

Modular Type Classes

Derek Dreyer

Toyota Technological Institute at Chicago
dreyer@tti-c.org

Robert Harper

Carnegie Mellon University
rwh@cs.cmu.edu

Manuel M.T. Chakravarty

University of New South Wales
chak@cse.unsw.edu.au

Abstract

ML modules and Haskell type classes have proven to be highly effective tools for program structuring. Modules emphasize explicit configuration of program components and the use of data abstraction. Type classes emphasize implicit program construction and *ad hoc* polymorphism. In this paper, we show how the implicitly-typed style of type class programming may be supported within the framework of an explicitly-typed module language by viewing type classes as a particular mode of use of modules. This view offers a harmonious integration of modules and type classes, where type class features, such as class hierarchies and associated types, arise naturally as uses of existing module-language constructs, such as module hierarchies and type components. In addition, programmers have explicit control over which type class instances are available for use by type inference in a given scope. We formalize our approach as a Harper-Stone-style elaboration relation, and provide a sound type inference algorithm as a guide to implementation.

Categories and Subject Descriptors D.3.1 [Programming Languages]: Formal Definitions and Theory; D.3.3 [Programming Languages]: Language Constructs and Features—Polymorphism, Modules; F.3.3 [Logics and Meanings of Programs]: Studies of Program Constructs—Type structure

General Terms Design, Languages, Theory

Keywords Type classes, modules, type inference, type systems

1. Introduction

The ML module system [17] and the Haskell type class system [23, 19] have proved, through more than 15 years of practical experience and theoretical analysis, to be effective linguistic tools for structuring programs. Each provides the means of *specifying* the functionality of program components, *abstracting* programs over such specifications, and *instantiating* programs with specific realizations of the specifications on which they depend. In ML such specifications are called *signatures*, abstraction is achieved through *functors*, and instantiation is achieved by *functor application* to *structures* that implement these signatures. In Haskell such specifications are called *type classes*, abstraction is achieved through *constrained polymorphism*, and instantiation is achieved through

polymorphic instantiation with *instances* of type classes. There is a clear correspondence between the highlighted concepts (see [24]), and consequently modules and type classes are sometimes regarded as opposing approaches to language design. We show that there is no opposition. Rather, type classes and modules are complementary aspects of a comprehensive framework of modularity.

Perhaps the most significant difference is the mode of use of the two concepts. The Haskell type class system is primarily intended to support *ad hoc* polymorphism in the context of a parametrically polymorphic language. It emphasizes the *implicit* inference of class constraints and *automatic* construction of instances during overload resolution, which makes it convenient to use in many common cases, but does not facilitate more general purposes of modular programming. Moreover, the emphasis on automatic generation of instances imposes inherent limitations on expressiveness—most importantly, there can be at most one instance of a type class at any particular type.

In contrast, the ML module system is designed to support the structuring of programs by forming hierarchies of components and imposing abstraction boundaries—both *client-side* abstraction, via functors, and *implementor-side* abstraction, via signature ascription (aka *sealing*). The module system emphasizes *explicit* manipulation of modules in the program, which makes it more flexible and general than the type class mechanism. Modules may be ascribed multiple signatures that reveal varying amounts of type information, signatures may be implemented by many modules, and neither modules nor signatures are restricted to have the rigid form that Haskell’s instances and classes have. On the other hand, ML lacks support for *implicit* module generation and *ad hoc* polymorphism, features which experience with Haskell has shown to be convenient and desirable.

There have been many proposals to increase the expressiveness of the original type class system as proposed by Wadler and Blott [23], including constructor classes [15], functional dependencies [12], named instances [16], and associated types [2, 1]. These may all be seen as adding functionality to the Haskell class system that mirrors aspects of the ML module system, while retaining the implicit style of usage of type classes. However, these (and other) extensions tend to complicate the type class system without alleviating the underlying need for a more expressive module system.

In fact, there are ways in which the Haskell type class mechanism impedes modularity. To support implicit instance generation while ensuring coherence of inference, Haskell insists that instances of type classes be drawn from a global set of instance declarations; in particular, instances are implicitly exported and imported, which puts their availability beyond programmer control. This can be quite inconvenient—for many type classes there is more than one useful instance of the class at a particular type, and the appropriate choice of instance depends on the context in which an overloaded operator is used. Hence, the Haskell Prelude must provide many functions in two versions: one using type classes and the other an explicit function argument—e.g., `sort` and `sortBy`.

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.

POPL’07 January 17–19, 2007, Nice, France.

Copyright © 2007 ACM 1-59593-575-4/07/0001...\$5.00.

In this paper we take a different tack. Rather than bolster the expressiveness of type classes, we instead propose that a more sensible approach to combining the benefits of type classes and modules is to *start* with modules as the fundamental concept, and *then* recover type classes as a particular mode of use of modularity. We retain the advantages of a fully expressive *explicit* module system, while also offering the conveniences of *implicit* type class programming, particularly the integration of *ad hoc* and parametric polymorphism. Moreover, the proposed design provides a clean separation between the *definition* of instances and their *availability for use* during type inference. This offers localization of instance scoping, enhanced readability, and the potential for instances to be compiled separately from their uses. The result is a harmonious integration of modules and type classes that provides the best features of both approaches in a single, consistent framework. The elegance of our approach stems from the observation that type class features, such as class hierarchies and associated types, arise naturally as uses of existing module-language constructs, such as module hierarchies and type components.

In summary, this paper makes the following contributions:

- We present a smooth integration of type classes and modules that provides a foundation for future work on incorporating type classes into ML and a proper module system into Haskell. We give an intuition of the integration of type classes into ML in Section 2.
- We highlight some interesting design issues that arose while developing the interpretation of type classes in terms of modules (Section 3).
- We specify the semantics of an extended module language that supports type classes. We formalize its elaboration (in the style of Harper and Stone [9]) into an explicitly-typed module type system. We also generalize Damas and Milner’s Algorithm \mathcal{W} [3] to an inference algorithm for modular type classes that we have proved sound with respect to the elaboration semantics. Due to space limitations, Section 4 only sketches the most salient and unusual features of our formalization. For the full formalization of our language, together with expanded technical discussion and theorem statements, we refer the reader to our companion technical report [7].

Our elaboration translation demonstrates that modules can serve as *evidence* in the sense of Jones [14]. Compared to the customary use of dictionary records as evidence, modules offer a cleaner way of handling extensions to the basic type class mechanism such as associated types. In addition, for the application to type classes, the use of modules as evidence makes clear that the construction of evidence respects the phase distinction [8], *i.e.*, it is based solely on compile-time information, not run-time information. We conclude in Section 5 with further discussion of related work.

2. Modular Type Classes: An Overview

In this section we summarize our approach to representing the main mechanisms of a Haskell-style type class system within the context of an ML-style module system. For readability, we employ ML-like syntax for our examples, although the formal design we describe later is syntactically more austere and leaves a number of (largely superficial) aspects of an actual ML extension to future work.

2.1 Classes are signatures, instances are modules

A type class in Haskell is essentially an interface describing a set of operations whose types mention a distinguished abstract type variable known as the *class parameter*. It is natural therefore to represent a class in the module setting as a signature (*i.e.*, an interface) with a distinguished type component (representing the class param-

eter). In particular, we insist that the distinguished type component be named “*t*”. It may be followed by any number of other type, value, or substructure components. We call such a signature a *class signature*, specifically an *atomic* class signature (in contrast to the *composite* ones that we describe below in Section 2.3.) For example, the class of equality types is represented by the atomic class signature EQ, defined as follows:

```
signature EQ = sig
  type t
  val eq : t * t -> bool
end
```

Note that class signatures like EQ are just ordinary ML signatures of a certain specified form.

Correspondingly, an instance of a type class is represented by a module. A monomorphic instance of a type class is represented by a structure, and a polymorphic instance is represented by a functor. For example, we can encode an `int` instance of the equality class as a structure whose signature is EQ where `type t = int`:

```
structure EqInt = struct
  type t = int
  val eq = Int.eq
end
```

As in Haskell, the instance for a compound type $t(t_1, \dots, t_n)$ is composed from instances of its component types, t_1, \dots, t_n , by a functor, `Eqt`, associated with its outermost type constructor, t . For example, here is an instance of equality for product types $t_1 * t_2$:

```
functor EqProd (X : EQ, Y : EQ) = struct
  type t = X.t * Y.t
  fun eq ((x1,y1), (x2, y2)) =
    X.eq(x1,x2) andalso Y.eq(y1,y2)
end
```

There is an evident correspondence with Haskell instance declarations, but rather than use Horn clause logic programs to specify closure conditions, we instead use functional programs (in the form of functors).

From the `EqInt` and `EqProd` modules we can construct an instance, say, of signature EQ where `type t=int*int`:

```
structure EqII = EqProd(EqInt,EqInt)
```

Of course, one of the main reasons for using type classes in the first place is so that we don’t have to write this functor application manually—it corresponds to the process known as dictionary construction in Haskell and can be performed automatically, behind the scenes, during type inference. In particular, such automatic functor application may occur in the elaboration of expressions that appear to be values, such as when a variable undergoes polymorphic instantiation (see below). Consequently, it is important that the application of an instance functor does not engender any computational effects, such as I/O or non-termination. We therefore require that instance functors be *total* in the sense that their bodies satisfy something akin to ML’s value restriction. This restriction appears necessary in order to ensure predictable program behavior.

2.2 Separating the definition of an instance from its use

In Haskell, an instance becomes immediately available for use by the type inference engine as soon as it is declared. As a consequence, due to the implicit global importing and exporting of instances, there can only ever be a single instance of a class at a certain type in one program. This is often a nuisance and leads to awkward workarounds. Proposals such as named instances [16] have attempted to alleviate this problem, but have not been generally accepted.

In contrast, our reconstruction of type classes in terms of modules provides a natural solution to this dilemma. Specifically, we require that an instance module only become available for use by the inference engine after it has been nominated for this purpose explicitly by a `using` declaration. This separates the *definition* of an instance from its *adoption* as a *canonical instance*, thus facilitating modular decomposition and constraining inference to make use only of a clearly specified set of instances. For example, the declaration

```
using EqInt, EqProd in mod
```

nominates the two instance modules defined earlier as available for canonical instance generation during elaboration of the module `mod`. The typing rule for `using` demands that `EqInt` and `EqProd` not overlap with any instances that have already been adopted as canonical. (A precise definition of overlapping instances is given in Section 3.2.)

In both our language and Haskell, canonical instance generation is implicitly invoked whenever overloading is resolved. In our language, we additionally provide a mechanism `canon(sig)` by which the programmer can explicitly request the canonical instance module implementing the class signature `sig`.¹ At whatever point within `mod` instance generation occurs, it will employ only those instances that have been adopted as canonical in that scope.

2.3 Class hierarchies via module hierarchies

In Haskell, one can extend a class *A* with additional operations to form a class *B*, at which point *A* is called a *superclass* of *B*. Class hierarchies arise in the module setting naturally from module hierarchies. This is easiest to illustrate by example.

Suppose we want to define a class called `ORD`, which extends the `EQ` class with a `lt` operation. We can do this by first defining an atomic class `LT` that only supports `lt`, and then defining `ORD` as a *composite* of `EQ` and `LT`:

```
signature ORD = sig
  structure E : EQ
  structure L : LT
  sharing type E.t = L.t
end
```

The sharing specification makes explicit that `ORD` is providing two different interpretations of the *same* type, as an equality type and as an ordered type. `ORD` is an example of what we call a *composite class signature*, i.e., a signature consisting of a collection of atomic signatures bound to submodules whose names are arbitrary.

Instances of composite class signatures are not written by the programmer directly, but rather are composed automatically by the inference engine from the instances for their atomic signature parts. For example, if we want to write instances of `ORD` for `int` and the `*` type constructor, what we do instead is to write instances of `LT`:

```
structure LtInt = struct
  type t = int
  val lt = Int.lt
end
functor LtProd (X : ORD, Y : LT) = struct
  type t = X.E.t * Y.t
  fun lt ((x1,y1), (x2,y2)) =
    X.L.lt(x1,x2) orelse
    (X.E.eq(x1,x2) andalso Y.L.lt(y1,y2))
end
```

Note that `LtProd` requires its first argument to be an instance of `ORD`, not `LT`. This is because the implementation of `lt` in the body

¹This feature is particularly useful in conjunction with our support for associated types; see Section 2.5.

of the functor depends on having both equality and ordering on the type `X.E.t` so that it can implement a lexicographic ordering on `X.E.t * Y.t`. For `Y.t`, only the `lt` operation is needed.

Now, let us assume these instances are made canonical (via the `using` declaration) in a certain scope. Then, during typechecking, if the inference engine demands a canonical module of signature `ORD` where `type E.t = int * int`, it will be computed to be

```
struct
  structure E = EqProd(EqInt,EqInt)
  structure L = LtProd(struct
    structure E = EqInt
    structure L = LtInt
  end,
  LtInt)
end
```

The fundamental reason that we do not allow instances for `ORD` to be adopted directly is that we wish to prevent the instances for `ORD` from having any overlap with existing instances that may have been adopted for `EQ`. If one were to define an instance for `ORD` where `type E.t = int` directly, one would implicitly provide an instance for `EQ` where `type t = int` through its `E` substructure; and if one tried to adopt such an `ORD` instance as canonical, it would overlap with any existing canonical instance of `EQ` where `type t = int`.

Under our approach, this sort of overlap is avoided. Moreover, the code one writes is ultimately very similar to the code one would write in Haskell (except that it is expressed entirely in terms of existing ML constructs). In particular, the instance declaration for `ORD` at `int` in Haskell is only permitted to provide a definition for the new operations (namely, `lt`) that are present in `ORD` but not in `EQ`. In other words, an instance declaration for `ORD` in Haskell is precisely what we would call an instance of `LT`.

2.4 Constrained polymorphism via functors

Under the Harper-Stone interpretation of Standard ML (hereafter, HS) [9], polymorphic functions in the *external* (source) language are elaborated into *functors* in an *internal* module type system. Specifically, a polymorphic value is viewed as a functor that takes a module consisting only of type components (representing the polymorphic type variables) as its argument and returns a module consisting of a single value component as its result.

The HS semantics supports the concept of *equality polymorphism* found in Standard ML by simply extending the class of signatures over which polymorphic functions may be abstracted to include the `EQ` signature defined above. For example, in the internal module type system of HS, the *ad hoc* polymorphic equality function is represented by the functor

```
functor eq (X:EQ) :> [X.t * X.t -> bool] = [X.eq]
```

where the brackets notation describes a module with a single value component. Polymorphic instantiation at a type τ consists of computing a canonical instance of `EQ` where `type t = τ` , as described above, applying the functor `eq` to it, and extracting the value component of the resulting module.

The present proposal is essentially a generalization of the HS treatment of equality polymorphism to arbitrary type classes. A functor that abstracts over a module representing an instance of a type class is reminiscent of the notion of a *qualified type* [11], except that we make use of the familiar concept of a functor from the ML module system, rather than introduce a new mechanism solely to support *ad hoc* polymorphism.

Of course, the programmer need not write the `eq` functor manually. Our external language provides an `overload` mechanism, and

the elaborator will generate the above functor automatically when the programmer writes

```
val eq = overload eq from EQ
```

Note that there is no need to bind the polymorphic function returned by the `overload` mechanism to the name `eq`; it can be called anything. In practice, it may be useful to be able to overload all the components of a class signature at once by writing `overload SIG` as syntactic sugar for a sequence of `overload`'s for the individual components of the signature.

The following are some examples of elaboration in the presence of the overloaded `eq` function:

```
using EqInt, EqProd in ...eq((2,3),(4,5))...
~> ...Val(eq(EqProd(EqInt,EqInt))) ((2,3),(4,5))...
```

```
fun refl y = eq(y,y)
~> functor refl (X : EQ) :> [[X.t -> bool]]
    = [fn y => Val(eq(X)) (y,y)]
```

(Note: the `Val` operator seen here is the mechanism in our internal module type system by which a value of type τ is extracted from a module of signature $[[\tau]]$.)

Our language also allows for the possibility that the programmer may wish to work with explicitly polymorphic functions in addition to implicit overloaded ones. In particular, by writing

```
functor Refl = explicit (refl :
  (X : EQ) -> sig val it : X.t -> bool end)
```

we convert the polymorphic function `refl` into an explicit functor `Refl`. The programmer can then apply it to an arbitrary module argument of signature `EQ` and project out the `it` component of the result. The reason we require a signature annotation on the `explicit` construct is that the implicitly-typed `refl` may be declaratively ascribed many different signatures. Whenever `refl` is used, type inference will compute the appropriate instance arguments for it regardless of the particular signature it has been ascribed. However, since it is the *programmer* who applies `Refl`, she needs to know exactly what shape `Refl`'s module argument is expected to have.

We also provide an `implicit` construct to coerce explicit functors into implicit ones. (See the technical report [7] for details.)

2.5 Associated types arise naturally

The experience with type classes in Haskell quickly led to the desire for type classes with more than one class parameter. However, these multi-parameter type classes are not generally very useful unless dependencies between the parameters can be expressed. This led in turn to the proposal of functional dependencies [12] and more recently associated types [2, 1] for Haskell.

An associated type is a type component provided by a class that is not the distinguished type component (class parameter). The associated types of a class do not play a role in determining the canonical instance of a class at a certain type—that is solely determined by the identity of the distinguished type.

Modular type classes immediately support associated types as additional type components of a class signature. An illustrative example is provided by a class of collection types:

```
signature COLLECTS = sig
  type t
  type elem
  val empty : t
  val insert : elem * t -> t
  val member : elem * t -> bool
  val toList : t -> elem list
end
```

The distinguished type `t` represents the collection type and the associated type component `elem` represents the type of elements. An instance for lists, where the elements are required to support equality for the membership test, would be defined as follows:

```
functor CollectsList (X : EQ) = struct
  type t = X.t list
  type elem = X.t
  val empty = []
  fun insert (x, L) = x::L
  fun member (x, []) = raise NotInCollection
    | member (x, y::L) = X.eq (x,y) orelse
      member (x,L)
  fun toList L = L
end
```

When using classes with associated types, it is common to need to place some constraints on the identities of the associated types. For example, suppose we write the following:

```
val toList = overload toList from COLLECTS
fun sumColl C = sum (toList C)
```

The `sumColl` function does not care what type of collection `C` is, so long as its element type is `int`. Correspondingly, the elaborator will assign `sumColl` the polymorphic type (*i.e.*, functor signature)

```
(X : COLLECTS where type elem = int) -> [[X.t -> int]]
```

Note that the constraint on the type `X.elem` is expressed completely naturally using ML's existing `where` type mechanism, which is just syntactic sugar for the transparent realization of an abstract type component in a signature. In contrast, the extension to handle associated type synonyms in Haskell [1] requires an additional mechanism called *equality constraints* in order to handle functions like `sumColl`.

As Chakravarty *et al.* [1] have demonstrated, it is useful in certain circumstances to be able to compute (statically) the identity of an associated type `assoc` in the canonical instance of a type class `SIG` at a given type τ . This is achieved in our setting via the `canon(sig)` construct, which we introduced above as a way of explicitly computing a canonical instance. In particular, we can write

```
canon(SIG where type t =  $\tau$ ).assoc
```

which constructs the canonical instance of `SIG` at τ and then projects the `assoc` type from it.² In the associated type extension to Haskell, one would instead write `assoc(τ)`.

While the ML syntax here is clearly less compact, there is a good reason for it. Specifically, the Haskell syntax only makes sense because Haskell ties each associated type name in the program to a single class (in this case, `assoc` would be tied to `SIG`). In contrast, in our setting, it is fine for several different class signatures to have an associated type component called `assoc`.

3. Design Considerations

In this section we examine some of the more subtle points in the design of modular type classes and explain our approach to handling them.

3.1 Coherence in the presence of scoped instances

The `using` mechanism described in the introduction separates the definition of instance functors from their adoption as canonical

²Note that, due to the principle of *phase separation* in the ML module system [8], the identity of the `assoc` type here can be determined purely statically, and elaboration does not actually need to construct the dynamic parts of `canon(SIG where type t = τ)`.

instances. It also raises questions of coherence stemming from the nondeterministic nature of polymorphic type inference. Suppose `EqInt1` and `EqInt2` are two observably distinct instances of `Eq` where `type t = int`. Consider the following code:

```
structure A = using EqInt1 in
  struct ...fun f x = eq(x,x)... end
structure B = using EqInt2 in
  struct ...val y = A.f(3)... end
```

The type inference algorithm is free to resolve the meaning of this program in two incompatible ways. On the one hand, it may choose to treat `A.f` as polymorphic over the class `Eq`; in this case, the application `A.f(3)` demands an instance of `Eq` where `type t=int`, which can only be resolved by `EqInt2`. On the other hand, type inference is free to assign the type `int -> bool` to `A.f` at the point where `f` is defined, in which case the demand for an instance of `Eq` can only be met by `EqInt1`. *These are both valid typings, but they lead to observably different behavior.*

An unattractive solution is to insist on a specific algorithm for type inference that arbitrarily chooses one resolution over another, but this sacrifices the elegant, declarative nature of a Hindley-Milner-style type system and, worse, imposes a specific resolution policy that may not be desired in practice. Instead, we prefer to take a different approach, which is to put the decision under programmer control, permitting either outcome at her discretion. We could achieve this by insisting that the scope of a `using` declaration be given an explicit signature, so that in the above example the programmer would have to specify whether `A.f` is to be polymorphic or monomorphic. However, this approach is awkward for nested `using` declarations, forcing repeated specifications of the same information.

Instead we propose that the `using` declaration be confined to an *outer* (or *top-level*) layer that consists only of module declarations, whose signatures are typically specified in any case. All core-level terms appear in the *inner* layer, where type inference proceeds without restriction, but no `using` clauses are allowed. Thus, the set of permissible instances is fixed in any inner context, but may vary across outer contexts. At the boundary of the two layers, a type or signature annotation is required. This ensures that the scope of a `using` declaration is explicitly typed without effecting duplication of annotations. The programmer who wishes to ignore type classes simply confines herself to the inner level, with no restrictions; only the use of type classes demands attention be paid to the distinction.

3.2 Overlapping instances

To ensure coherence of type inference, the set of available instances in any context must be *non-overlapping*. Roughly speaking, this means that there should only be one way to compute the canonical instance of any given class at any given type. There is considerable leeway, though, in determining the precise definition of overlap, and indeed this remains a subject of debate in the Haskell community. For the purposes of this paper we follow the guidelines used in Haskell 98. In particular, we insist that there be one instance per type constructor, so that instance resolution proceeds by a simple inductive analysis of the structure of the instance type, composing instance functors to obtain the desired result.

However, in the modular approach suggested here, there is an additional complication. Just as a module may satisfy several different signatures, so a single module may qualify as an instance of several different type classes. For example, the module

```
struct type t = int; fun f(x:t) = x end
```

may be seen as an instance of the class

```
sig type t; val f : t -> t end
```

and also of the class

```
sig type t; val f : t -> int end.
```

Thus, to check if two instances `A` and `B` (with the same `t` component) are non-overlapping, we need to ensure that the set of *all* classes to which `A` could belong is disjoint from the set of *all* classes to which `B` could also belong.

A simple, but practical, criterion to ensure this is to define two instances to be non-overlapping iff either (1) they differ on their distinguished `t` component, so that no overlap is possible, or (2) in the case that they have the same `t` component, that they be *structurally dissimilar*, which we define to mean that their components do not all have the same names and appear in the same order. While other, more refined definitions are possible, we opt here for simplicity until evidence of the need for a more permissive criterion is available.

3.3 Unconstrained type components in class signatures

In order to support ordinary ML-style polymorphism, we need a way to include unconstrained type components in a class signature. We could use the class signature `sig type t end` for this purpose. However, since our policy is that the only canonical instances of atomic class signatures are those that have been explicitly adopted as canonical by a `using` declaration, this would amount to treating `sig type t end` as a special case.

We choose instead to allow composite class signatures to contain arbitrary unconstrained type components, so long as they are named something other than `t`. For example, under our approach, the divergent function

```
fun f x = f x
```

can be assigned the polymorphic type

```
(X : sig type a; type b end) -> [X.a -> X.b]
```

(The choice of the particular names `a` and `b` here is arbitrary.)

In our formal system, we refer to the union of the `t` components and the unconstrained components of a class signature `S` as the *parameters* of `S`.

3.4 Multi-parameter and constructor classes

Two extensions to Wadler & Blott's [23] type class system that have received considerable attention are *multi-parameter type classes* and *constructor classes*. We have chosen not to cover these extensions in this paper. Concerning multi-parameter classes, most uses of them require functional dependencies [12], which when rewritten to use associated types (which we support), turn into single-parameter classes. Hence, we expect the need for multi-parameter classes to be greatly diminished in our case.

As for constructor classes, we see no fundamental problems in supporting them in an extension of our framework since type components of ML modules may have higher kind. However, we view them as an orthogonal extension, and thus have opted to omit them in the interest of a clearer and more compact presentation.

4. Formal System

In this section, we will give a brief sketch of our type-theoretic formalization of modular type classes, highlighting its most distinctive features. For full details, see the companion technical report [7].

4.1 Declarative elaboration semantics

Following Harper and Stone [9], we define our language of modular type classes using an *elaboration semantics*, in which *external language* (EL) source programs are interpreted by translation into an *internal language* (IL) type system. The elaboration translation is

syntax-directed, but it is also nondeterministic with respect to polymorphic generalization and instantiation. This style of definition is the standard method of giving meanings to programs involving type classes, although in the context of Haskell it is often referred to as *evidence translation* [14].

The IL we use is a simplified variant of the type system for modules given in Dreyer’s thesis [4], which in turn is based on the higher-order module calculus of Dreyer, Cray and Harper [6]. For defining the semantics of type classes, the most salient feature of this IL is that it distinguishes between two kinds of functors (and functor signatures): *total* and *partial*. Total functor signatures, written $\forall X:S_1.S_2$, classify functors with argument signature S_1 and result signature S_2 , whose bodies are judged syntactically to be free of computational effects. Partial functor signatures, written $\Pi X:S_1.S_2$, classify functors whose bodies may contain effects. In general, since ML is not purely functional, ML functors may be partial. However, as we explained in Section 2.1, we require that instance functors be total in order to ensure predictable program behavior at points of polymorphic instantiation. (A technical aside: the notation S_2 for the result signature of a total functor indicates that it is required syntactically to be transparent. This restriction is demanded by Dreyer *et al.*’s treatment of data abstraction as a computational effect, but it is in no way a hindrance—functors corresponding to Haskell instances are naturally transparent.)

As for our external language, we have already described most of its novel constructs informally in Section 2. One feature of our EL that we have not discussed is its mechanism for inducing polymorphic generalization. Traditionally, generalization is performed implicitly as part of typechecking a term-level let construct, let $x=exp_1$ in exp_2 (hence the name *let-polymorphism*). After typechecking exp_1 , the principal type scheme of exp_1 is generalized into a polymorphic type (or *polytype*), to which x is bound during the typechecking of exp_2 .

As explained in Section 2.4, polymorphic types are modeled in our language as a special case of functor signatures in which the argument has a class signature and the result signature specifies a single value component. Thus, instead of tying generalization to let, we opt instead to induce it via an *orthogonal* construct $[exp]$ that coerces an EL core term exp into a module, thereby generalizing its monotype (*i.e.*, type) into a polytype (*i.e.*, functor signature). Likewise, polymorphic instantiation occurs implicitly when a module path P —a module variable X followed by zero or more component projections—is used as a core-language expression. Under this approach, traditional let-polymorphism is modeled as a composition of the generalizing $[exp]$ and a non-generalizing let-construct let $X=mod$ in exp . In particular, the let-polymorphic let $x=exp_1$ in exp_2 is encodable as let $X=[exp_1]$ in $\{x \mapsto X\} exp_2$.

Following Harper and Stone, the main translation judgments in our elaboration semantics all have the form

$$\Theta; \Gamma \vdash EL\text{-phrase} \rightsquigarrow IL\text{-phrase} : IL\text{-classifier}$$

Here, *EL-phrase* and *IL-phrase* range over EL and IL modules, terms and type constructors, and *IL-classifier* ranges correspondingly over IL signatures, types and kinds. We design the inference rules so that the output *IL-phrase* is guaranteed to have the output *IL-classifier* in the IL type system. The context Γ may contain bindings of type variables α to kinds K , term variables x to types τ , and module variables X to signatures S . The novel element in this judgment form is Θ , which is a set of paths to structures and functors that are to be considered canonical instances within *EL-phrase*. We call Θ the canonical instance set.

Most of the rules in our elaboration semantics are similar to corresponding rules in the Harper-Stone semantics, and thus do not interact with the canonical instance set Θ . There are three major rules that do. One is the rule for the using mechanism, which has

the effect of adding a given path to Θ (under the condition that it does not overlap with any instance modules already in Θ). Most of the work in formalizing this rule is in specifying what it means for two instances to overlap. See the technical report for details [7].

Two other rules that interact with Θ are those for polymorphic generalization and instantiation. Rule 1 formalizes the polymorphic instantiation that occurs when a module path P is used as a term:

$$\frac{\Gamma \vdash P : \forall X:S.[\tau] \quad \Gamma \vdash S \leq S \quad \Theta; \Gamma \vdash_{\text{can}} \mathbb{V} : S}{\Theta; \Gamma \vdash P \rightsquigarrow \text{Val}(P(\mathbb{V})) : \tau[\mathbb{V}/X]} \quad (1)$$

The first premise checks that P is in fact a polymorphic value (represented as a total functor). Instantiation then consists of finding the canonical instance module of the class signature S to which P will be applied. Since the parameters of S are abstract, the choice of which instance module is nondeterministic. Consequently, the second premise picks a transparent signature \mathbb{S} that is a *subtype* of S , meaning that it realizes the parameters of S with some choices τ_1, \dots, τ_n . (Signature subtyping, a common judgment in module type systems, is defined formally in [7].) Lastly, the third premise computes the canonical IL module \mathbb{V} of signature \mathbb{S} using the *canonical module* judgment $\Theta; \Gamma \vdash_{\text{can}} \mathbb{V} : \mathbb{S}$. Note that all of this is done in terms of module and signature judgments, without ever explicitly mentioning the instantiating types τ_1, \dots, τ_n !

The canonical module judgment $\Theta; \Gamma \vdash_{\text{can}} \mathbb{V} : \mathbb{S}$ is straightforward to define. In short, a composite instance module is canonical if all its atomic instance components are canonical; an atomic instance module is canonical if it is either a canonical instance structure (from the set Θ) or the result of applying a canonical instance functor from Θ to a canonical argument. Canonical modules may also contain arbitrary unconstrained type components (named something other than τ , as per the discussion in Section 3.3).

Rule 2 formalizes polymorphic generalization for $[exp]$:

$$\frac{X \notin \text{FV}(exp) \quad \Gamma \vdash_{\text{class}} X : S \rightsquigarrow \Theta' \quad \Theta, \Theta'; \Gamma, X : S \vdash exp \rightsquigarrow v : \tau}{\Theta; \Gamma \vdash [exp] \rightsquigarrow \Lambda X : S.[v] : \forall X : S.[\tau]} \quad (2)$$

One can view this rule as “guessing” a polymorphic type $\forall X : S.[\tau]$ to assign to exp . Suppose that S is an atomic class signature like EQ. In order to see whether exp can be elaborated with this type, we add the class constraint $X : S$ to the context and make it a canonical instance of the signature S where type $\tau = X.t$ (by adding X to Θ) before typechecking exp . The last step is critical: if X is *not* added to Θ , then the canonical module judgment will have no way of knowing that X is the canonical module of signature S where type $\tau = X.t$ at polymorphic instantiation time.

However, in the case that S is a composite class, the elaborator does not permit X to be added directly to the instance set Θ . To simplify the formalization of other judgments, we require all the instance structures in Θ to have atomic signature. Thus, in general we need a way of parsing the class constraint $X : S$ in order to produce a *set* of paths Θ' (all of which are rooted at X) that represent the atomic instance modules contained within X .

This class parsing is achieved via the *class elaboration* judgment $\Gamma \vdash_{\text{class}} X : S \rightsquigarrow \Theta'$ used in the second premise of Rule 2. For example, if S were the composite class ORD from Section 2.3, then Θ' would be the set $\{X.E, X.L\}$. In the case that there are multiple paths in X to atomic instances of the same signature, Θ' will include exactly one. (It doesn’t matter which one, since X represents a canonical instance module, and the \vdash_{can} judgment guarantees that any two submodules of a canonical module that have the same transparent signature must be the same module value.)

The class elaboration judgment also checks that S is a valid class signature. A signature is considered a valid class signature if it is a collection of unconstrained type components and atomic instance components (whose first component is τ), in which the

unconstrained and t components—*i.e.*, the *parameters* of S —are all abstract (although possibly, as in the case of ORD, subject to type sharing constraints). The requirement that the parameters of S be abstract ensures that the instances in the set Θ' all concern abstract type components of the freshly chosen variable X , which in turn guarantees that the instances in Θ' do not overlap with any instances in the input instance set Θ .

4.2 Type inference algorithm

The elaboration semantics sketched above is nondeterministic, and hence is not directly implementable without backtracking. In order to guide implementation, we therefore also provide a type inference algorithm in the style of Algorithm \mathcal{W} [3]. This section describes some highlights of our algorithm.

Following Damas and Milner, we thread through the inference rules a substitution δ whose domain consists of *unification variables*, denoted by bold α . In addition, polymorphic instantiation in the presence of type classes generates *constraints*, which we denote Σ . Constraints are sets of $X:S$ bindings, in which the X 's do not appear free in the S 's. Each $X:S$ represents a demand generated by the algorithm for a canonical module of signature \mathbb{S} to be substituted for X in the term or module that is output by elaboration.

For example, the inference judgment for terms has the form $\Theta; \Gamma \vdash \text{exp} \Rightarrow e : \tau / (\Sigma; \delta)$, where everything to the left of the \Rightarrow is input to the algorithm and everything to the right of the \Rightarrow is output. Rule 3 is the inference rule for polymorphic instantiation:

$$\frac{\Gamma \vdash P : \downarrow \forall X:S. [\tau] \quad S \Rightarrow \exists \bar{\alpha}. \mathbb{S} \quad \Gamma, X:S \vdash \tau \downarrow \tau'}{\Theta; \Gamma \vdash P \Rightarrow \text{Val}(P(X)) : \tau' / (X:S; \text{id})} \quad (3)$$

Given a path P of polymorphic signature $\forall X:S. [\tau]$, the second premise uses the auxiliary judgment $S \Rightarrow \exists \bar{\alpha}. \mathbb{S}$ to generate fresh unification variables $\bar{\alpha}$ corresponding to the abstract type components of S . It then applies P to an unknown canonical module X of signature \mathbb{S} , and projects out the value component. This in turn effects a demand for $X:S$ in the output constraint. For example, if S were the class EQ, then the output constraint would be $X : \text{EQ where type } t = \alpha$. (Note: the “ \downarrow ” judgment used in the first premise indicates $\forall X:S. [\tau]$ is the normal form signature of P , and the last premise normalizes τ so that references to type components of X become references to the corresponding $\bar{\alpha}$.)

As type inference uncovers the identity of certain unification variables, it becomes possible (and at certain points necessary) to eliminate some of the constraints amassed in Σ through a process we call *constraint normalization*. This process takes zero or more steps of *constraint reduction* until the input constraint has been converted to a normal form in which all residual constraints are instances of atomic classes at unification variables. The normalization judgment has the form $\Theta; \Gamma \vdash \Sigma_1 \downarrow (\Sigma_2; \sigma; \delta)$, where σ is a module substitution whose domain is that of Σ_1 . The relation between the input and output of normalization is summarized by the following invariant:

$$\text{If } \Theta; \Gamma \vdash \Sigma_1 \downarrow (\Sigma_2; \sigma; \delta), \\ \text{then } \forall X:S \in \Sigma_1. \Theta, \text{dom}(\Sigma_2); \delta\Gamma, \Sigma_2 \vdash_{\text{can}} \sigma X : \delta S.$$

That is, if we treat the domain of the normalized constraint Σ_2 as a set of canonical instances, then from those instances together with the canonical instances already in Θ , the substitution σ shows how to construct canonical modules to satisfy all the demands of the original constraint Σ_1 (subject to type substitution δ).

To make this concrete, suppose Θ contains the EqInt and EqProd instance modules given in Section 2.1, and suppose that Σ_1 is $X : \text{EQ where type } t = \text{int} * \alpha$. Then the normalized Σ_2 would be $Y : \text{EQ where type } t = \alpha$, and the substitution σ would map X to EqProd(EqInt, Y). (In this case, δ would simply be the identity substitution id .)

The constraint normalization algorithm may be viewed as a backchaining implementation of the canonical module judgment. It is essentially a combination of Haskell-style *context reduction* (aka *simplification*) and *constraint improvement* [10], except that it is formalized entirely in terms of modules and signatures.

4.3 Soundness

We have proven that our inference algorithm is sound with respect to the elaboration semantics. For space reasons, we collect the main results here and refer the reader to the companion technical report [7] for the full statement of the soundness theorem and its auxiliary definitions, including the precise meaning of the theorem’s preconditions.

Theorem (Soundness)

Suppose $(\Theta; \Gamma)$ is valid for inference, $\Theta' \supseteq \Theta$, $\Gamma' \vdash \delta' : \delta\Gamma$, $\vdash (\Theta'; \Gamma')$ ok, and $\forall X:S \in \Sigma. \Theta'; \Gamma' \vdash_{\text{can}} \sigma' X : \delta' S$. Then:

1. If $\Theta; \Gamma \vdash \text{exp} \Rightarrow e : \tau / (\Sigma; \delta)$,
then $\Theta'; \Gamma' \vdash \text{exp} \rightsquigarrow \sigma' \delta' e : \delta' \tau$.
2. If $\Theta; \Gamma \vdash \text{mod} \Rightarrow M : S / (\Sigma; \delta)$,
then $\Theta'; \Gamma' \vdash \text{mod} \rightsquigarrow \sigma' \delta' M : \delta' S$.

Consider part 1. Informally, Θ , Γ and exp are inputs. If type inference on exp succeeds, it produces an IL term e , along with a constraint Σ and a substitution δ . If in any “future world” $(\Theta'; \Gamma')$ the constraint Σ can be solved by substitutions σ' and δ' , then exp will declaratively elaborate to $\sigma' \delta' e$ in that world. The theorem statement for modules (part 2) is analogous.

4.4 Incompleteness

While the inference algorithm is sound, it is *not* complete, for reasons that arise independently of the present work. One source of incompleteness is inherited from Haskell and concerns a fundamental problem with type classes, namely the problem of *ambiguity* [14]. The canonical example uses the following two signatures:

```
signature SHOW = sig
  type t
  val show : t -> string
end
signature READ = sig
  type t
  val read : string -> t
end
val show = overload show from SHOW
val read = overload read from READ
```

Given this overloading, the expression `show (read "1")` is ambiguous, as the result type of `read` and argument type of `show` are completely unconstrained. This is problematic because, depending on the available canonical instances, two or more valid elaborations with observably different behaviour may exist. Hence, ambiguous programs need to be rejected. This can be done easily during inference, but for inference to be complete the completeness theorem has to be formulated in such a way that ambiguous programs are excluded from consideration. We have avoided this issue here entirely in the interest of a clearer presentation.

Another source of incompleteness is inherited from ML, and arises from the interaction between modules and type inference. Consider the following Standard ML program:

```
functor F(X : sig type t end) = struct
  val f = (print "Hello"; fn x => x)
end
structure Y1 = F(struct type t = int end)
structure Y2 = F(struct type t = bool end)
val z1 = Y1.f(3)
val z2 = Y2.f(true)
```

The binding of f in F is chosen to have an effect, so that it cannot be given a polymorphic type. This raises the question of what signature should be assigned to F . According to the Definition of Standard ML [18] (and the HS semantics as well), the above program is well-typed because f may be assigned the type $X.t \rightarrow X.t$, which is consistent with both subsequent uses of F . But in order to figure this out, a compiler would have to do a form of higher-order unification—once we leave the scope of $X.t$, the unification variable in the type of f should be skolemized over $X.t$.

As a result, nearly all existing implementations of Standard ML reject this program, as do we. (The only one that accepts it is MLton, but MLton also *accepts* similar programs that the Definition *rejects* [5].) This example points out that the interactions between type inference and modules are still not fully understood, and merit further investigation beyond the scope of this paper.

5. Related Work

Type classes in Haskell. Since Wadler and Blott’s seminal paper [23], the basic system of type classes has been extended in a number of ways. Of these, Jones’ framework of *qualified types* [11] and the resulting generalizations to constructor classes [15], multi-parameter type classes, and functional dependencies [12] are the most widely used. We discussed the option of supporting multi-parameter and constructor classes in the modular setting in Section 3.4. Instead of functional dependencies, we support associated types, as they arise naturally from type components in modules.

Achieving a separation between instance declaration and instance adoption, so that instance declarations need not have global scope, is still an open problem in the Haskell setting. There exists an experimental proposal by Kahl and Scheffczyk [16] that is motivated by a comparison with ML modules. Their basic idea is to allow constrained polymorphic functions to be given explicit instance arguments instead of having their instance arguments computed automatically. We support this functionality by providing the ability to coerce back and forth between polymorphic functions and functors, the latter of which may be given explicit module arguments (Section 2.4). Moreover, we permit different instances of the same signature to be made canonical in different scopes, which Kahl and Scheffczyk do not.

Comparing type classes and modules. The only formal comparison between ML modules and Haskell type classes is by Wehr [24]. He formalizes a translation from type classes to modules and vice versa, proves that both translations are type-preserving, and uses the translations as the basis for a comparison of the expressiveness of the language features. Wehr concludes that his encoding can help Haskell programmers to emulate certain aspects of modules in Haskell, but that the module encoding of type classes in ML is too heavyweight to be used for realistic programs. Not surprisingly, Wehr’s encoding of type classes as modules uses signatures for classes and modules for instances, as we do. In fact, his translation can be regarded as an elaboration from a Haskell core language to a fragment of ML. However, the fundamental difference between our work and his is that he performs elaboration in the *non-modular* context of Haskell, whereas we demonstrate how to perform elaboration and type inference in the *modular* context of ML.

Type classes for ML. Schneider [20] has proposed to extend ML with type classes as a feature independent of modules. This leads to significant duplication of mechanism and a number of technical problems, which we avoid by expressing type classes via modules. More recently, Siek and Lumsdaine [22] have described a language F^G that integrates *concepts*, which are closely related to type classes, into System F. However, F^G does not support type inference. Siek’s thesis [21] defines a related language \mathcal{G} , which supports inference for type applications, but not type abstractions.

Concepts in \mathcal{G} are treated as a distinct construct, unrelated to modules, and \mathcal{G} does not support parameterized modules (*i.e.*, functors).

Parameterized signatures. Jones [13] has proposed a way of supporting modular programming in a Haskell-like language, in which a signature is encoded as a record type parameterized over the abstract type components of the signature. However, he does not consider the interaction with type classes.

Acknowledgments

We thank Stefan Wehr for stimulating discussions on ways of representing type classes with modules.

References

- [1] Manuel M. T. Chakravarty, Gabriele Keller, and Simon Peyton Jones. Associated type synonyms. In *ICFP '05*.
- [2] Manuel M. T. Chakravarty, Gabriele Keller, Simon Peyton Jones, and Simon Marlow. Associated types with class. In *POPL '05*.
- [3] Luis Damas and Robin Milner. Principal type schemes for functional programs. In *POPL '82*.
- [4] Derek Dreyer. *Understanding and Evolving the ML Module System*. PhD thesis, Carnegie Mellon University, Pittsburgh, PA, May 2005.
- [5] Derek Dreyer and Matthias Blume. Principal type schemes for modular programs, October 2006. Submitted for publication.
- [6] Derek Dreyer, Karl Crary, and Robert Harper. A type system for higher-order modules. In *POPL '03*.
- [7] Derek Dreyer, Robert Harper, Manuel M.T. Chakravarty, and Gabriele Keller. Modular type classes. Technical Report TR-2006-XX, University of Chicago Computer Science Department, 2006.
- [8] Robert Harper, John C. Mitchell, and Eugenio Moggi. Higher-order modules and the phase distinction. In *POPL '90*.
- [9] Robert Harper and Chris Stone. A type-theoretic interpretation of Standard ML. In G. Plotkin, C. Stirling, and M. Tofte, editors, *Proof, Language, and Interaction: Essays in Honor of Robin Milner*. MIT Press, 2000.
- [10] Mark P. Jones. Simplifying and improving qualified types. In *FPCA '95*.
- [11] Mark P. Jones. A theory of qualified types. In *ESOP '92*.
- [12] Mark P. Jones. Type classes with functional dependencies. In *ESOP '00*.
- [13] Mark P. Jones. Using parameterized signatures to express modular structure. In *POPL '96*.
- [14] Mark P. Jones. *Qualified Types: Theory and Practice*. Cambridge University Press, 1994.
- [15] Mark P. Jones. A system of constructor classes: Overloading and implicit higher-order polymorphism. *Journal of Functional Programming*, 5(1), 1995.
- [16] Wolfram Kahl and Jan Scheffczyk. Named instances for Haskell type classes. In *Haskell Workshop*, 2001.
- [17] David MacQueen. Modules for Standard ML. In *LFP '84*.
- [18] Robin Milner, Mads Tofte, Robert Harper, and David MacQueen. *The Definition of Standard ML (Revised)*. MIT Press, 1997.
- [19] Simon Peyton Jones et al. Haskell 98 language and libraries: the revised report. *Journal of Functional Programming*, 13(1), 2003.
- [20] Gerhard Schneider. ML mit Typklassen. Master’s thesis, June 2000.
- [21] Jeremy Siek. *A Language for Generic Programming*. PhD thesis, Indiana University, August 2005.
- [22] Jeremy Siek and Andrew Lumsdaine. Essential language support for generic programming. In *PLDI '05*.
- [23] P. Wadler and S. Blott. How to make ad-hoc polymorphism less ad-hoc. In *POPL '89*.
- [24] Stefan Wehr. ML modules and Haskell type classes: A constructive comparison. Master’s thesis, Albert-Ludwigs-Universität Freiburg, Institut für Informatik, 2005.