

1ML – Core and modules united (*F-ing first-class modules*)

Andreas Rossberg

Google

rossberg@mpi-sws.org

Abstract

ML is two languages in one: there is the core, with types and expressions, and there are modules, with signatures, structures and functors. Modules form a separate, higher-order functional language on top of the core. There are both practical and technical reasons for this stratification; yet, it creates substantial duplication in syntax and semantics, and it reduces expressiveness. For example, selecting a module cannot be made a dynamic decision. Language extensions allowing modules to be packaged up as first-class values have been proposed and implemented in different variations. However, they remedy expressiveness only to some extent, are syntactically cumbersome, and do not alleviate redundancy.

We propose a redesign of ML in which modules are truly first-class values, and core and module layer are unified into one language. In this “1ML”, functions, functors, and even type constructors are one and the same construct; likewise, no distinction is made between structures, records, or tuples. Or viewed the other way round, everything is just (“a mode of use of”) modules. Yet, 1ML does not require dependent types, and its type structure is expressible in terms of plain System F_ω , in a minor variation of our *F-ing modules* approach. We introduce both an explicitly typed version of 1ML, and an extension with Damas/Milner-style implicit quantification. Type inference for this language is not complete, but, we argue, not substantially worse than for Standard ML.

An alternative view is that 1ML is a user-friendly surface syntax for System F_ω that allows combining term and type abstraction in a more compositional manner than the bare calculus.

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

General Terms Languages, Design, Theory

Keywords ML modules, first-class modules, type systems, abstract data types, existential types, System F, elaboration

1. Introduction

The ML family of languages is defined by two splendid innovations: parametric polymorphism with Damas/Milner-style type in-

ference [18, 3], and an advanced module system based on concepts from dependent type theory [17]. Although both have contributed to the success of ML, they exist in almost entirely distinct parts of the language. In particular, the convenience of type inference is available only in ML’s so-called *core language*, whereas the *module language* has more expressive types, but for the price of being painfully verbose. Modules form a separate language layered on top of the core. Effectively, ML is two languages in one.

This stratification makes sense from a historical perspective. Modules were introduced for programming-in-the-large, when the core language already existed. The dependent type machinery that was the central innovation of the original module design was alien to the core language, and could not have been integrated easily.

However, we have since discovered that dependent types are not actually necessary to explain modules. In particular, Russo [26, 28] demonstrated that module types can be readily expressed using only System-F-style quantification. The *F-ing modules* approach later showed that the entire ML module system can in fact be understood as a form of syntactic sugar over System F_ω [25].

Meanwhile, the second-class nature of modules has increasingly been perceived as a practical limitation. The standard example being that it is not possible to select modules at runtime:

```
module Table = if size > threshold then HashMap else TreeMap
```

A definition like this, where the choice of an implementation is dependent on dynamics, is entirely natural in object-oriented languages, but not expressible with ordinary ML modules.

1.1 Packaged Modules

It comes to no surprise, then, that various proposals have been made (and implemented) that enrich ML modules with the ability to package them up as first-class values [27, 22, 6, 25, 7]. Such *packaged modules* address the most imminent needs, but they are not to be confused with truly first-class modules. They require explicit injection into and projection from first-class core values, accompanied with heavy annotations. For example, in OCaml 4 the above example would have to be written as follows:

```
module Table = (val (if size > threshold
  then (module HashMap : MAP)
  else (module TreeMap : MAP))) : MAP
```

which, arguably, is neither natural nor pretty. Packaged modules have limited expressiveness as well. In particular, type sharing with a packaged module is only possible via a detour through core-level polymorphism, such as in:

```
f : (module S with type t = 'a) → (module T with type u = 'a) → 'a
```

Because core-level polymorphism is first-order, this approach cannot express type sharing between type *constructors* – a complaint that has come up several times on the OCaml mailing list; for example, if one were to abstract over a monad:

```
map : (module MONAD with type 'a t = ?) → ('a → 'b) → ? → ?
```

There is nothing that can be put in place of the ?'s to complete this function signature. The programmer is forced to either use weaker types (if possible at all), or drop the use of packaged modules and lift the function (and potentially a lot of downstream code) to the functor level – which not only is very inconvenient, it also severely restricts the possible computational behaviour of such code.

1.2 First-Class Modules

Can we overcome this situation and make modules more equal citizens of the language? The answer from the literature has been: no, because first-class modules make type-checking undecidable and type inference infeasible.

The most relevant work is Harper & Lillibridge's calculus of *translucent sums* [9] (a precursor of later work on *singleton types* [31]). It can be viewed as an idealised functional language that allows types as components of (dependent) records, so that they can express modules. In the type of such a record, individual type members can occur as either transparent or opaque (hence, *translucent*), which is the defining feature of ML module typing.

Harper & Lillibridge prove that type-checking this language is undecidable. Their result applies to any language that has (a) contravariant functions, (b) both transparent and opaque types, and (c) allows opaque types to be subtyped with arbitrary transparent types. The latter feature usually manifests in a subtyping rule like

$$\frac{\{D_1[\tau/t]\} \leq \{D_2[\tau/t]\}}{\{\mathbf{type} \ t=\tau; D_1\} \leq \{\mathbf{type} \ t; D_2\}}_{\text{FORGET}}$$

which is, in some variation, at the heart of every definition of signature matching. In the premise the concrete type τ is substituted for the abstract t . Obviously, this rule is not inductive. The substitution can arbitrarily grow the types, and thus potentially require infinite derivations. A concrete example triggering non-termination is the following, adapted from Harper & Lillibridge's paper [9]:

```
type T = {type A; f : A → ()}
type U = {type A; f : (T where type A = A) → ()}
type V = T where type A = U
g (X : V) = X : U (* V ≤ U ? *)
```

Checking $V \leq U$ would match **type** A with **type** A=U, substituting U for A accordingly, and then requires checking that the types of f are in a subtyping relation – which contravariantly requires checking that $(T \text{ where } \mathbf{type} \ A = A)[U/A] \leq A[U/A]$, but that is the same as the $V \leq U$ we wanted to check in the first place.

In fewer words, signature matching is no longer decidable when module types can be abstracted over, which is the case if module types are simply collapsed into ordinary types. It also arises if “abstract signatures” are added to the language, as in OCaml, where the same example can be constructed on the module type level.

Some may consider decidability a rather theoretical concern. However, there also is the – quite practical – issue that the introduction of signature matching into the core language makes ML-style type inference impossible. Obviously, Milner's algorithm \mathcal{W} [18] is far too weak to handle dependent types. Moreover, modules introduce subtyping, which breaks unification as the basic algorithmic tool for solving type constraints. And while inference algorithms for subtyping exist, they have much less satisfactory properties than our beloved Hindley/Milner sweet spot.

Worse, module types do not even form a lattice under subtyping:

```
f1 : {type t a; x : t int} → int
f2 : {type t a; x : int} → int
g = if condition then f1 else f2
```

There are at least two possible types for g:

```
g : {type t a = int; x : int} → int
g : {type t a = a; x : int} → int
```

Neither is more specific than the other, so no least upper bound exists. Consequently, annotations are necessary to regain principal types for constructs like conditionals, in order to restore any hope for compositional type *checking*, let alone inference.

1.3 F-ing Modules

In our work on *F-ing modules* with Russo & Dreyer [25] we have demonstrated that ML modules can be expressed and encoded entirely in vanilla System F (or F_ω , depending on the concrete core language and the desired semantics for functors). Effectively, the F-ing semantics defines a type-directed desugaring of module syntax into System F types and terms, and inversely, interprets a stylised subset of System F types as module signatures.

The core language that we assume in that paper is System $F_{(\omega)}$ itself, leading to the seemingly paradoxical situation that the core language appears to have *more* expressive types than the module language. That makes sense because the translation rules manipulate the sublanguage of module types in ways that would not generalise to arbitrary System F types. In particular, the rules *implicitly* introduce and eliminate universal and existential quantifiers, which is key to making modules a usable means of abstraction. But the process is guided by, and only meaningful for, module syntax; likewise, the built-in subtyping relation is only “complete” for the specific occurrences of quantifiers in module types.

Nevertheless, the observation that modules are just sugar for certain kinds of constructs that the core language can already express (even if less concisely), raises the question: what necessitates modules to be second-class in that system?

1.4 IML

The answer to that question is: very little! And the present paper is motivated by exploring that answer.

In essence, the F-ing modules semantics reveals that the syntactic stratification between ML core and module language is merely a rather coarse means to enforce *predicativity* for module types: it prevents that abstract types themselves can be instantiated with binders for abstract types. But this heavy *syntactic* restriction can be replaced by a more surgical *semantic* restriction! It is enough to employ a simple universe distinction between *small* and *large* types (reminiscent of Harper & Mitchell's XML [10]), and limit the equivalent of the FORGET rule shown earlier to only allow small types for substitution, which serves to exclude problematic quantifiers.

That would settle decidability, but what about type inference? Well, we can use the same distinction! A quick inspection of the subtyping rules in the F-ing modules semantics reveals that they, almost, degenerate to type equivalence when applied to *small* types — the main exception being width subtyping on structures. If we are willing to accept that inference is not going to be complete for records (which it already isn't in Standard ML), then a simple restriction to inferring only small types is sufficient to make type inference work almost as usual.

In this spirit, this paper presents *IML*, an ML-dialect in which modules are truly first-class values. The name is both short for “1st-class module language” and a pun on the fact that it unifies core and modules of ML into one language.

We see several benefits with this redesign: it produces a language that is more *expressive* and *concise*, and at the same time, more *minimal* and *uniform*. “Modules” become a natural way to express all forms of (first-class) polymorphism, and can be freely intermixed with “computational” code and data. Type inference integrates in a rather seamless manner, reducing the need for explicit annotations to large types, module or not. Every programming concept is derived from a small set of orthogonal constructs, over which general and uniform syntactic sugar can be defined.

2. IML with Explicit Types

To separate concerns a little, we will start out by introducing IML_{ex} , a sublanguage of IML proper that is explicitly typed and does not support any type inference. Its kernel syntax is given in Figure 1. Let us take a little tour of IML_{ex} by way of examples.

Functional Core A major part of IML_{ex} consists of fairly conventional functional language constructs. On the expression level, as a representative for a base type, we have Booleans (in examples that follow, we will often assume the presence of an integer type and respective constructs as well). Then there are records, which consist of a sequence of bindings. And of course, it wouldn't be a functional language without functions.

In a first approximation, these forms are reflected on the type level as one would expect, except that for functions we allow two forms of arrows, distinguishing pure function types (\Rightarrow) from impure ones (\rightarrow) (discussed later).

Like in the F-ing modules paper [25], most elimination forms in the kernel syntax only allow variables as subexpressions. However, the general expression forms are all definable as straightforward syntactic sugar, as shown in the lower half of Figure 1. For example,

$(\text{fun } (n : \text{int}) \Rightarrow n + n) \ 3$

desugars into

$\text{let } f = \text{fun } (n : \text{int}) \Rightarrow n + n; \ x = 3 \ \text{in } f \ x$

and further into

$\{f = \text{fun } (n : \text{int}) \Rightarrow n + n; \ x = 3; \ \text{body} = f \ x\}.\text{body}$

This works because records actually behave like ML structures, such that every bound identifier is in scope for later bindings – which enables encoding let-expressions.

Also, notably, if-expressions require a type annotation in IML_{ex} . As we will see, the type language subsumes module types, and as discussed in Section 1.2 there wouldn't generally be a unique least upper bound otherwise. However, in Section 4 we show that this annotation can usually be omitted in full IML.

Reified Types The core feature that makes IML_{ex} able to express modules is the ability to embed types in a first-class manner: the expression $\text{type } T$ reifies the type T as a value.¹ Such an expression has type type , and thereby can be abstracted over. For example,

$\text{id} = \text{fun } (a : \text{type}) \Rightarrow \text{fun } (x : a) \Rightarrow x$

defines a polymorphic identity function, similar to how it would be written in dependent type theories. Note in particular that a is a *term* variable, but it is used as a *type* in the annotation for x . This is enabled by the “path” form E in the syntax of types, which expresses the (implicit) projection of a type from a term, provided this term has type type . Consequently, all variables are term variables in IML, there is no separate notion of type variable.

More interestingly, a function can *return* types, too. Consider

$\text{pair} = \text{fun } (a : \text{type}) \Rightarrow \text{fun } (b : \text{type}) \Rightarrow \text{type } \{\text{fst} : a; \ \text{snd} : b\}$

which takes a type and returns a type, and effectively defines a type *constructor*. Applied to a reified type it yields a reified type. Again, the implicit projection from “paths” enables using this as a type:

$\text{second} = \text{fun } (a : \text{type}) \Rightarrow \text{fun } (b : \text{type}) \Rightarrow \text{fun } (p : \text{pair } a \ b) \Rightarrow p.\text{snd}$

In this example, the whole of “pair a b” is a term of type type .

¹Ideally, “ $\text{type } T$ ” should be written just “ T ”, like in dependently typed systems. However, that would create various syntactic ambiguities, e.g. for phrases like “ $\{\}$ ”, which could only be avoided by moving to a more artificial syntax for types themselves. Nevertheless, we at least allow writing “ $E \ T$ ” for the application “ $E (\text{type } T)$ ” if T unambiguously is a type.

Figure 1 also defines a bit of syntactic sugar to make function and type definitions look more like in traditional ML. For example, the previous functions could equivalently be written as

$\text{id } a \ (x : a) = x$
 $\text{type pair } a \ b = \{\text{fst} : a; \ \text{snd} : b\}$
 $\text{second } a \ b \ (p : \text{pair } a \ b) = p.\text{snd}$

It may seem surprising that we can just reify types as first-class values. But reified types (or “atomic type modules”) have been common in module calculi for a long time [16, 6, 24, 25]. We are merely making them available in the source language directly. For the most part, this is just a notational simplification over what first-class modules already offer: instead of having to define $T = \{\text{type } t = \text{int}\} : \{\text{type } t\}$ and then refer to $T.t$, we allow injecting types into modules (i.e., values) *anonymously*, without wrapping them into a structure; thus $T = (\text{type } \text{int}) : \text{type}$, which can be referred to as just T .

Translucency The type type allows classifying types abstractly: given a value of type type , nothing is known about *what* type it is. But for modular programming it is essential that types can selectively be specified *transparently*, which enables expressing the vital concept of *type sharing* [12].

As a simple example, consider these type aliases:

$\text{type size} = \text{int}$
 $\text{type pair } a \ b = \{\text{fst} : a; \ \text{snd} : b\}$

According to the idea of translucency, the variables defined by these definitions can be classified in one of two ways. Either opaquely:

size : type
 $\text{pair} : (a : \text{type}) \Rightarrow (b : \text{type}) \Rightarrow \text{type}$

Or transparently:

size : $(= \text{type } \text{int})$
 $\text{pair} : (a : \text{type}) \Rightarrow (b : \text{type}) \Rightarrow (= \text{type } \{\text{fst} : a; \ \text{snd} : b\})$

The latter use a variant of *singleton types* [31, 6] to reveal the definitions: a type of the form “ $=E$ ” is inhabited only by values that are “structurally equivalent” to E , in particular, with respect to parts of type type . It allows the type system to infer, for example, that the application pair size size is equivalent to the (reified) type $\{\text{fst} : \text{int}; \ \text{snd} : \text{int}\}$. A type $=E$ is a subtype of the type of E itself, and consequently, transparent classifications define subtypes of opaque ones, which is the crux of ML signature matching.

Translucent types usually occur as part of module type declarations, where IML can abbreviate the above to the more familiar

type size or, respectively, $\text{type size} = \text{int}$
 $\text{type pair } a \ b$ or, respectively, $\text{type pair } a \ b = \{\text{fst} : a; \ \text{snd} : b\}$

(i.e., as in ML, transparent declarations look just like definitions).

Singletons can be formed over arbitrary values. This gives the ability to express *module sharing* and *aliases*. In the basic semantics described in this paper, this is effectively a shorthand for sharing all types contained in the module (including those defined inside transparent functors, see below). We leave the extension to full *value* equivalence (including primitive types like Booleans), as in our F-ing semantics for applicative functors [25], to future work.

Functors Returning to the IML grammar, the remaining constructs of the language are typical for ML modules, although they are perhaps a bit more general than what is usually seen. Let us explain them using an example that demonstrates that our language can readily express “real” modules as well. Here is the (unavoidable, it seems) functor that defines a simple map ADT:

$\text{type EQ} =$
 $\{$
 $\quad \text{type } t;$
 $\quad \text{eq} : t \rightarrow t \rightarrow \text{bool}$
 $\}$

(identifiers)	X		
(types)	$T ::= E \mid \mathbf{bool} \mid \{D\} \mid (X:T) \Rightarrow T \mid \mathbf{type} \mid =E \mid T \mathbf{where} (\overline{X}:T)$		
(declarations)	$D ::= X : T \mid \mathbf{include} T \mid D;D \mid \epsilon$		
(expressions)	$E ::= X \mid \mathbf{true} \mid \mathbf{false} \mid \mathbf{if} X \mathbf{then} E \mathbf{else} E:T \mid \{B\} \mid E.X \mid \mathbf{fun} (X:T) \Rightarrow E \mid X X \mid \mathbf{type} T \mid X:>T$		
(bindings)	$B ::= X=E \mid \mathbf{include} E \mid B;B \mid \epsilon$		
(types)		(expressions)	
$\mathbf{let} B \mathbf{in} T$	$:= \{B;X=\mathbf{type} T\}.X$	$\mathbf{let} B \mathbf{in} E$	$:= \{B;X=E\}.X$
$T_1 \Rightarrow T_2$	$:= (X:T_1) \Rightarrow T_2$	$\mathbf{if} E_1 \mathbf{then} E_2 \mathbf{else} E_3:T$	$:= \mathbf{let} X=E_1 \mathbf{in} \mathbf{if} X \mathbf{then} E_2 \mathbf{else} E_3:T$
$T \mathbf{where} (\overline{X} \overline{P}=E)$	$:= T \mathbf{where} (\overline{X}:\overline{P} \Rightarrow (=E))$	$E_1 E_2$	$:= \mathbf{let} X_1=E_1; X_2=E_2 \mathbf{in} X_1 X_2$
$T \mathbf{where} (\mathbf{type} \overline{X} \overline{P}=T')$	$:= T \mathbf{where} (\overline{X}:\overline{P} \Rightarrow (= \mathbf{type} T'))$	$E T$	$:= E(\mathbf{type} T)$ (if T unambiguous)
(declarations)		$E:T$	$:= ((X:T) \Rightarrow X) E$
$\mathbf{local} B \mathbf{in} D$	$:= \mathbf{include} (\mathbf{let} B \mathbf{in} \{D\})$	$E:>T$	$:= \mathbf{let} X=E \mathbf{in} X:>T$
$X \overline{P}:T$	$:= X:\overline{P} \Rightarrow T$	$\mathbf{fun} \overline{P} \Rightarrow E$	$:= \overline{\mathbf{fun}} \overline{P} \Rightarrow E$
$X \overline{P}=E$	$:= X:\overline{P} \Rightarrow (=E)$	(bindings)	
$\mathbf{type} X \overline{P}$	$:= X:\overline{P} \Rightarrow \mathbf{type}$	$\mathbf{local} B \mathbf{in} B'$	$:= \mathbf{include} (\mathbf{let} B \mathbf{in} \{B'\})$
$\mathbf{type} X \overline{P}=T$	$:= X:\overline{P} \Rightarrow (= \mathbf{type} T)$	$X \overline{P}:T' :>T''=E$	$:= X = \mathbf{fun} \overline{P} \Rightarrow E:T' :>T''$
where: (parameter) $P ::= (X:T)$	$X := (X:\mathbf{type})$	$\mathbf{type} X \overline{P}=T$	$:= X = \mathbf{fun} \overline{P} \Rightarrow \mathbf{type} T$

(Identifiers X only occurring on the right-hand side are considered fresh)

Figure 1. IML_{ex} syntax and syntactic abbreviations

```

};
type MAP =
{
  type key;
  type map a;
  empty a : map a;
  add a : key → a → map a → map a;
  lookup a : key → map a → opt a
};
Map (Key : EQ) :> MAP where (type .key = Key.t) =
{
  type key = Key.t;
  type map a = key → opt a;
  empty a = fun (k : key) ⇒ none a;
  lookup a (k : key) (m : map a) = m k;
  add a (k : key) (v : a) (m : map a) =
    fun (x : key) ⇒ if Key.eq x k then some a v else m x : opt a
}

```

The record type EQ amounts to a module signature, since it contains an abstract type component t . It is referred to in the type of eq , which shows that record types are *dependent*: like for terms, earlier components are in scope for later components.

Similarly, MAP defines a signature with abstract key and map types. Note how type parameters on the left-hand side conveniently and uniformly generalise to value declarations, avoiding the need for brittle implicit scoping rules like in conventional ML: as shown in Figure 1, “empty $a : T$ ” means “empty : ($a : \mathbf{type}$) $\Rightarrow T$ ”.

The Map function is a functor: it takes a value of type EQ, i.e., a module. From that it constructs a naive implementation of maps. “ $X:>T$ ” is the usual *sealing* operator that opaquely ascribes a type (i.e., signature) to a value (a.k.a. module). The *type refinement* syntax “ $T \mathbf{where} (\mathbf{type} \overline{X}=T)$ ” should be familiar from ML, but here it actually is derived from a more general construct: “ $T \mathbf{where} (\overline{X}:U)$ ” refines T ’s subcomponent at path \overline{X} to type U , which can be any subtype of what’s declared by T . That form subsumes module sharing as well as other forms of refinement.

Applicative vs. Generative In this paper, we stick to a relatively simple semantics for functor-like functions, in which Map is *generative* [28, 4, 25]. That is, like in Standard ML, each application will yield a fresh map ADT, because sealing occurs inside the functor:

```

M1 = Map IntEq;
M2 = Map IntEq;
m = M1.add int 7 M2.empty (* ill-typed: M1.map ≠ M2.map *)

```

But as we saw earlier, type constructors like `pair` or `map` are essentially functors, too! Sealing the body of the Map functor hence implies higher-order sealing of the nested map “functor”, as if performing `map :> type ⇒ type`. It is vital that the resulting functor has *applicative* semantics [15, 25], so that

```

type map a = M1.map a; type t = map int; type u = map int

```

yields $t = u$, as one would expect from a proper type constructor.

We hence need applicative functors as well. To keep things simple, we restrict ourselves to the simplest possible semantics in this paper, in which we distinguish between pure (\Rightarrow , i.e. applicative) and impure (\rightarrow , i.e. generative) function types, but sealing is always impure (or *strong* [6]). That is, sealing *inside* a functor always makes it generative. The only way to produce an applicative functor is by sealing a (fully transparent) functor *as a whole*, with applicative functor type, as for the map type constructor above.

For example, consider:

```

F = (fun (a : type) ⇒ type {x : a}) :> type ⇒ type
G = (fun (a : type) ⇒ type {x : a}) :> type → type
H = fun (a : type) ⇒ (type {x : a} :> type)
J = G :> type ⇒ type (* ill-typed! *)

```

F is an applicative functor, such that $F \text{ int} = F \text{ int}$. G and H on the other hand are generative functors; the former because it is sealed with impure functor type, the latter because sealing occurs inside its body. Consequently, G int or H int are impure expressions and invalid as type paths (though it is fine to bind their result to a name, e.g., “**type** w = G int”, and use the constant w as a type). Lastly, J is ill-typed, because applicative functor types are subtypes of generative ones, but not the other way round.

This semantics for applicative functors (which is very similar to the applicative functors of Shao [30]) is somewhat limited, but just enough to encode sealing over type constructors and hence recover the ability to express type definitions as in conventional ML. An extension of IML to applicative functors with *pure* sealing à la F-ing modules [25] is given in the Technical Appendix [23].

Higher Polymorphism So far, we have only shown how IML recovers constructs well-known from ML. As a first example of something that cannot directly be expressed in conventional ML, consider first-class polymorphic arguments:

```
f (id : (a : type) ⇒ a → a) = {x = id int 5; y = id bool true}
```

Similarly, existential types are directly expressible:

```
type SHAPE = {type t; area : t → float; v : t}
volume (height : int) (x : SHAPE) = height * x.area (x.v)
```

SHAPE can either be read as a module signature or an existential type, both are indistinguishable. The function volume is agnostic about the actual type of the shape it's passed.

It turns out that the previous examples can still be expressed with packaged modules (Section 1.1). But now consider:

```
type COLL c =
{
  type key;
  type val;
  empty : c;
  add : c → key → val → c;
  lookup : c → key → opt val;
  keys : c → list key
};
entries c (C : COLL c) (xs : c) : list (C.key × C.val) = ...
```

COLL amounts to a *parameterised signature*, and is akin to a Haskell-style type class [34]. It contains two abstract type specifications, which are known as *associated types* in the type class literature (or in C++ land). The function entries is parameterised over a corresponding module C – an (explicit) type class instance if you want. Its result type depends directly on C's definition of the associated types. Such a dependency can be expressed in ML on the module level, but not on the core level.²

Moving to higher kinds, things become even more interesting:s

```
type MONAD (m : type ⇒ type) =
{
  return a : a → m a;
  bind a b : m a → (a → m b) → m b
};
map a b (m : type ⇒ type) (M : MONAD m) (f : a → b) (mx : m a) =
M.bind a b mx (fun (x : a) ⇒ M.return b (f x)) (* : m b *)
```

Here, MONAD is again akin to a type class, but over a type constructor. As explained in Section 1.1, this kind of polymorphism cannot be expressed even in MLs with packaged modules.

Computed Modules Just for completeness, we should mention that the motivating example from Section 1 can of course be written (almost) as is in IML_{ex}:

```
Table = if size > threshold then HashMap else TreeMap : MAP
```

The only minor nuisance is the need to annotate the type of the conditional, as explained earlier.

Predicativity What is the restriction we employ to maintain decidability? It is simple: during subtyping (a.k.a. signature matching) the type **type** can only be matched by *small* types, which are those that do not themselves contain the type **type** opaquely; or in other words, monomorphic types. This restriction affects annotations, parameterisation over types, and the formation of abstract

² In OCaml 4, this example can be approximated with heavy fibration:

```
module type COLL = sig type coll type key type val ... end
let entries (type c) (type k) (type v)
  (module C : COLL with
   type coll = c and type key = k and type value = v)
  (xs : c) : (k * v) list = ...
```

types. For example, for any of the following *large* types T_i ,

```
type T1 = type;           type T4 = (x : {}) → type;
type T2 = {type u};       type T5 = (a : type) ⇒ {};
type T3 = {type u = T2}; type T6 = {type u a = bool};
```

the following definitions are all ill-typed:

```
type U = pair Ti Ti; (* error *)
A = (type Ti) : type; (* error *)
B = {type u = Ti} :> {type u}; (* error *)
C = if b then Ti else int : type (* error *)
```

Notably, the case A with T_1 literally implies **type type / type** (although **type type** itself *is* a well-formed expression!). The main challenge with first-class modules is preventing such a type:type situation, and the separation into a small universe (denoted by **type**) and a large one (for which no syntax exists) achieves that.

A *transparent* type is small as long as it reveals a small type:

```
type T'1 (= type int);   type T'2 = {type u = int};
```

would *not* cause an error when inserted into the above definitions.

Recursion The IML_{ex} syntax we give in Figure 1 omits a couple of constructs that one can rightfully expect from any serious ML contender: in particular, there is no form of recursion, neither for terms nor for types. It turns out that those are largely orthogonal to the overall design of IML, so we only sketch them here.

ML-style recursive functions can be added simply by throwing in a primitive polymorphic fixpoint operator

$$\text{fix } a \ b : (a \rightarrow b) \rightarrow (a \rightarrow b)$$

plus perhaps some suitable syntactic sugar:

```
rec X  $\bar{Y}$  (Z:T) : U=E :=
X = fun  $\bar{Y}$  ⇒ fix T U (fun (X:(Z:T) → T') ⇒ fun (Z:T) ⇒ E)
```

Given an appropriate fixpoint operator, this generalises to mutually recursive functions in the usual ways. Note how the need to specify the result type b (respectively, U) prevents using the operator to construct transparent recursive types, because U has no way of referring to the result of the fixpoint. Moreover, fix yields an impure function, so even an attempt to define an abstract type recursively,

```
rec stream (a : type) : type = type {head : a; tail : stream a}
```

won't type-check, because stream wouldn't be an applicative functor, and so the term $\text{stream } a$ on the right-hand side is not a valid type — fortunately, because there would be no way to translate such a definition into System F_ω with a conventional fixpoint operator.

Recursive (data)types have to be added separately. One approach, that has been used by Harper & Stone's type-theoretic account of Standard ML [13], is to interpret a recursive datatype like

```
datatype t = A | B of T
```

as a module defining a primitive ADT with the signature

```
{type t; A : t; B : T ⇒ t; expose a : ({ } → a) ⇒ (T → a) ⇒ t → a}
```

where expose is a case-operator accessed by pattern matching compilation. We refer to [13] for more details on this approach.

Impredicativity Reloaded Predicativity is a severe restriction. Can we enable impredicative type abstraction without breaking decidability? Yes we can. One possibility is the usual trick of piggy-backing datatypes: we can allow their data constructors to have large parameters. Because datatypes are *nominal* in ML, impredicativity is “hidden away” and does not interfere with subtyping.

Structural impredicative types are also possible, as long as large types are injected into the small universe *explicitly*, by way of a special type “**wrap T**”. The gist of this approach is that subtyping does not extend to wrapped types. It is an easy extension, the Technical Appendix [23] gives the details.

(kinds)	$\kappa ::= \Omega \mid \kappa \rightarrow \kappa$
(types)	$\tau ::= \alpha \mid \tau \rightarrow \tau \mid \{\overline{l:\tau}\} \mid \forall \alpha:\kappa.\tau \mid \exists \alpha:\kappa.\tau \mid \lambda \alpha:\kappa.\tau \mid \tau \tau$
(terms)	$e, f ::= x \mid \lambda x:\tau.e \mid e e \mid \{\overline{l=e}\} \mid e.l \mid \lambda \alpha:\kappa.e \mid e \tau \mid \text{pack } \langle \tau, e \rangle_\tau \mid \text{unpack } \langle \alpha, x \rangle = e \text{ in } e$
(environ's)	$\Gamma ::= \cdot \mid \Gamma, \alpha:\kappa \mid \Gamma, x:\tau$

Figure 2. Syntax of F_ω

(abstracted)	$\Xi ::= \exists \bar{\alpha}.\Sigma$
(large)	$\Sigma ::= \pi \mid \text{bool} \mid [= \Xi] \mid \{\overline{l:\Sigma}\} \mid \forall \bar{\alpha}.\Sigma \rightarrow_\iota \Xi$
(small)	$\sigma ::= \pi \mid \text{bool} \mid [= \sigma] \mid \{\overline{l:\sigma}\} \mid \sigma \rightarrow_I \sigma$
(paths)	$\pi ::= \alpha \mid \pi \bar{\sigma}$
(purity)	$\iota ::= P \mid I$

Desugarings into F_ω :

(types)	$[= \tau] ::= \{\text{typ} : \tau \rightarrow \{\}\}$	(terms)	$[\tau] ::= \{\text{typ} = \lambda x:\tau.\{\}\}$
$\tau_1 \rightarrow_\iota \tau_2$	$:= \tau_1 \rightarrow \{l : \tau_2\}$	$\lambda l x:\tau.e$	$:= \lambda x:\tau.\{l : e\}$

Notation:	$\iota \leq \iota$	$\iota \vee \iota := \iota$	$\iota(\Sigma) = P$
	$P \leq I$	$P \vee I := I \vee P := I$	$\iota(\exists \bar{\alpha}.\Sigma) = I$
	$\tau.\bar{l} := \tau$	$\tau[\bar{l}=\tau_2] := \tau_2$	$(\bar{l} = \epsilon)$
	$\{l:\tau, \dots\}.\bar{l} := \tau.\bar{l}'$	$\{l:\tau, \dots\}[\bar{l}=\tau_2] := \{l:\tau[\bar{l}=\tau_2], \dots\}$	$(\bar{l} = l.\bar{l}')$

Figure 3. Semantic Types

3. Type System and Elaboration

So much for leisure, now for work. The general recipe for IML_{ex} is simple: take the semantics from F-ing modules [25], collapse the levels of modules and core, and impose the predicativity restriction needed to maintain decidability. This requires surprisingly few changes to the whole system. Unfortunately, space does not permit explaining all of the F-ing semantics in detail, so we encourage the reader to refer to [25] (mostly Section 4) for background, and will focus primarily on the differences and novelties in what follows.

3.1 Internal Language

System F_ω The semantics is defined by elaborating IML_{ex} types and terms into types and terms of (call-by-value, impredicative) System F_ω , the higher-order polymorphic λ -calculus [1], extended with simple record types (Figure 2). We assume obvious encodings of let-expressions and n -ary universal and existential types. The semantics is completely standard; we omit it here and reuse the formulation from [25]. The only point of note is that it allows term (but not type) variables in the environment Γ to be shadowed without α -renaming, which is convenient for translating bindings.

To ease notation we often drop type annotations from let, pack, and unpack where clear from context. We will also omit kind annotations on type variables, and where necessary, use the notation κ_α to refer to the kind implicitly associated with α .

Semantic Types Elaboration translates IML_{ex} types directly into “equivalent” System F_ω types. The shape of these *semantic* types is given by the grammar in Figure 3.

The main magic of the elaboration is that it inserts appropriate quantifiers to bind abstract types. Following Mitchell & Plotkin [20], abstract types are represented by existentials: an *abstracted* type $\Xi = \exists \bar{\alpha}.\Sigma$ quantifies over all the abstract types (i.e., components of type **type**) from the underlying *concretised* type Σ , by naming them $\bar{\alpha}$. Inside Σ they can hence be represented as transparent types, equal to those $\bar{\alpha}$'s. A sketch of the mapping between

syntactic types T and semantic types Ξ is as follows:

type	$\rightsquigarrow \exists \alpha.[= \alpha]$
(= type T)	$\rightsquigarrow [= \Xi]$
$\{X_1:T_1; X_2:T_2\}$	$\rightsquigarrow \exists \bar{\alpha}_1.\bar{\alpha}_2.\{X_1:\Sigma_1, X_2:\Sigma_2\}$
$(X:T_1) \rightarrow T_2$	$\rightsquigarrow \forall \bar{\alpha}_1.\Sigma_1 \rightarrow_I \exists \bar{\alpha}_2.\Sigma_2$
$(X:T_1) \Rightarrow T_2$	$\rightsquigarrow \exists \bar{\alpha}_2.\forall \bar{\alpha}_1.\Sigma_1 \rightarrow_P \Sigma_2$
A.t	$\rightsquigarrow \alpha_{A.t}$
$F(A).u$	$\rightsquigarrow \alpha_{F(-).u} \bar{\sigma}_{A.t}$

That is, (transparent) reified types are represented as $[= \Xi]$, using a similar coding trick as in [25]. With all abstract types being named, they always appear as transparent as well, albeit quantified. Because all type constructors are represented as functors, we have no need for reified types of higher kind.

Records, no surprise, map to records. We assume an implicit injection from IML identifiers X into both F_ω variables x and labels l , so we can conveniently treat any X as a variable or label.

Function types map to polymorphic functions in F_ω . Being in negative position, the existential quantifier for the abstract types $\bar{\alpha}_1$ from the parameter type Σ_1 turns into a universal quantifier, scoping over the whole type, and allowing the result type Σ_2 to refer to the parameter types. Functions are also annotated by a simple *effect* ι , which distinguishes pure from impure functions, and thus, applicative from generative functors. Pure function types encode applicative semantics for the abstract types they return by having their existential quantifiers $\bar{\alpha}_2$ “lifted” over their parameters. To capture potential dependencies, the $\bar{\alpha}_2$ are skolemised over $\bar{\alpha}_1$ [2, 28, 25]. That is, the kinds of $\bar{\alpha}_2$ are of the form $\kappa_{\bar{\alpha}_1} \rightarrow \kappa$, which is where higher kinds come into play. We impose the syntactic invariant that a pure function type never has an existential quantifier on the right.

Abstract types are denoted by their type variables, but may generally take the form of a *semantic path* π if they have parameters. Projecting an abstract type from an application of a pure function (applicative functor) becomes the application of a higher-kinded type variable to the concrete types from its argument. Because we enforce predicativity, these argument types have to be small.

Figure 3 also defines the subgrammar of small types, which cannot have quantifiers in them. Moreover, small functions are required to be impure, which will simplify type inference (Section 5).

3.2 Elaboration

The complete elaboration rules for IML_{ex} are collected in Figure 4. There is one judgement for each syntactic class, plus an auxiliary judgement for subtyping. If you are merely interested in typing IML then you can ignore the greyed out parts “ $\rightsquigarrow e$ ” in the rules – they are concerned with the translation of terms, and are only relevant to define the operational semantics of the language.

Types and Declarations The main job of the elaboration rules for types is to name all abstract type components with type variables, collect them, and bind them hoisted to an outermost existential (or universal, in the case of functions) quantifier. The rules are mostly identical to [25], except that **type** is a free-standing construct instead of being tied to the syntax of bindings, and IML’s “**where**” construct requires a slightly more general rule.

Also, we drop the side condition for Σ to be *explicit* in rule TSING (corresponding to rule S-LIKE in [25]), as explained below.

Expressions and Bindings The elaboration of expressions closely follows the rules from the first part of [25], but adds the tracking of purity as in Section 7 of that paper. However, to keep the current paper simple, we left out the ability to perform pure sealing, or to create pure functions around it. That avoids some of the notational contortions necessary for the applicative functor semantics from [25]. An extension of IML_{ex} with pure sealing can be found in the Technical Appendix [23].

Types

$$\begin{array}{c}
\frac{\Gamma \vdash E :_{\text{p}} [= \Xi] \rightsquigarrow e}{\Gamma \vdash E \rightsquigarrow \Xi} \text{TPATH} \quad \frac{\kappa_{\alpha} = \Omega}{\Gamma \vdash \mathbf{type} \rightsquigarrow \exists \alpha. [= \alpha]} \text{TTYPE} \quad \frac{}{\Gamma \vdash \mathbf{bool} \rightsquigarrow \mathbf{bool}} \text{TBOOL} \quad \frac{\Gamma \vdash D \rightsquigarrow \Xi}{\Gamma \vdash \{D\} \rightsquigarrow \Xi} \text{TSTR} \\
\frac{\Gamma \vdash T_1 \rightsquigarrow \exists \bar{\alpha}_1. \Sigma_1 \quad \Gamma, \bar{\alpha}_1, X : \Sigma_1 \vdash T_2 \rightsquigarrow \exists \bar{\alpha}_2. \Sigma_2}{\Gamma \vdash (X : T_1) \rightarrow T_2 \rightsquigarrow \forall \bar{\alpha}_1. \Sigma_1 \rightarrow_1 \exists \bar{\alpha}_2. \Sigma_2} \text{TFUN} \quad \frac{\Gamma \vdash T_1 \rightsquigarrow \exists \bar{\alpha}_1. \Sigma_1 \quad \Gamma, \bar{\alpha}_1, X : \Sigma_1 \vdash T_2 \rightsquigarrow \exists \bar{\alpha}_2. \Sigma_2 \quad \kappa_{\alpha'_2} = \bar{\kappa}_{\alpha_1} \rightarrow \kappa_{\alpha_2}}{\Gamma \vdash (X : T_1) \Rightarrow T_2 \rightsquigarrow \exists \bar{\alpha}'_2. \forall \bar{\alpha}_1. \Sigma_1 \rightarrow_{\text{p}} \Sigma_2 [\bar{\alpha}'_2 \bar{\alpha}_1 / \bar{\alpha}_2]} \text{TPFUN} \\
\frac{\Gamma \vdash E :_{\text{p}} \Sigma \rightsquigarrow e}{\Gamma \vdash (= E) \rightsquigarrow \Sigma} \text{TSING} \quad \frac{\Gamma \vdash T_1 \rightsquigarrow \exists \bar{\alpha}_1. \Sigma_1 \quad \bar{\alpha}_1 = \bar{\alpha}_{11} \uplus \bar{\alpha}_{12} \quad \Gamma \vdash T_2 \rightsquigarrow \exists \bar{\alpha}_2. \Sigma_2 \quad \Gamma, \bar{\alpha}_{11}, \bar{\alpha}_2 \vdash \Sigma_2 \leq_{\bar{\alpha}_{12}} \Sigma_1. \bar{X} \rightsquigarrow \delta; f}{\Gamma \vdash T_1 \mathbf{where} (\bar{X} : T_2) \rightsquigarrow \exists \bar{\alpha}_{11} \bar{\alpha}_2. \delta \Sigma_1 [\bar{X} = \Sigma_2]} \text{TWHERE}
\end{array}$$

Declarations

$$\begin{array}{c}
\frac{\Gamma \vdash T \rightsquigarrow \exists \bar{\alpha}. \Sigma}{\Gamma \vdash X : T \rightsquigarrow \exists \bar{\alpha}. \{X : \Sigma\}} \text{DVAR} \quad \frac{\Gamma \vdash T \rightsquigarrow \exists \bar{\alpha}. \{\bar{X} : \bar{\Sigma}\}}{\Gamma \vdash \mathbf{include} T \rightsquigarrow \exists \bar{\alpha}. \{X : \Sigma\}} \text{DINCL} \\
\frac{\Gamma \vdash D_1 \rightsquigarrow \exists \bar{\alpha}_1. \{\bar{X}_1 : \bar{\Sigma}_1\} \quad \Gamma, \bar{\alpha}_1, \bar{X}_1 : \bar{\Sigma}_1 \vdash D_2 \rightsquigarrow \exists \bar{\alpha}_2. \{\bar{X}_2 : \bar{\Sigma}_2\} \quad \bar{X}_1 \cap \bar{X}_2 = \emptyset}{\Gamma \vdash D_1; D_2 \rightsquigarrow \exists \bar{\alpha}_1 \bar{\alpha}_2. \{\bar{X}_1 : \bar{\Sigma}_1, \bar{X}_2 : \bar{\Sigma}_2\}} \text{DSEQ} \quad \frac{}{\Gamma \vdash \epsilon \rightsquigarrow \{\}} \text{DEMPTY}
\end{array}$$

Expressions

$$\begin{array}{c}
\frac{\Gamma(X) = \Sigma}{\Gamma \vdash X :_{\text{p}} \Sigma \rightsquigarrow X} \text{EVAR} \quad \frac{\Gamma \vdash T \rightsquigarrow \Xi}{\Gamma \vdash \mathbf{type} T :_{\text{p}} [= \Xi] \rightsquigarrow [\Xi]} \text{ETYPE} \quad \frac{}{\Gamma \vdash \mathbf{true} :_{\text{p}} \mathbf{bool} \rightsquigarrow \mathbf{true}} \text{ETRUE} \\
\frac{}{\Gamma \vdash \mathbf{false} :_{\text{p}} \mathbf{bool} \rightsquigarrow \mathbf{false}} \text{EFALSE} \quad \frac{\Gamma \vdash X :_{\text{p}} \mathbf{bool} \rightsquigarrow e \quad \Gamma \vdash E_1 :_{i_1} \Xi_1 \rightsquigarrow e_1 \quad \Gamma \vdash \Xi_1 \leq \Xi \rightsquigarrow f_1 \quad \Gamma \vdash T \rightsquigarrow \Xi \quad \Gamma \vdash E_2 :_{i_2} \Xi_2 \rightsquigarrow e_2 \quad \Gamma \vdash \Xi_2 \leq \Xi \rightsquigarrow f_2}{\Gamma \vdash \mathbf{if} X \mathbf{then} E_1 \mathbf{else} E_2 :_{i_1 \vee i_2 \vee i(\Xi)} \Xi \rightsquigarrow \mathbf{if} e \mathbf{then} f_1 e_1 \mathbf{else} f_2 e_2} \text{EIF} \\
\frac{\Gamma \vdash B :_{i} \Xi \rightsquigarrow e}{\Gamma \vdash \{B\} :_{i} \Xi \rightsquigarrow e} \text{ESTR} \quad \frac{\Gamma \vdash E :_{i} \exists \bar{\alpha}. \{\bar{X}' : \bar{\Sigma}'\} \rightsquigarrow e \quad X : \Sigma \in \bar{X}' : \bar{\Sigma}'}{\Gamma \vdash E.X :_{i} \exists \bar{\alpha}. \Sigma \rightsquigarrow \mathbf{unpack} (\bar{\alpha}, y) = e \mathbf{in} \mathbf{pack} (\bar{\alpha}, y.X)} \text{EDOT} \\
\frac{\Gamma \vdash T \rightsquigarrow \exists \bar{\alpha}. \Sigma \quad \Gamma, \bar{\alpha}, X : \Sigma \vdash E :_{i} \Xi \rightsquigarrow e}{\Gamma \vdash \mathbf{fun} (X : T) \Rightarrow E :_{\text{p}} \forall \bar{\alpha}. \Sigma \rightarrow_i \Xi \rightsquigarrow \lambda \bar{\alpha}. \lambda_i X : \Sigma. e} \text{EFUN} \quad \frac{\Gamma \vdash X_1 :_{\text{p}} \forall \bar{\alpha}. \Sigma_1 \rightarrow_i \Xi \rightsquigarrow e_1 \quad \Gamma \vdash X_2 :_{\text{p}} \Sigma_2 \rightsquigarrow e_2 \quad \Gamma \vdash \Sigma_2 \leq_{\bar{\alpha}} \Sigma_1 \rightsquigarrow \delta; f}{\Gamma \vdash X_1 X_2 :_{i} \delta \Xi \rightsquigarrow (e_1 (\delta \bar{\alpha})) (f e_2). i} \text{EAPP} \\
\frac{\Gamma \vdash X :_{\text{p}} \Sigma_1 \rightsquigarrow e \quad \Gamma \vdash T \rightsquigarrow \exists \bar{\alpha}. \Sigma_2 \quad \Gamma \vdash \Sigma_1 \leq_{\bar{\alpha}} \Sigma_2 \rightsquigarrow \delta; f}{\Gamma \vdash X :_{>T} :_{i(\exists \bar{\alpha}. \Sigma_2)} \exists \bar{\alpha}. \Sigma_2 \rightsquigarrow \mathbf{pack} (\delta \bar{\alpha}, f e)} \text{ESEAL}
\end{array}$$

Bindings

$$\begin{array}{c}
\frac{\Gamma \vdash E :_{i} \exists \bar{\alpha}. \Sigma \rightsquigarrow e}{\Gamma \vdash X = E :_{i} \exists \bar{\alpha}. \{X : \Sigma\} \rightsquigarrow \mathbf{unpack} (\bar{\alpha}, x) = e \mathbf{in} \mathbf{pack} (\bar{\alpha}, \{X = x\})} \text{BVAR} \quad \frac{\Gamma \vdash E :_{i} \exists \bar{\alpha}. \{\bar{X} : \bar{\Sigma}\} \rightsquigarrow e}{\Gamma \vdash \mathbf{include} E :_{i} \exists \bar{\alpha}. \{X : \Sigma\} \rightsquigarrow e} \text{BINCL} \\
\frac{\Gamma \vdash B_1 :_{i_1} \exists \bar{\alpha}_1. \{\bar{X}_1 : \bar{\Sigma}_1\} \rightsquigarrow e_1 \quad \Gamma, \bar{\alpha}_1, \bar{X}_1 : \bar{\Sigma}_1 \vdash B_2 :_{i_2} \exists \bar{\alpha}_2. \{\bar{X}_2 : \bar{\Sigma}_2\} \rightsquigarrow e_2 \quad \bar{X}'_1 = \bar{X}_1 - \bar{X}_2 \quad \bar{X}'_1 : \bar{\Sigma}'_1 \subseteq \bar{X}_1 : \bar{\Sigma}_1}{\Gamma \vdash B_1; B_2 :_{i_1 \vee i_2} \exists \bar{\alpha}_1 \bar{\alpha}_2. \{\bar{X}'_1 : \bar{\Sigma}'_1, \bar{X}_2 : \bar{\Sigma}_2\} \rightsquigarrow \mathbf{unpack} (\bar{\alpha}_1, y_1) = e_1 \mathbf{in} \mathbf{let} \bar{X}_1 = y_1. \bar{X}_1 \mathbf{in} \mathbf{unpack} (\bar{\alpha}_2, y_2) = e_2 \mathbf{in} \mathbf{pack} (\bar{\alpha}_1 \bar{\alpha}_2, \{\bar{X}'_1 = y_1. \bar{X}'_1, \bar{X}_2 = y_2. \bar{X}_2\})} \text{BSEQ} \quad \frac{}{\Gamma \vdash \epsilon :_{\text{p}} \{\} \rightsquigarrow \{\}} \text{BEMPTY}
\end{array}$$

Subtyping

$$\begin{array}{c}
\Gamma \vdash \Xi \leq \Xi' \rightsquigarrow f := \Gamma \vdash \Xi \leq_{\epsilon} \Xi' \rightsquigarrow \text{id}; f \\
\frac{}{\Gamma \vdash \pi \leq \pi \rightsquigarrow \lambda x. x} \text{SPATH} \quad \frac{}{\Gamma \vdash \mathbf{bool} \leq \mathbf{bool} \rightsquigarrow \lambda x. x} \text{SBOOL} \\
\frac{\Gamma \vdash \Xi' \leq \Xi \rightsquigarrow f \quad \Gamma \vdash \Xi \leq \Xi' \rightsquigarrow f'}{\Gamma \vdash [= \Xi'] \leq [= \Xi] \rightsquigarrow \lambda x. [\Xi]} \text{STYPE} \quad \frac{\pi = \alpha \bar{\alpha}'}{\Gamma \vdash [= \sigma] \leq_{\pi} [= \pi] \rightsquigarrow [\lambda \bar{\alpha}'. \sigma / \alpha]; \lambda x. x} \text{SFORGET} \\
\frac{}{\Gamma \vdash \{\bar{l} : \bar{\Sigma}'\} \leq \{\} \rightsquigarrow \lambda x. \{\}} \text{SEMPY} \quad \frac{\Gamma \vdash \Sigma'_1 \leq_{\bar{\pi}_1} \Sigma_1 \rightsquigarrow \delta_1; f_1 \quad \Gamma \vdash \{\bar{l}' : \bar{\Sigma}'\} \leq_{\bar{\pi}_2} \{\bar{l} : \bar{\Sigma}\} \rightsquigarrow \delta_2; f_2 \quad \delta_2 \Sigma_1 = \Sigma_1}{\Gamma \vdash \{l_1 : \Sigma'_1, \bar{l}' : \bar{\Sigma}'\} \leq_{\bar{\pi}_1 \bar{\pi}_2} \{l_1 : \Sigma_1, \bar{l} : \bar{\Sigma}\} \rightsquigarrow \delta_1 \delta_2; \lambda x. \{l_1 = f_1(x.l_1), \bar{l} = (f_2 x). \bar{l}\}} \text{SSTR} \\
\frac{\Gamma, \bar{\alpha} \vdash \Sigma \leq_{\bar{\alpha}'} \Sigma' \rightsquigarrow \delta_1; f_1 \quad i' \leq i \quad \Gamma, \bar{\alpha} \vdash \delta_1 \Xi' \leq_{\bar{\pi} \bar{\alpha}} \Xi \rightsquigarrow \delta_2; f_2 \quad \delta_2 \Sigma = \Sigma}{\Gamma \vdash (\forall \bar{\alpha}'. \Sigma' \rightarrow_{i'} \Xi') \leq_{\bar{\pi}} (\forall \bar{\alpha}. \Sigma \rightarrow_i \Xi) \rightsquigarrow \delta_2; \lambda x. \lambda \bar{\alpha}. \lambda_i y : \Sigma. f_2 ((x (\delta_1 \bar{\alpha}')) (f_1 y)). i'} \text{SFUN} \quad \frac{\Gamma, \bar{\alpha}' \vdash \Sigma' \leq_{\bar{\alpha}} \Sigma \rightsquigarrow \delta; f \quad \bar{\alpha}' \bar{\alpha} \neq \epsilon}{\Gamma \vdash \exists \bar{\alpha}'. \Sigma' \leq \exists \bar{\alpha}. \Sigma \rightsquigarrow \lambda x. \mathbf{unpack} (\bar{\alpha}', y) = x \mathbf{in} \mathbf{pack} (\delta \bar{\alpha}, f y)} \text{SABS}
\end{array}$$

Figure 4. Elaboration of IML_{ex}

The only other non-editorial changes over [25] are that “**type** T ” is now handled as a first-class value, no longer tied to bindings, and that Booleans have been added as representatives of the core.

The rules collect all abstract types generated by an expression (e.g. by sealing or by functor application) into an existential package. This requires repeated unpacking and repacking of existentials created by constituent expressions. Moreover, the sequencing rule BSEQ combines two (n -ary) existentials into one.

It is an invariant of the expression elaboration judgement that $\iota = \text{I}$ if Ξ is not a concrete type Σ — i.e., abstract type “generation” is impure. Without this invariant, rule EFUN might form an invalid function type that is marked pure but yet has an inner existential quantifier (i.e., is “generative”). To maintain the invariant, both sealing (rule ESEAL) and conditionals (rule EIF) have to be deemed impure if they generate abstract types — enforced by the notation $\iota(\Xi)$ defined in Figure 3. In that sense, our notion of purity actually corresponds to the stronger property of *valuability* in the parlance of Dreyer [4], which also implies *phase separation*, i.e., the ability to separate static type information from dynamic computation, key to avoiding the need for dependent types.

Subtyping The subtyping judgement is defined on semantic types. It generates a coercion function f as computational evidence of the subtyping relation. The domain of that function always is the left-hand type Ξ' ; to avoid clutter, we omit its explicit annotation from the λ -terms in the rules. The rules mostly follow the structure from [25], merely adding a straightforward rule for abstract type paths π , which now may occur as “module types”.

However, we make one structural change: instead of guessing the substitution for the right-hand side’s abstract types non-deterministically in a separate rule (rule U-MATCH in [25]), the current formulation looks them up algorithmically as it goes, using the new rule SFORGET to match an individual abstract type. The reason for this change is merely a technical one: it eliminates the need for any significant meta-theory about decidability, which was somewhat non-trivial before, at least with applicative functors.

To this end, the judgement is indexed by a vector $\bar{\pi}$ of abstract paths that correspond to the abstract types from the right-hand Ξ . The counterparts of those types have to be looked up in the left-hand Ξ' , which happens in rule SFORGET. And that’s where the predicativity restriction materialises: the rule only allows a small type on the left. Lookup produces a substitution δ whose domain corresponds to the root variables of the abstract paths $\bar{\pi}$. Normally, each of $\bar{\pi}$ is just a plain abstract type variable (which occur free in Ξ in this judgement). But in the formation rule TPFUN for pure function types, lifting produces more complex paths. So when subtyping goes inside a pure functor in rule SFUN, the same abstract paths with skolem parameters have to be formed for lookup, so that rule SFORGET can match them accordingly.

The move to deterministic subtyping allows us to drop the auxiliary notion of *explicit* types, which was present in [25] to ensure that non-deterministic lookup can be made deterministic. There is one side effect from dropping the “explicitness” side condition from rule TSING, though: subtyping is no longer reflexive. There are now “monster” types that cannot be matched, not even by themselves. For example, take $\{\} \rightarrow_1 \exists \alpha. \alpha$, which is created by

$(= (\text{fun } (x : \{\}) \Rightarrow (\{\text{type } t = \text{int}; v = 0\} \rightarrow \{\text{type } t; v : t\}).v))$

However, this does not break anything else, so we make that simplification anyway (if desired, explicitness could easily be revived).

3.3 Meta-Theory

It is relatively straightforward to verify that elaboration is correct:

PROPOSITION 3.1 (Correctness of IML_{ex} Elaboration).
Let Γ be a well-formed F_ω environment.

1. If $\Gamma \vdash T/D \rightsquigarrow \Xi$, then $\Gamma \vdash \Xi : \Omega$.
2. If $\Gamma \vdash E/B : \iota \Xi \rightsquigarrow e$, then $\Gamma \vdash e : \Xi$, and if $\iota = \text{P}$ then $\Xi = \Sigma$.
3. If $\Gamma \vdash \Xi' \leq_{\bar{\alpha}\bar{\alpha}'} \Xi \rightsquigarrow \delta; f$ and $\Gamma \vdash \Xi' : \Omega$ and $\Gamma, \bar{\alpha} \vdash \Xi : \Omega$, then $\text{dom}(\delta) = \bar{\alpha}$ and $\Gamma \vdash \delta : \Gamma, \bar{\alpha}$ and $\Gamma \vdash f : \Xi' \rightarrow \delta \Xi$.

Together with the standard soundness result for F_ω we can tell that IML_{ex} is sound, i.e., a well-typed IML_{ex} program will either diverge or terminate with a value of the right type:

THEOREM 3.2 (Soundness of IML_{ex}). If $\cdot \vdash E : \Xi \rightsquigarrow e$, then either $e \uparrow$ or $e \hookrightarrow^* v$ such that $\cdot \vdash v : \Xi$.

More interestingly, the IML_{ex} type system is also decidable:

THEOREM 3.3 (Decidability of IML_{ex} Elaboration).
All IML_{ex} elaboration judgements are decidable.

This is immediate for all but the subtyping judgement, since they are syntax-directed and inductive, with no complicated side conditions. The rules can be read directly as an inductive algorithm. (In the case of **where**, it seems necessary to find a partitioning $\bar{\alpha}_1 = \bar{\alpha}_{11} \uplus \bar{\alpha}_{12}$, but it is not hard to see that the subtyping premise can only possibly succeed when picking $\bar{\alpha}_{12} = \text{fv}(\Sigma_1) \cap \bar{\alpha}_1$.)

The only tricky judgement is subtyping. Although it is syntax-directed as well, the rules are not actually inductive: some of their premises apply a substitution δ to the inspected types. Alas, that is exactly what can cause undecidability (see Section 1.2).

The restriction to substituting small types saves the day. We can define a weight metric over semantic types such that a quantified type variable has more weight than any possible substitution of that variable *with a small type*. We can then show that the overall weight of types involved decreases in all subtyping rules. For space reasons, the details appear in the Technical Appendix [23].

4. Full IML

A language without type inference is not worth naming ML. Because that is so, Figure 5 shows the minimal extension to IML_{ex} necessary to recover ML-style implicit polymorphism. Syntactically, there are merely two new forms of type expression.

First, “ $_$ ” stands for a type that is to be inferred from context. The crucial restriction here is that this can only be a *small* type. This fits nicely with the notion of a *monotype* in core ML, and prevents the need to infer polymorphic types in an analogous manner.

On top of this new piece of kernel syntax we allow a type annotation “ $_$ ” on a function parameter or conditional to be omitted, thereby recovering the implicitly typed expression syntax familiar from ML. (At the same time we drop the IML_{ex} sugar interpreting an unannotated parameter as a type; we only keep that interpretation in **type** declarations or bindings.)

Second, there is a new type of *implicit* function, distinguished by a leading tick ‘ (a choice that will become clear in a moment). This corresponds to an ML-style polymorphic type. The parameter has to be of type **type**, whose being small fits nicely with the fact that ML can only abstract monotypes, and no type constructors. For obvious reasons, an implicit function has to be pure. We write the semantic type of implicit functions with an arrow \rightarrow_A , in order to reuse notation. It is distinct from \rightarrow_ι , however, and A not an effect.

As the name would suggest, there are no explicit introduction or elimination forms for implicit functions. Instead, they are introduced and eliminated implicitly. The respective typing rules (EGEN and EINST) match common formulations of ML-style polymorphism [3]. Any pure expression can have its type generalised, which is more liberal than ML’s *value restriction* [35] (recall that purity also implies that no abstract types are produced).

Subtyping allows the implicit elimination of implicit functions as well, via instantiation on the left, or skolemisation on the right (rules SIMPLL and SIMPLR). This closely corresponds to ML’s

Syntax	(expressions)	$\text{if } E_1 \text{ then } E_2 \text{ else } E_3$	$:=$	$\text{if } E_1 \text{ then } E_2 \text{ else } E_3; _$	
(types)	T	$::= \dots \mid _ \mid '(X:\text{type}) \Rightarrow T$	(types)	$\text{fun } X \Rightarrow E$	$::= \text{fun } (X: _) \Rightarrow E$
			(types)	$'X \Rightarrow T$	$::= '(X:\text{type}) \Rightarrow T$
			(declarations)	$X 'Y:T$	$::= X : '(Y:\text{type}) \Rightarrow T$
Semantic Types	(large signatures)	Σ	$::= \dots \mid \forall \bar{\alpha}. \{ \} \rightarrow_A \Sigma$		
Types	$\frac{\Gamma \vdash \sigma : \Omega}{\Gamma \vdash _ \rightsquigarrow \sigma} \text{TINFER}$	$\frac{\Gamma, \alpha, X:[= \alpha] \vdash T \rightsquigarrow \Sigma \quad \kappa_\alpha = \Omega}{\Gamma \vdash '(X:\text{type}) \Rightarrow T \rightsquigarrow \forall \alpha. \{ \} \rightarrow_A \Sigma} \text{TIMPL}$	$\boxed{\Gamma \vdash T \rightsquigarrow \Xi}$		
Expressions	$\frac{\Gamma, \bar{\alpha} \vdash E :_p \Sigma \rightsquigarrow e \quad \overline{\kappa_\alpha = \Omega}}{\Gamma \vdash E :_p \forall \bar{\alpha}. \{ \} \rightarrow_A \Sigma \rightsquigarrow \lambda \bar{\alpha}. \lambda_A x. \{ \}. e} \text{EGEN}$	$\frac{\Gamma \vdash E :_i \exists \bar{\alpha}. \forall \bar{\alpha}'. \{ \} \rightarrow_A \Sigma \rightsquigarrow e \quad \overline{\Gamma, \bar{\alpha} \vdash \sigma : \kappa_{\alpha'}}}{\Gamma \vdash E :_i \exists \bar{\alpha}. \Sigma[\bar{\sigma}/\bar{\alpha}'] \rightsquigarrow \text{unpack } \langle \bar{\alpha}, x \rangle = e \text{ in pack } \langle \bar{\alpha}, (x \bar{\sigma} \{ \}). A \rangle} \text{EINST}$	$\boxed{\Gamma \vdash E :_i \Xi \rightsquigarrow e}$		
Subtyping	$\frac{\overline{\Gamma \vdash \sigma : \kappa_{\alpha'}} \quad \Gamma \vdash \Sigma'[\bar{\sigma}/\bar{\alpha}'] \leq_{\bar{\pi}} \Sigma \rightsquigarrow \delta; f}{\Gamma \vdash \forall \bar{\alpha}'. \{ \} \rightarrow_A \Sigma' \leq_{\bar{\pi}} \Sigma \rightsquigarrow \delta; \lambda x. f((x \bar{\sigma} \{ \}). A)} \text{SIMPLL}$	$\frac{\Gamma, \bar{\alpha} \vdash \Sigma' \leq_{\bar{\pi}} \Sigma \rightsquigarrow \delta; f \quad \overline{\text{fv}(\delta \pi) \not\cap \bar{\alpha}}}{\Gamma \vdash \Sigma' \leq_{\bar{\pi}} \forall \bar{\alpha}. \{ \} \rightarrow_A \Sigma \rightsquigarrow \delta; \lambda x. \lambda \bar{\alpha}. \lambda_A y. \{ \}. f x} \text{SIMPLR}$	$\boxed{\Gamma \vdash \Xi' \leq_{\bar{\pi}} \Xi \rightsquigarrow \delta; f}$		

Figure 5. Extension to Full IML

signature matching rules, which allow any value to be matched by a value of more polymorphic type. However, this behaviour can now be intermixed with proper “module” types. In particular, that means that we allow looking up types from an implicit function, similar to other pure functions. For example, the following subtyping holds, by implicitly instantiating the parameter a with int :

$$'(a : \text{type}) \Rightarrow \{ \text{type } t = a; f : a \rightarrow t \} \leq \{ \text{type } t; f : \text{int} \rightarrow \text{int} \}$$

With these few extensions, the Map functor from Section 2 can now be written in IML very much like in traditional ML:

```

type MAP =
{
  type key;
  type map a;
  empty 'a : map a;
  lookup 'a : key  $\rightarrow$  map a  $\rightarrow$  opt a;
  add 'a : key  $\rightarrow$  a  $\rightarrow$  map a  $\rightarrow$  map a
};
Map (Key : EQ) :> MAP where (type .key = Key.t) =
{
  type key = Key.t;
  type map a = key  $\rightarrow$  opt a;
  empty = fun x  $\Rightarrow$  none;
  lookup x m = m x;
  add x y m = fun z  $\Rightarrow$  if Key.eq z x then some y else m z
}

```

The MAP signature here uses one last bit of syntactic sugar defined in Figure 5, which is to allow implicit parameters on the left-hand side of declarations, like we already do for explicit parameters (cf. Figure 1). The tick becomes a pun on ML’s type variable syntax, but without relying on brittle implicit scoping rules.

Space reasons forbid more extensive examples, but it should be clear from the rules that there is nothing preventing the use of implicit functions as first-class values, given sufficient annotations for their (large) types. For example:

$$(\text{fun } (\text{id} : 'a \Rightarrow a \rightarrow a) \Rightarrow \{ x = \text{id } 3; y = \text{id } \text{true} \}) (\text{fun } x \Rightarrow x)$$

5. Type Inference

With the additions from Figure 5 we have turned the deterministic typing and elaboration judgements of IML_{ex} non-deterministic. They have to guess types (in rules TINFER, EINST, SIMPLL) and quantifiers (in rule EGEN). Clearly, an algorithm is needed.

Fortunately, what’s going on is not fundamentally different from core ML. Where core ML would require type equivalence (and type inference would use unification), the IML rules use subtyping.

That may seem scary at first, but a closer inspection of the subtyping rules reveals that, when applied to small types, subtyping almost degenerates to type equivalence! The only exception is width subtyping for records. The IML type system only promises to infer small types, so we are not far away from conventional ML. That is, we can still formulate an algorithm based on *inference variables* (which we write v) holding place for small types.

5.1 Algorithm

Figure 6 shows the essence of this algorithm, formulated via inference rules. The basic idea is to modify the declarative typing rules such that wherever they have to guess a (small) type, we simply introduce a (free) inference variable. Furthermore, the rules are augmented with outputting a substitution θ for resolved inference variables: all judgements have the form $\Gamma \vdash_{\theta} \mathcal{J}$, which, roughly, implies the respective declarative judgement $\bar{v}, \theta \Gamma \vdash \theta \mathcal{J}$, where \bar{v} binds the unresolved inference variables that still appear free in $\theta \Gamma$ or $\theta \mathcal{J}$. Notation is simplified by abbreviations of the form

$$\Gamma \vdash_{\theta'} \mathcal{J} \quad := \quad \theta \Gamma \vdash_{\theta''} \theta \mathcal{J} \wedge \theta' = \theta'' \circ \theta$$

where $\theta \mathcal{J}$ is meant to apply θ to \mathcal{J} ’s “inputs”. It’s used to thread and compose substitutions through multiple premises (e.g. rule IEIF).

There are two main complications, both due to the fact that, unlike in old ML, small types can be intermixed with large ones.

First, it may be necessary to infer a small type from a large one via subtyping. For example, we might encounter the inequation

$$\forall \alpha. [= \alpha] \rightarrow_p [= \alpha] \leq v$$

which can be solved just fine with $v = [= \sigma] \rightarrow_I [= \sigma]$ for any σ ; through contravariance, similar situations can arise with an inference variable on the left. Because of this, it is not enough to just consider the cases $v \leq \sigma$ or $\sigma \leq v$ for resolving v . Instead, when the subtyping algorithm hits $v \leq \Sigma$ or $\Sigma \leq v$ (rules ISRESL and ISRESR, where Σ may or may not be small) it invokes the auxiliary *Resolution* judgement $\Gamma \vdash_{\theta} v \approx \Sigma$, which only resolves v so far as to match the shape of Σ and inserts fresh inference variables for its subcomponents. After that, subtyping “tries again”.

Second, an inference variable v can be introduced in the scope of abstract types (i.e., regular type variables). In general, it would be incorrect to resolve v to a type containing type variables that are

Types	$\frac{\Gamma \vdash_{\theta}^! E :_{\mathbb{P}} [= \Xi]}{\Gamma \vdash_{\theta} E \rightsquigarrow \Xi} \text{ITPATH} \quad \frac{v \text{ fresh} \quad \Delta_v = \text{dom}(\Gamma)}{\Gamma \vdash_{\square} _ \rightsquigarrow v} \text{ITINFER} \quad \frac{}{\Gamma \vdash_{\square} \text{type} \rightsquigarrow \exists \alpha. [= \alpha]} \text{ITTYPE} \quad \frac{\Gamma \vdash_{\theta} E : \Sigma}{\Gamma \vdash_{\theta} (= E) \rightsquigarrow \Sigma} \text{ITSING}$	$\boxed{\Gamma \vdash_{\theta} T \rightsquigarrow \Xi}$
	$\frac{\Gamma \vdash_{\theta_1} T_1 \rightsquigarrow \exists \bar{\alpha}_1. \Sigma_1 \quad \Gamma; \bar{\alpha}_1, X : \Sigma_1 \vdash_{\theta_2} T_2 \rightsquigarrow \exists \bar{\alpha}_2. \Sigma_2 \quad \overline{\kappa_{\alpha'_2} = \bar{\kappa}_{\alpha_1} \rightarrow \bar{\kappa}_{\alpha_2}}}{\Gamma \vdash_{\theta_2} (X : T_1) \Rightarrow T_2 \rightsquigarrow \exists \bar{\alpha}'_2. \forall \bar{\alpha}_1. \Sigma_1 \rightarrow_{\mathbb{P}} \Sigma_2[\bar{\alpha}'_2 \bar{\alpha}_1 / \bar{\alpha}_2]} \text{ITPFUN} \quad \frac{\Gamma; \alpha, X : [= \alpha] \vdash_{\theta} T \rightsquigarrow \Sigma \quad \overline{\kappa_{\alpha} = \Omega}}{\Gamma \vdash_{\theta} '(X : \text{type}) \Rightarrow T \rightsquigarrow \forall \alpha. \{ \} \rightarrow_{\mathbb{A}} \Sigma} \text{ITIMPL}$	
Expressions	$\frac{\Gamma(X) = \Sigma}{\Gamma \vdash_{\square} X :_{\mathbb{P}} \Sigma} \text{IEVAR} \quad \frac{\Gamma \vdash_{\theta_0}^! X :_{\mathbb{P}} \text{bool} \quad \Gamma \vdash_{\theta_1} E_1 :_{\iota_1} \Xi_1 \quad \Gamma \vdash_{\theta_2} E_2 :_{\iota_2} \Xi_2 \quad \Gamma \vdash_{\theta_3} E_3 :_{\iota_3} \Xi_3 \leq \Xi_1 \quad \Gamma \vdash_{\theta_4} E_4 :_{\iota_4} \Xi_4 \leq \Xi_2}{\Gamma \vdash_{\theta_5} \text{if } X \text{ then } E_1 \text{ else } E_2 : T :_{\iota_1 \vee \iota_2 \vee \iota_3(\Xi)} \Xi} \text{IEIF} \quad \frac{\Gamma \vdash_{\theta}^! E :_{\iota} \exists \bar{\alpha}. \{ X : \Sigma, X' : \Sigma' \}}{\Gamma \vdash_{\theta} E.X :_{\iota} \exists \bar{\alpha}. \Sigma} \text{IEDOT}$	$\boxed{\Gamma \vdash_{\theta} E :_{\iota} \Xi}$
	$\frac{\Gamma \vdash_{\theta_1} T \rightsquigarrow \exists \bar{\alpha}. \Sigma \quad \Gamma; \bar{\alpha}, X : \Sigma \vdash_{\theta_2} E :_{\iota} \Xi}{\Gamma \vdash_{\theta_2} \text{fun } (X : T) \Rightarrow E :_{\mathbb{P}} \forall \bar{\alpha}. \Sigma \rightarrow_{\iota} \Xi} \text{IEFUN} \quad \frac{\Gamma \vdash_{\theta_1}^! X_1 :_{\mathbb{P}} \forall \bar{\alpha}. \Sigma_1 \rightarrow_{\iota} \Xi \quad \Gamma \vdash_{\theta_2} X_2 :_{\mathbb{P}} \Sigma_2 \quad \Gamma \vdash_{\theta_3} X_3 X_2 :_{\iota} \delta \leq \bar{\alpha}. \Sigma_1 \rightsquigarrow \delta}{\Gamma \vdash_{\theta_3} X_1 X_2 :_{\iota} \delta \Xi} \text{IEAPP}$	
Bindings	$\frac{\Gamma \vdash_{\theta} E :_{\mathbb{I}} \exists \bar{\alpha}. \Sigma}{\Gamma \vdash_{\theta} X = E :_{\mathbb{I}} \exists \bar{\alpha}. \{ X : \Sigma \}} \text{IBVAR} \quad \frac{\Gamma \vdash_{\theta} E :_{\mathbb{P}} \Sigma \quad \bar{v} = \text{undet}(\theta \Sigma) - \text{undet}(\theta \Gamma) \quad \overline{\kappa_{\alpha} = \Omega}}{\Gamma \vdash_{\theta} X = E :_{\mathbb{P}} \{ X : \forall \bar{\alpha}. \{ \} \rightarrow_{\mathbb{A}} \Sigma[\bar{\alpha} / \bar{v}] \}} \text{IBPVAR}$	$\boxed{\Gamma \vdash_{\theta} B :_{\iota} \Xi}$
Subtyping	$\frac{}{\Gamma \vdash_{\square} v \leq v} \text{ISREFL} \quad \frac{\Gamma \vdash_{\theta}^! v \approx \Sigma \quad \Gamma \vdash_{\theta'} v \leq \Sigma}{\Gamma \vdash_{\theta'} v \leq \Sigma} \text{ISRESL} \quad \frac{\Gamma \vdash_{\theta}^! v \approx \Sigma' \quad \Gamma \vdash_{\theta'} \Sigma' \leq v}{\Gamma \vdash_{\theta'} \Sigma' \leq v} \text{ISRESR}$	$\boxed{\Gamma \vdash_{\theta} \Xi \leq_{\bar{\pi}} \Xi' \rightsquigarrow \delta}$
	$\frac{\Gamma; \bar{\alpha} \vdash_{\theta_1} \Sigma \leq_{\bar{\alpha}'} \Sigma' \rightsquigarrow \delta_1 \quad \bar{v}' \leq \iota \quad \Gamma; \bar{\alpha} \vdash_{\theta_2} \delta_1 \Xi' \leq_{\bar{\pi} \bar{\alpha}} \Xi \rightsquigarrow \delta_2 \quad \theta_2 \delta_2 \Sigma = \theta_2 \Sigma}{\Gamma \vdash_{\theta_2} (\forall \bar{\alpha}'. \Sigma' \rightarrow_{\iota} \Xi') \leq_{\bar{\pi}} (\forall \bar{\alpha}. \Sigma \rightarrow_{\iota} \Xi) \rightsquigarrow \delta_2} \text{ISFUN} \quad \frac{\Gamma; \bar{\alpha} \vdash_{\theta} \Sigma' \leq_{\bar{\alpha}} \Sigma \rightsquigarrow \delta \quad \bar{\alpha}' \bar{\alpha} \neq \epsilon}{\Gamma \vdash_{\theta} \exists \bar{\alpha}'. \Sigma' \leq \exists \bar{\alpha}. \Sigma} \text{ISABS}$	
	$\frac{\bar{v} \text{ fresh} \quad \overline{\Delta_v = \text{dom}(\Gamma)} \quad \Gamma \vdash_{\theta} \Sigma'[\bar{v} / \bar{\alpha}'] \leq_{\bar{\pi}} \Sigma \rightsquigarrow \delta}{\Gamma \vdash_{\theta} \forall \bar{\alpha}'. \{ \} \rightarrow_{\mathbb{A}} \Sigma' \leq_{\bar{\pi}} \Sigma \rightsquigarrow \delta} \text{ISIMPLL} \quad \frac{\Gamma; \bar{\alpha} \vdash_{\theta} \Sigma' \leq_{\bar{\pi}} \Sigma \rightsquigarrow \delta; f \quad \bar{\alpha} \not\uparrow \text{fv}(\theta \delta)}{\Gamma \vdash_{\theta} \Sigma' \leq_{\bar{\pi}} \forall \bar{\alpha}. \{ \} \rightarrow_{\mathbb{A}} \Sigma \rightsquigarrow \delta} \text{ISIMPLR}$	
Resolution	$\Gamma \vdash_{\theta}^! v \approx \Sigma \quad := \quad v \notin \text{undet}(\Sigma) \wedge \Gamma \vdash_{\theta} v \approx \Sigma \quad \boxed{\Gamma \vdash_{\theta} v \approx \Sigma}$	
	$\frac{v'' \text{ fresh} \quad \Delta_{v''} = \Delta_v \cap \Delta_{v'}}{\Gamma \vdash_{[v''/v, v''/v']} v \approx v'} \text{IRINFER} \quad \frac{\alpha \in \Delta_v \quad \overline{v' \text{ fresh}} \quad \overline{\Delta_{v'} = \Delta_v}}{\Gamma \vdash_{[\alpha \bar{v}'/v]} v \approx \alpha \bar{\sigma}} \text{IRPATH}$	
	$\frac{}{\Gamma \vdash_{[\text{bool}/v]} v \approx \text{bool}} \text{IRBOOL} \quad \frac{v' \text{ fresh} \quad \Delta_{v'} = \Delta_v}{\Gamma \vdash_{[v'/v]} v \approx [= \Xi]} \text{IRTYPE} \quad \frac{v_1, v_2 \text{ fresh} \quad \Delta_{v_1} = \Delta_{v_2} = \Delta_v}{\Gamma \vdash_{[(v_1 \rightarrow_{\mathbb{I}} v_2)/v]} v \approx \forall \bar{\alpha}. \Sigma \rightarrow_{\iota} \Xi} \text{IRFUN}$	
Instantiation	$\Gamma \vdash_{\theta}^! E :_{\iota} \Xi \quad := \quad \Gamma \vdash_{\theta} E :_{\iota} \Xi' \wedge \Gamma \vdash_{\theta'} \Xi' \preceq \Xi \quad \boxed{\Gamma \vdash_{\theta} \Xi \preceq \Xi'}$	
	$\frac{}{\Gamma \vdash_{\theta} \Xi \preceq \Xi} \text{INREFL} \quad \frac{\Gamma; \bar{\alpha} \vdash_{\theta} v \approx \Sigma}{\Gamma \vdash_{\theta} \exists \bar{\alpha}. v \preceq \exists \bar{\alpha}. \Sigma} \text{INRES} \quad \frac{\bar{v} \text{ fresh} \quad \overline{\Delta_v = \text{dom}(\Gamma, \bar{\alpha})} \quad \Gamma \vdash_{\theta} \exists \bar{\alpha}. \Sigma[\bar{v} / \bar{\alpha}'] \preceq \exists \bar{\alpha}. \Sigma'}{\Gamma \vdash_{\theta} \exists \bar{\alpha}. \forall \bar{\alpha}'. \{ \} \rightarrow_{\mathbb{A}} \Sigma \preceq \exists \bar{\alpha}. \Sigma'} \text{INIMPL}$	

Figure 6. Type Inference for IML (Excerpt)

not in scope for *all* occurrences of v in a derivation. To prevent that, each v is associated with a set Δ_v of type variables that are known to be in scope for v everywhere. The set is verified when resolving v (see rule IRPATH in particular). The set also is propagated to any other v' the original v is unified with, by intersecting $\Delta_{v'}$ with Δ_v — or more precisely, by introducing a new variable v'' with the intersected $\Delta_{v''}$, and replacing both v and v' with it (see e.g. rule IRINFER); that way, we can treat Δ_v as a globally fixed set for each v , and do not need to maintain those sets separately. Inference variables also have to be updated when type variables go out of scope. That is achieved by employing the following notation in rules locally extending Γ with type variables (we write $\text{undet}(\Xi)$ to denote the free inference variables of Ξ):

$$\Gamma; \Gamma' \vdash_{\theta'} \mathcal{J} \quad := \quad \Gamma, \Gamma' \vdash_{\theta'} \mathcal{J} \wedge \theta' = [\bar{v}' / \bar{v}] \circ \theta''$$

where $\bar{v} = \text{undet}(\theta'' \mathcal{J})$

$$\overline{\bar{v}' \text{ fresh with } \Delta_{v'} = \Delta_v \cap \text{dom}(\Gamma)}$$

The net effect is that all local α 's from Γ' are removed from all Δ -sets of inference variable remaining after executing $\Gamma, \Gamma' \vdash \mathcal{J}$. We omit θ in this notation when it is the identity.

Implicit functions work mostly like in ML. Like with let-polymorphism, generalisation is deferred to the point where an expression is bound — in this case, in rule IBPVAR.

Similarly, instantiation is deferred to rules corresponding to elimination forms (e.g. IEIF, IEDOT, IEAPP, but also ITPATH). There, the auxiliary *Instantiation* judgement is invoked (as part of the notation $\Gamma \vdash_{\theta}^! \mathcal{J}$). This does not only instantiate implicit functions (possibly under existential binders), it also may resolve inference variables to create a type whose shape matches the shape that is expected by the invoking rule.

Instantiation can also happen implicitly as part of subtyping (rule ISIMPLL), which covers the case where a polymorphic value is supplied as the argument to a function expecting a monomorphic (or less polymorphic) parameter.

5.2 Incompleteness

There are a couple of sources of incompleteness in this algorithm:

Width subtyping Subtyping like $v \leq \{l : \sigma\}$ does not determine the shape of the record type that v stands for: the set of labels can

still vary. Consequently, the Resolution judgement has no rule for structures — instead a structure type must be determined by the previous context.

This is, in fact, similar to Standard ML [19], where record types cannot be inferred either, and require type annotation. However, SML implementations typically ensure that type inference is still order-independent, i.e., the information may be supplied *after* the point of use. They do so by employing a simple form of row inference. A similar approach would be possible for IML, but subtyping would still make more programs fail to type-check. For the sake of presentation, we decided to err on the side of simplicity.

The real solution of course would be to incorporate not just row inference but *row polymorphism* [21], so that width subtyping on structures can be recast as universal and existential quantification. We leave investigating such an extension for future work (though we note that **include** would still represent a challenge).

Type Scoping Tracking of the sets Δ_v is conservative: after leaving the scope of a type variable α , we exclude any solution for v that would still involve α , even if v only appears inside a type binder for α . Consider, for example [5]:

```
G (x : int) = {M = {type t = int; v = x} :> {type t; v : t}; f = id id};
C = G 3;
x = C.f (C.M.v);
```

and assume $\text{id} : '(a : \mathbf{type}) \Rightarrow a \rightarrow a$. Because id is impure, the definition of f is impure, and its type cannot be generalised; moreover, G is impure too. The algorithm will infer G 's type as

$$\text{int} \rightarrow \exists \beta. \{M : \{t : [= \beta], v : \beta\}, f : v \rightarrow_{\mathbf{I}} v\}$$

with $\beta \notin \Delta_v$ (because β goes out of scope the moment we bind it with a local quantifier), and then generalises to

$$G : \forall \alpha. \{ \} \rightarrow_{\mathbf{A}} \text{int} \rightarrow \exists \beta. \{M : \{t : [= \beta], v : \beta\}, f : \alpha \rightarrow_{\mathbf{I}} \alpha\}$$

But it's too late, the solution $v = \beta$, which would make x well-typed, is already precluded. When typing C , instantiating α with β is not possible either, because β can only come into scope again *after* having applied an argument for α already.

Although not well-known, this very problem is already present in good old ML, as Dreyer & Blume point out [5]: existing type inference implementations are incomplete, because combinations of functors and the value restriction (like above) do not have principal types. Interestingly, a variation of the solution suggested by Dreyer & Blume (implicitly generalising the types of functors) is implied by the IML typing rules: since functors are just functions, their types can already be generalised. However, generalisation happens outside the abstraction, which is more rigid than what they propose (but which is not expressible in System F_ω). Consequently, IML can type some examples from their paper, but not all.

Purity Annotations Due to effect subtyping, a function type as an upper bound does not determine the purity of a smaller type. Technically, that does not affect completeness, because we defined small types to only include impure functions: the resolution rule IRFUN can always pick \mathbf{I} . But arguably, that is cheating a little by side-stepping the issue, and it prevents the natural use of pure function types to specify “core-like” functions.

Again, the solution would be more polymorphism, in this case a simple form of effect polymorphism [32]. That will be future work.

Despite these limitations, we found IML inference quite usable. In practice, MLs have long given up on complete type inference. In our limited experience with a prototype, IML is not substantially worse, at least not when used in the same manner as traditional ML.

5.3 Metatheory

If the inference algorithm isn't complete, then at least it is sound. That is, we can show the following result:

THEOREM 5.1 (Correctness of IML Inference).

Let \bar{v}, Γ be a well-formed F_ω environment.

1. *If $\Gamma \vdash_{\theta} T/D \rightsquigarrow \Xi$, then $\bar{v}', \theta\Gamma \vdash T/D \rightsquigarrow \theta\Xi$.*
2. *If $\Gamma \vdash_{\theta} E/B :_i \Xi \rightsquigarrow e$, then $\bar{v}', \theta\Gamma \vdash E/B :_i \theta\Xi \rightsquigarrow \theta e$.*
3. *If $\Gamma \vdash_{\theta} \Xi' \leq_{\pi} \Xi \rightsquigarrow \delta; f$ and $\bar{v}, \Gamma \vdash \Xi' : \Omega$ and $\bar{v}, \Gamma, \bar{\alpha} \vdash \Xi : \Omega$, then $\bar{v}', \theta\Gamma \vdash \theta\Xi' \leq_{\pi} \theta\Xi \rightsquigarrow \theta\delta; \theta f$.*

THEOREM 5.2 (Termination of IML Inference).

All IML type inference judgements terminate.

We have to defer the details to the Technical Appendix [23].

6. Related Work

Packaged Modules The first concrete proposal for extending ML with packaged modules was by Russo [27], and is implemented in Moscow ML. Later work on type systems for modules routinely included them [6, 4, 24, 25], and variations have been implemented in other ML dialects, such as Alice ML [22] and OCaml [7].

To avoid soundness issues in the combination with applicative functors, Russo's original proposal conservatively allowed unpacking a module only local to core-level expressions, but this restriction has been lifted in later systems, restricting only the occurrence of unpacking inside applicative functors.

First-Class Modules The first to unify ML's stratified type system into one language was Harper & Mitchell's XML calculus [10]. It is a dependent type theory modeling modules as terms of Martin-Löf-style Σ and Π types, closely following MacQueen's original ideas [17]. The system enforces predicativity through the introduction of two universes U_1 and U_2 , which correspond directly to our notion of small and large type, and both systems allow both $U_1 : U_2$ and $U_1 \subseteq U_2$. XML lacks any account of either sealing or translucency, which makes it fall short as a foundation for modern ML.

That gap was closed by Harper & Lillibridge's calculus of *translucent sums* [9, 16], which also was a dependently typed language of first-class modules. Its main novelty were records with both opaque and transparent type components, directly modeling ML structures. However, unlike XML, the calculus is impredicative, which renders it undecidable.

Translucent sums were later superseded by the notion of *singleton types* [31]; they formed the foundation of Dreyer et al.'s type theory for higher-order modules [6]. However, to avoid undecidability, this system went back to second-class modules.

One concern in dependently typed theories is *phase separation*: to enable compile-time checking without requiring core-level computation, such theories must be sufficiently restricted. For example, Harper et al. [11] investigate phase separation for the XML calculus. The beauty of the F-ing approach is that it enjoys phase separation by construction, since it does not use dependent types.

Applicative Functors Leroy proposed applicative semantics for functors [15], as implemented in OCaml. Russo later combined both generative and applicative functors in one language [28] and implemented them in Moscow ML; others followed [30, 6, 4, 25].

A system like Leroy's, where *all* functors are applicative, would be incompatible with first-class modules, because the application in type paths like $F(A).t$ needs to be phase-separable to enable type checking, but not all functions are. Russo's system has similar problems, because it allows converting generative functors into applicative ones. Like Dreyer [4] or F-ing modules [25], IML hence combines applicative (pure) and generative (impure) functors such that applicative semantics is only allowed for functors whose body is both pure *and* separable. In F-ing modules, applicativity is even inferred from purity, and sealing itself not considered impure; the Technical Appendix [23] shows a similar extension to IML.

In the version of IML shown in the main paper, an applicative functor can only be created by sealing a fully transparent functor with pure function type, very much like in Shao’s system [30].

Type Inference There has been little work that has considered type inference for modules. Russo examined the interplay between core-level inference and modules [28], elegantly dealing with variable scoping via unification under a mixed prefix. Dreyer & Blume investigated how functors interfere with the value restriction [5].

At the same time, there have been ambitious extensions of ML-style type inference with higher-rank or impredicative types [8, 14, 33, 29]. Unlike those systems, IML never tries to infer a polymorphic type annotation: all guessed types are monomorphic and polymorphic parameters require annotation.

On the other hand, IML allows bundling types and terms together into structures. While it is necessary to explicitly annotate terms that contain types, associated type *quantifiers* (both universal and existential) and their actual introduction and elimination are implicit and effectively inferred as part of the elaboration process.

7. Future Work

IML, as shown here, is but a first step. There are several possible improvements and extensions.

Implementation We have implemented a simple prototype interpreter for IML (mpi-sws.org/~rossberg/1ml/), but it would be great to gather more experience with a “real” implementation.

Applicative Functors We would like to extend IML’s rather basic notion of applicative functor with *pure sealing* à la F-ing modules (see the Technical Appendix [23]), but more importantly, make it properly *abstraction-safe* by tracking value identities [25].

Implicits The domain of implicit functions in IML is limited to type **type**. Allowing richer types would be a natural extension, and might provide functionality like Haskell-style *type classes* [34].

Type Inference Despite the ability to express first-class and higher-order polymorphism, inference in IML is rather simple. Perhaps it is possible to combine IML elaboration with some of the more advanced approaches to inference described in literature.

More Polymorphism Replacing more of subtyping with polymorphism might lead to better inference: *row polymorphism* [21] could express width subtyping, and simple *effect polymorphism* [32] would allow more extensive use of pure function types.

Dependent Types Finally, IML goes to length to push the boundaries of non-dependent typing. It’s a legitimate question to ask, what for? Why not go fully dependent? Well, even then sealing necessitates some equivalent of weak sums (existential types). Incorporating them, along with the quantifier pushing of our elaboration, into a dependent type system might pose an interesting challenge.

References

- [1] H. Barendregt. Lambda calculi with types. In S. Abramsky, D. Gabbay, and T. Maibaum, editors, *Handbook of Logic in Computer Science*, vol. 2, chapter 2, pages 117–309. Oxford University Press, 1992.
- [2] S. K. Biswas. Higher-order functors with transparent signatures. In *POPL*, 1995.
- [3] L. Damas and R. Milner. Principal type-schemes for functional programs. In *POPL*, 1982.
- [4] D. Dreyer. *Understanding and Evolving the ML Module System*. PhD thesis, CMU, 2005.
- [5] D. Dreyer and M. Blume. Principal type schemes for modular programs. In *ESOP*, 2007.
- [6] D. Dreyer, K. Crary, and R. Harper. A type system for higher-order modules. In *POPL*, 2003.
- [7] J. Garrigue and A. Frisch. First-class modules and composable signatures in Objective Caml 3.12. In *ML*, 2010.
- [8] J. Garrigue and D. Rémy. Semi-explicit first-class polymorphism for ML. *Information and Computation*, 155(1-2), 1999.
- [9] R. Harper and M. Lillibridge. A type-theoretic approach to higher-order modules with sharing. In *POPL*, 1994.
- [10] R. Harper and J. C. Mitchell. On the type structure of Standard ML. In *ACM TOPLAS*, volume 15(2), 1993.
- [11] R. Harper, J. C. Mitchell, and E. Moggi. Higher-order modules and the phase distinction. In *POPL*, 1990.
- [12] R. Harper and B. Pierce. Design considerations for ML-style module systems. In B. C. Pierce, editor, *Advanced Topics in Types and Programming Languages*, chapter 8, pages 293–346. MIT Press, 2005.
- [13] R. Harper and C. Stone. A type-theoretic interpretation of Standard ML. In *Proof, Language, and Interaction: Essays in Honor of Robin Milner*. MIT Press, 2000.
- [14] D. Le Botlan and D. Rémy. MLF: Raising ML to the power of System F. In *ICFP*, 2003.
- [15] X. Leroy. Applicative functors and fully transparent higher-order modules. In *POPL*, 1995.
- [16] M. Lillibridge. *Translucent Sums: A Foundation for Higher-Order Module Systems*. PhD thesis, CMU, 1997.
- [17] D. MacQueen. Using dependent types to express modular structure. In *POPL*, 1986.
- [18] R. Milner. A theory of type polymorphism in programming languages. *JCSS*, 17:348–375, 1978.
- [19] R. Milner, M. Tofte, R. Harper, and D. MacQueen. *The Definition of Standard ML (Revised)*. MIT Press, 1997.
- [20] J. C. Mitchell and G. D. Plotkin. Abstract types have existential type. *ACM TOPLAS*, 10(3):470–502, July 1988.
- [21] D. Rémy. Records and variants as a natural extension of ML. In *POPL*, 1989.
- [22] A. Rossberg. The Missing Link – Dynamic components for ML. In *ICFP*, 2006.
- [23] A. Rossberg. IML – Core and modules as one (Technical Appendix), 2015. mpi-sws.org/~rossberg/1ml/.
- [24] A. Rossberg and D. Dreyer. Mixin’ up the ML module system. *ACM TOPLAS*, 35(1), 2013.
- [25] A. Rossberg, C. Russo, and D. Dreyer. F-ing modules. *JFP*, 24(5):529–607, 2014.
- [26] C. Russo. Non-dependent types for Standard ML modules. In *PPDP*, 1999.
- [27] C. Russo. First-class structures for Standard ML. *Nordic Journal of Computing*, 7(4):348–374, 2000.
- [28] C. Russo. Types for Modules. *ENTCS*, 60, 2003.
- [29] C. Russo and D. Vytiniotis. QML: Explicit first-class polymorphism for ML. In *ML*, 2009.
- [30] Z. Shao. Transparent modules with fully syntactic signatures. In *ICFP*, 1999.
- [31] C. A. Stone and R. Harper. Extensional equivalence and singleton types. *ACM TOCL*, 7(4):676–722, 2006.
- [32] J.-P. Talpin and P. Jouvelot. Polymorphic type, region and effect inference. *JFP*, 2(3):245271, 1992.
- [33] D. Vytiniotis, S. Weirich, and S. Peyton Jones. FPH: First-class polymorphism for Haskell. In *ICFP*, 2008.
- [34] P. Wadler and S. Blott. How to make ad-hoc polymorphism less ad hoc. In *POPL*, 1989.
- [35] A. Wright. Simple imperative polymorphism. *LASC*, 8:343–356, 1995.