\tau_1]\rho, \dots, x_n : \mathcal{T}[\tau_n]\rho$$

Finally, if we take the basic constants and the primitive operations in  $\lambda_{<}^{\text{DM}}$  and assume that  $\pi_p$  is sound with respect to  $\delta$ , then the translation can be used to assign types to the constants and operations such that they are sound in the target calculus. We first extend the definition of  $\mathcal{T}$  to  $\pi_b$  and  $\pi_p$  in the obvious way:

$$\begin{aligned} \mathcal{T}[\pi_b] &= \pi'_b \text{ where } \pi'_b(c) = \mathcal{T}[\pi_b(c)] \\ \mathcal{T}[\pi_p] &= \pi'_p \text{ where } \pi'_p(c) = \{\mathcal{T}[\tau] \mid \tau \in \pi_p(c)\} \end{aligned}$$

We can further show that the translated types do not allow us to “misuse” the constants in  $\lambda_{\top}^{\text{DM}}$ :

**Theorem 3** *If  $\pi_p$  is sound with respect to  $\delta$  in  $\lambda_{<}^{\text{DM}}$ , then  $\mathcal{T}[\pi_p]$  is sound with respect to  $\delta$  in  $\lambda_{\top}^{\text{DM}}$ .*

We can now define the translation of expressions via a translation of typing derivations,  $\mathcal{E}$ , taking care to respect the types given by the above type translation. We note that the translation below only works if the concrete encodings being used do not contain free type variables. Again, the translation is parameterized by an environment  $\rho$ , as in the type translation.

$$\begin{aligned} \mathcal{E}[\Delta; \Gamma \vdash_{<} x : \tau]\rho &= x \\ \mathcal{E}[\Delta; \Gamma \vdash_{<} c : \tau]\rho &= c \\ \mathcal{E}[\Delta; \Gamma \vdash_{<} \lambda x : \tau'. e : \tau]\rho &= \lambda x : \mathcal{T}[\tau]\rho. \mathcal{E}[e]\rho \\ \mathcal{E}[\Delta; \Gamma \vdash_{<} e_1 e_2 : \tau]\rho &= (\mathcal{E}[e_1]\rho) \mathcal{E}[e_2]\rho \\ \mathcal{E}[\Delta; \Gamma \vdash_{<} \text{let } x = p \text{ in } e : \tau]\rho &= \text{let } x = \mathcal{E}[p]\rho \text{ in } \mathcal{E}[e]\rho \\ \mathcal{E}[\Delta; \Gamma \vdash_{<} p[\tau_1, \dots, \tau_n] : \tau]\rho &= \\ & \quad (\mathcal{E}[p]\rho)[\tau_{11}, \dots, \tau_{1k_1}, \dots, \tau_{n1}, \dots, \tau_{nk_n}] \\ & \quad \text{where } \mathcal{B}[p]\Gamma = (\alpha_1, \tau_1^B), \dots, (\alpha_n, \tau_n^B) \text{ and } \tau_i^A = \mathcal{A}[\tau_i^B] \\ & \quad \text{and } FV(\tau_i^B) = \langle \alpha_{i1}, \dots, \alpha_{ik_i} \rangle \text{ and } \tau_i^T = \mathcal{T}[\tau_i]\rho \\ & \quad \text{and } \text{unify}(\tau_i^A, \tau_i^T) = \langle (\alpha_{i1}, \tau_{i1}), \dots, (\alpha_{ik_i}, \tau_{ik_i}), \dots \rangle \end{aligned}$$

$$\begin{aligned} \mathcal{E}[\Delta; \Gamma \vdash_{<} x : \sigma]\rho &= x \\ \mathcal{E}[\Delta; \Gamma \vdash_{<} \Lambda \alpha_1 <: \tau_1, \dots, \alpha_n <: \tau_n. e : \sigma]\rho &= \\ & \quad \Lambda \alpha_{11}, \dots, \alpha_{1k_1}, \dots, \alpha_{n1}, \dots, \alpha_{nk_n}. \mathcal{E}[e]\rho[\alpha_i \mapsto \tau_i^A] \\ & \quad \text{where } \tau_i^A = \mathcal{A}[\tau_i] \text{ and } FV(\tau_i^A) = \langle \alpha_{i1}, \dots, \alpha_{ik_i} \rangle \end{aligned}$$

Again, if  $\rho$  is empty, we simply write  $\mathcal{E}[e]$ . The function  $\mathcal{B}$  returns the bounds of a type abstraction, using the environment  $\Gamma$  to resolve variables. It is defined as follows:

$$\begin{aligned} \mathcal{B}[x]\Gamma &= \langle (\alpha_1, \tau_1), \dots, (\alpha_n, \tau_n) \rangle \\ \text{where } \Gamma(x) &= \forall \alpha_1 <: \tau_1, \dots, \alpha_n <: \tau_n. \tau \\ \mathcal{B}[\Lambda \alpha_1 <: \tau_1, \dots, \alpha_n <: \tau_n. e]\Gamma &= \langle (\alpha_1, \tau_1), \dots, (\alpha_n, \tau_n) \rangle \end{aligned}$$

We use  $\mathcal{B}$  and *unify* to perform unification “by hand.” In most programming languages, type inference performs this automatically.

We can verify that this translation is type-preserving:

**Theorem 4** *If  $\vdash_{<} e : \tau$ , then  $\vdash_{\tau} \mathcal{E}[\vdash_{<} e : \tau] : \mathcal{T}[\tau]$ .*

Theorem 4 is interesting in that it shows that the translation, in a sense, captures the right notion of subtyping, particularly when designing an interface. Given a set of constants making up the interface, suppose we can assign types to those constants in  $\lambda_{<}^{\text{DM}}$  in a way that gives the desired subtyping; that is, we can write type correct expressions of the form  $\Lambda\alpha <: t. \lambda x : \alpha. c x$  with type  $\forall\alpha <: t. \alpha \rightarrow \tau$ . In other words, the typing  $\pi_p$  is sound with respect to the semantics of  $\delta$ . By Theorem 1, this means that  $\lambda_{<}^{\text{DM}}$  with these constants is sound and we can safely use these constants in  $\lambda_{<}^{\text{DM}}$ . In particular, we can write the program:

$$\begin{aligned} & \text{let } f_1 = \Lambda\alpha <: t_{i_1}. \lambda x : \alpha. c_1 x \text{ in} \\ & \quad \vdots \\ & \text{let } f_n = \Lambda\alpha <: t_{i_n}. \lambda x : \alpha. c_n x \text{ in} \\ & e \end{aligned}$$

By Theorem 4, the translation of the above program executes without run-time errors. Furthermore, by Theorem 3, the phantom types encoding of the types of these constants are sound with respect to  $\delta$  in  $\lambda_{\tau}^{\text{DM}}$ . Hence, by Theorem 2,  $\lambda_{\tau}^{\text{DM}}$  with these constants is sound and we can safely use these constants in  $\lambda_{\tau}^{\text{DM}}$ . Therefore, we can replace to the body of the translated program with an arbitrary  $\lambda_{\tau}^{\text{DM}}$  expression that type-checks in that context and the resulting program will still execute without run-time errors. Essentially, the translation of the let bindings corresponds to a “safe” interface to the primitives; programs that use this interface in a type-safe manner are guaranteed to execute without run-time errors.

## 6. Conclusion

Essentially, the phantom types technique uses the definition of type equivalence in ML to encode information in a free type variable of a type. Unification can then be used to enforce a particular structure on the information carried by two such types. In this paper, we have focused on encoding subtyping information. We also showed how to extend the techniques we developed to encode a form of prenex bounded polymorphism, with subsumption occurring only at type application. It goes without saying that this approach to encoding subtyping is not without its problems from a practical point of view. As the encodings in this paper show, the types involved can become quite large. Type abbreviations can help simplify the presentation of concrete types, but for abstract encodings, which require type variables, those type variables must appear in the interface.

We also note that the source language of Section 5 provides only a lower bound on the power of phantom types. For example, one can use features of the *specific* encoding used to further constrain operations [13]. One can also capture programming invariants associated with user-defined datatypes [7]. An interesting direction for future work is formalizing these additional applications of the phantom types technique.

## Acknowledgments

We have benefitted from discussions with Greg Morrisett and Dave MacQueen. John Reppy pointed out the work of Burton. Stephanie Weirich, Vicky Weissman, and Steve Zdancewic provided helpful comments on an early draft of this paper. Thanks also to the anonymous referees. The second author was partially supported by ONR grant N00014-00-1-03-41.

## References

- [1] M. Blume. No-Longer-Foreign: Teaching an ML compiler to speak C “natively”. In *Electronic Notes in Theoretical Computer Science*, volume 59. Elsevier Science Publishers, 2001.
- [2] F. Burton. Type extension through polymorphism. *ACM Transactions on Programming Languages and Systems*, 12(1):135–138, January 1990.
- [3] L. Cardelli, S. Martini, J. C. Mitchell, and A. Scedrov. An extension of System F with subtyping. *Information and Computation*, 109(1–2):4–56, 1994.
- [4] L. Damas and R. Milner. Principal type-schemes for functional programs. In *Conference Record of the Ninth Annual ACM Symposium on Principles of Programming Languages*, pages 207–212. ACM Press, 1982.
- [5] C. Elliott, S. Finne, and O. de Moor. Compiling embedded languages. In *Workshop on Semantics, Applications, and Implementation of Program Generation*, 2000.
- [6] S. Finne, D. Leijen, E. Meijer, and S. Peyton Jones. Calling hell from heaven and heaven from hell. In *Proceedings of the 1999 ACM SIGPLAN International Conference on Functional Programming*, pages 114–125. ACM Press, 1999.
- [7] S. Kahrs. Red-black trees with types. *Journal of Functional Programming*, 11(3):425–432, 2001.
- [8] D. Leijen and E. Meijer. Domain specific embedded compilers. In *Proceedings of the Second Conference on Domain-Specific Languages (DSL’99)*, pages 109–122, 1999.
- [9] R. Milner. A theory of type polymorphism in programming. *Journal of Computer and Systems Sciences*, 17(3):348–375, 1978.
- [10] R. Milner, M. Toft, R. Harper, and D. MacQueen. *The Definition of Standard ML (Revised)*. The MIT Press, Cambridge, Mass., 1997.
- [11] F. Pessaux and X. Leroy. Type-based analysis of uncaught exceptions. In *Conference Record of the Twenty-Sixth Annual ACM Symposium on Principles of Programming Languages*, pages 276–290. ACM Press, 1999.
- [12] D. Rémy. Records and variants as a natural extension of ML. In *Conference Record of the Sixteenth Annual ACM Symposium on Principles of Programming Languages*, pages 77–88. ACM Press, 1989.
- [13] J. H. Reppy. A safe interface to sockets. Technical memorandum, AT&T Bell Laboratories, 1996.
- [14] M. Wand. Complete type inference for simple objects. In *Proceedings of the 2nd Annual IEEE Symposium on Logic in Computer Science*, 1987.