

# *IML – 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 imposes seemingly unnecessary limits on expressiveness because it makes modules second-class citizens of the language. For example, selecting one among several possible modules implementing a given interface 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 and tend to be even more syntactically heavyweight than using second-class modules alone.

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 “IML”, functions, functors, and even type constructors are one and the same construct; likewise, no distinction is needed between structures, records, or tuples. Or viewed the other way round, everything is just (“a mode of use of”) modules. Yet, IML does not require dependent types: its type structure is expressible in terms of plain System  $F_\omega$ , in a minor variation of our *F-ing modules* approach.

We introduce both an explicitly typed version of IML, 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.

---

## 1 Introduction

The ML family of languages is defined by two splendid innovations: parametric polymorphism with Damas/Milner-style type inference (Milner, 1978; Damas & Milner, 1982), and an advanced module system based on concepts from dependent type theory (MacQueen, 1986). 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 total 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 demonstrated that ML-style module types can be readily expressed using only System-F-style quantification (Russo, 1999; Russo, 2003). 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  $F_\omega$  (Rossberg *et al.*, 2014).

Meanwhile, the second-class nature of modules is increasingly perceived as a practical limitation. The standard example is the desire 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. Yet, it is not expressible with conventional ML modules. Isn't that a shame!

### 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 (Russo, 2000; Dreyer *et al.*, 2003; Rossberg, 2006; Dreyer & Rossberg, 2008; Garrigue & Frisch, 2010; Rossberg *et al.*, 2014). 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 by heavy annotations. For example, in OCaml 4, which implements a simple form of packaged modules, 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.

But its not just notation, 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. For example, with proper modules, we can express the type of a functor that abstracts over *two* modules of a given signature with the requirement that their types  $t$  are equivalent:

```
signature S = {type t ...}
F : (X : S) → (Y : S with type t = X.t) → U
```

In the case of a core-level function abstracting over packaged modules the dependency on a parameter is not possible, so it gets a bit more complicated. We need to introduce an auxiliary polymorphic type variable (which is implicitly quantified):

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

Worse, 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 – it would require a type variable of higher kind, which are not supported in ML. In some cases the programmer may be able to work around this limitation resorting to weaker types. But in the general case they are forced to drop the use of packaged modules and raise 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, because all module-level computations are effectively “static”.

One could imagine addressing this particular limitation by introducing higher-kinded polymorphism into the ML core. But with such an extension type inference would require higher-order unification and hence become undecidable – unless accompanied by significant restrictions that are likely to defeat this example (or others).

### 1.2 First-Class Modules

Why not true first-class modules, then? Can we overcome the segregation and make modules more equal citizens of the language? This idea has been explored, of course. The answer from the literature so far has been: no, because first-class modules make type-checking undecidable and type inference infeasible.

The most directly relevant work is the calculus of *translucent sums* of Harper & Lillibridge (1994) (a precursor of later work on *singleton types* (Stone & Harper, 2006)). 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 by arbitrary transparent types. The latter feature usually manifests in a subtyping rule like the following,

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

which is, in some variation, at the heart of every definition of signature matching. In the premise the concrete type  $\tau$  is substituted for the abstract  $t$  in the remaining declarations  $D_1$  and  $D_2$  of the signatures. Obviously, this rule is not inductive. And that is a real problem: 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 (1994):

```
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 divergent example can be constructed on the module type level alone.

Some readers 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  $W$  (Milner, 1978) 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:

$$\begin{aligned} f_1 &: \{\mathbf{type} \ t \ a; \ x : t \ \mathit{int}\} \rightarrow \mathit{int} \\ f_2 &: \{\mathbf{type} \ t \ a; \ x : \mathit{int}\} \rightarrow \mathit{int} \\ g &= \mathbf{if} \ \mathit{condition} \ \mathbf{then} \ f_1 \ \mathbf{else} \ f_2 \end{aligned}$$

There are at least two possible types for  $g$ :

$$\begin{aligned} g &: \{\mathbf{type} \ t \ a = \mathit{int}; \ x : \mathit{int}\} \rightarrow \mathit{int} \\ g &: \{\mathbf{type} \ t \ a = a; \ x : \mathit{int}\} \rightarrow \mathit{int} \end{aligned}$$

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 the work on *F-ing modules* with Russo & Dreyer (Rossberg *et al.*, 2014) we have demonstrated that ML modules can be expressed and encoded entirely in vanilla System F (or  $F_\omega$ , 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.

Interestingly, in that paper we assume that the core language that modules sit on is simply System F (respectively,  $F_\omega$ ) itself. That leads to the seemingly paradoxical situation that the core language has *more* expressive types than the module language. But that makes sense when considering that the module 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” (and decidable) 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 the previous question is: very little! And the present article is all about explaining, exploring and exploiting that answer.

In essence, the F-ing modules semantics reveals that the syntactic stratification between ML core and module language is merely a rather blunt approach to enforce *predicativity* for module types: it prevents abstract types themselves from being instantiated with binders for abstract types. But this heavy *syntactic* restriction can be replaced by a much more surgical *semantic* restriction!

It is in fact enough to introduce a simple universe distinction between *small* and *large* types – reminiscent of Harper & Mitchell’s XML (1993) – and limit the equivalent of the FORGET rule shown earlier to only allow small types for substitution, which serves to exclude problematic occurrences of quantifiers:

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

In particular, a type like  $\{\mathbf{type} \ t; D\}$  is *not* small. The side condition thus prevents the substitution from introducing new occurrences of an opaque type specification, which would necessitate further (potentially never-ending) substitutions. That’s all, really.

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. If we limit instantiation of implicit polymorphic type variables (and thus inference and unification) to small types then type inference works almost as usual. The only exception to subtyping coinciding with equivalence on small types is width subtyping on structures – we hence need to be willing to accept that inference is not going to be complete for records (which it already isn’t in Standard ML or OCaml).

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. Our contributions are as follows:

- We present a *decidable type system* for a language of first-class modules that subsumes conventional second-class ML modules.
- We give an elaboration of this language into *plain System  $F_\omega$* .
- Conversely, we demonstrate that  $F_\omega$  can be *embedded* into this language.
- We show how Damas/Milner-style *type inference* can be integrated; it is incomplete, but only in ways that are already present in existing ML implementations.
- Besides theoretic considerations, we develop the basis for a *practical design* of an ML-like language in which the distinction between core and modules is eliminated.

This redesign has several benefits: it produces a language that is more *expressive* and *concise*, and at the same time, more *minimal* and *uniform*. “Modules” become a natural means of expressing all forms of polymorphism, both universal and existential. They can be freely intermixed with “computational” code and data, that is, polymorphism is fully first-class. 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

(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)	$D ::= X:T \mid \mathbf{include} T \mid D;D \mid \varepsilon$
(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)	$B ::= X=E \mid \mathbf{include} E \mid B;B \mid \varepsilon$

Fig. 1. Kernel syntax of  $\text{1ML}_{\text{ex}}$ 

from a small set of orthogonal constructs, over which general and uniform syntactic sugar can be defined.

## 2 1ML with Explicit Types

Despite the simplicity of the basic ideas, the dear reader will probably have guessed that the whole shebang isn't quite as trivial. To separate concerns a little, we will first introduce  $\text{1ML}_{\text{ex}}$ , a sublanguage of 1ML proper that is explicitly typed. We hold off fanciness like implicit typing and type inference until Sections 6 and 7.

The kernel syntax of  $\text{1ML}_{\text{ex}}$  is given in Figure 1. This kernel language is intentionally small, but features everything that is semantically relevant. In addition, Figure 2 defines a rich set of syntactic sugar over this kernel that provides many of the syntactic forms familiar from the ML family of languages.<sup>1</sup> In both figures, as well as in other places throughout this article, we use an overbar,  $\overline{\textit{phrase}}$ , to stand for an arbitrary number of repetitions of *phrase* (including zero of them).

Let us take a little tour of  $\text{1ML}_{\text{ex}}$  by way of examples.

**Functional Core** A major part of  $\text{1ML}_{\text{ex}}$  consists of fairly conventional functional language constructs. On the expression level, as a representative for a base type, we have Booleans (**true**, **false**, **if**) – in examples that follow, we will often assume the presence of an integer type and respective constructs as well. Then there are records  $\{\overline{B}\}$ , which consist of a sequence of bindings, and can be accessed via dot notation  $E.X$ . And of course, it wouldn't be a functional language without functions (**fun**) and their application  $(X_1 X_2)$ .

In a first approximation, these constructs are mirrored 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 (Rossberg *et al.*, 2014), 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 given in Figure 2. For example, the application

```
(fun (n : int) => n + n) 3
```

desugars into

```
let f = fun (n : int) => n + n; x = 3 in f x
```

<sup>1</sup> We omit further sugar that is standard and too boring to mention, such as tuples as records  $\{_{-}1 : T_1 ; \dots ; _n : T_n\}$ , pattern matching, or multi-argument functions as functions over tuples.

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

where (parameter)  $P ::= (X:T)$  with abbreviation (parameter)  $X ::= (X:\mathbf{type})$

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

Fig. 2. Syntactic abbreviations for  $\text{IML}_{\text{ex}}$

and further into

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

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  $\text{IML}_{\text{ex}}$ . 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 without a type annotation. However, in Section 6 we show that this annotation can usually be omitted in full IML with type inference.

**Reified Types** The core feature that makes  $\text{IML}_{\text{ex}}$  able to express modules is the ability to embed types in a first-class manner: the expression  $\mathbf{type} T$  reifies the type  $T$  as a value.<sup>2</sup> Such an expression has type  $\mathbf{type}$ , and thereby can be abstracted over. For example,

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

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

<sup>2</sup> Ideally, “ $\mathbf{type} T$ ” should be written just “ $T$ ”, like in dependently typed languages. However, that would create various syntactic ambiguities here, e.g. for phrases like “ $\{\}$ ”. The ambiguity 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 “ $E (\mathbf{type} T)$ ” if  $T$  unambiguously parses as a type.

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

```
pair = fun (a : type) ⇒ fun (b : type) ⇒ type {fst : a; 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, as a first class value. Again, the implicit projection from “path” expressions enables using this as a type:

```
second = fun (a : type) ⇒ fun (b : type) ⇒ fun (p : pair a b) ⇒ p.snd
```

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

Figure 2 also defines a bit of syntactic sugar for function and type definitions in “equational” form with parameters on the left, which makes them look more like in traditional ML. For example, the previous functions could equivalently be written as

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

For now, Figure 2 defines that omitting a type annotation on a parameter, as in the case of  $a$  and  $b$  above, makes it default to **type** – a rule that we will refine somewhat when we introduce implicit types in Section 6.

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 (Lillibridge, 1997; Dreyer *et al.*, 2003; Rossberg & Dreyer, 2013; Rossberg *et al.*, 2014). 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 a spurious module  $T = \{\mathbf{type} \ t = \text{int}\} : \{\mathbf{type} \ t\}$  and then refer to  $T.t$ , we allow injecting types into modules (i.e., values in IML) *anonymously*, without wrapping them into a structure; thus  $t = (\mathbf{type} \ \text{int}) : \mathbf{type}$ , which can be referred to as just  $t$ .

**Declarations** In the ML module system, *bindings* – i.e., definitions – in the module syntax are mirrored by *declarations* – i.e., specifications – in the signature syntax. That is no different in IML, where modules become records, such that bindings define record fields: analogously, signatures become record types, and declarations *specify* record fields.

For example, the previous definitions could be collected in a record with a type as follows:

```
type S =
{
  id : (a : type) ⇒ a → a
  pair : (a : type) ⇒ (b : type) ⇒ type
  second : (a : type) ⇒ (b : type) ⇒ pair a b → b
}
```

This amounts to a module signature, since it contains the type (constructor) `pair` as a component. And that is referred to in the type of `second`, which shows that record types



are seemingly “dependent”: like for terms, earlier components are in scope for later components – the key insight of the F-ing approach is that this dependency is benign, however, and can be translated away, as we will see in Section 4.

As for bindings, Figure 2 defines analogous syntactic sugar for declarations, which enables writing the above declarations in a more familiar style:

```
id a : a → a;
type pair a b;
second a b : pair a b → b
```

In particular, we introduce the same kind of sugar that allows putting parameters on the left of declarations as we do for bindings. And like before, type annotations on parameters can be omitted, defaulting to **type**. This design not just reconstructs the way type specifications are written in ML (as for pair).<sup>3</sup> The left-hand side parameter syntax uniformly generalises to ordinary value specifications and allows to rid ourselves of the brittle implicit scoping rules that conventional ML defines for polymorphic type variables (as for id and second).

**Translucency** The type **type** allows classifying types *opaquely*: given a value of type **type**, nothing is known about *what* type it is – it is abstract. But for modular programming it is essential that types can selectively be specified *transparently*, with a concrete definition. In particular, that enables expressing the vital concept of *type sharing* (Harper & Pierce, 2005).

As a simple example, consider these type aliases:

```
type size = int
type pair a b = {fst : a; snd : b}
```

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

```
size : type
pair : (a : type) ⇒ (b : type) ⇒ type
```

Or transparently:

```
size : (= type int)
pair : (a : type) ⇒ (b : type) ⇒ (= type {fst : a; snd : b})
```

The latter use a variant of *singleton types* (Stone & Harper, 2006; Dreyer *et al.*, 2003) 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 components of type **type**. It allows the type system to infer, for example, that the application pair size size is equivalent to the (reified) type {fst : int; snd : 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 declarations in signatures, where according to Figure 2 once more, they can be abbreviated to the more familiar

<sup>3</sup> Minus ML’s orthodox backwards type application syntax that is slightly less than popular.

```
type size = int
type pair a b = {fst : a; 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 article, this is effectively a shorthand for sharing all types contained in the module (including those defined inside transparent functors, see below). For example, given a signature

```
type T = {type t a; A : {type u; x : u}; f : A.u → t bool}
```

and a module  $M : T$ , the signature

```
type S = {B = M : T}
```

with its transparent module component is equivalent to the more explicit formulation

```
type S = {B : {type t a = M.t a; A : {type u = M.A.u; x : u}; f : A.u → t bool}}
```

which only makes the individual type components transparent, while not implying anything about the identity of its values. This behaviour matches the meaning of structure sharing in Standard ML '97 (Milner *et al.*, 1997).<sup>4</sup>

We leave the extension of 1ML to tracking full value equivalence, like in our F-ing semantics for applicative functors (Rossberg *et al.*, 2014), to future work.

**Functors** Returning to the 1ML grammar from Figure 1, 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. Here is the (unavoidable, it seems) functor that defines a simple map ADT:

```
type EQ =
{
  type t;
  eq : t → t → bool
};
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;
```

<sup>4</sup> But without disallowing transparent types in the signature of  $M$ , which is one of the more annoying limitations of SML

```

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 with an abstract type component  $t$ . Similarly, MAP defines a signature with abstract key and map types. 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 “ $T$  **where** (**type**  $X=T$ )” should be familiar from ML, but here it actually is derived from a more general construct: “ $T$  **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, with some appropriate sugar defined in Figure 2.

**Applicative vs. Generative** In this article, we stick to a relatively simple semantics for functor-like functions, in which Map is *generative* (Russo, 2003; Dreyer, 2005; Rossberg *et al.*, 2014). 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**  $\Rightarrow$  **type**. It is vital that the so sealed map functor – which returns an abstract type – has *applicative* semantics (Leroy, 1995; Rossberg *et al.*, 2014), 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* (Dreyer *et al.*, 2003). That is, the use of sealing *inside* a functor 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. Given:

```

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 = F :> type → type
K = 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 function type, the latter because

sealing occurs inside its body. Consequently,  $G \text{ int}$  or  $H \text{ int}$  are impure expressions and invalid as type paths (though it is fine to bind their result to a name, e.g., “**type**  $w = G \text{ int}$ ”, and use the constant  $w$  as a type). Applicative functor types are subtypes of generative ones, so that  $J$  turns  $F$  into a generative functor. Lastly,  $K$  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 (1999) – 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. This choice is mainly to keep this article focussed on the novelties of IML. More fully-featured applicative functors à la F-ing modules (Rossberg *et al.*, 2014), where sealing is pure, could easily be incorporated into IML, but complicate elaboration in ways that are largely orthogonal to our present goals.

**Effects** The ability to distinguish pure from impure, for computations or functions, is motivated by dealing with type abstraction correctly. But even in a traditional “core language” context this distinction is useful. For example, if we add impure constructs to the language, such as an ML-style type  $\text{ref } T$  of mutable references with the obvious operators:

$$\begin{aligned} \text{new} &: (a : \mathbf{type}) \Rightarrow (x : a) \rightarrow \text{ref } a \\ \text{rd} &: (a : \mathbf{type}) \Rightarrow (r : \text{ref } a) \rightarrow a \\ \text{wr} &: (a : \mathbf{type}) \Rightarrow (r : \text{ref } a) \Rightarrow a \rightarrow \{\} \end{aligned}$$

Their types make explicit where effects can occur, by mixing pure and impure arrows. A type parameter – expressing (explicit) polymorphism – is best considered a pure function: “instantiating” a polymorphic type should not have any effect. Similarly, an impure function with curried value parameters (like  $\text{wr}$ ) releases its effect only when the last argument is provided.

Expressing effects in function types is not just useful documentation, it can also aid type inference (see Section 6.1).

The simplistic effect system could naturally be refined to track relevant effects in a manner more fine-grained than with just a black & white distinction. However, we do not explore that space further here. Clearly, there is plenty of room for future work. We discuss one interesting direction, effect polymorphism, in a spin-off paper (Rossberg, 2016), see Section 8..

**Higher Polymorphism** So far, we have focussed on how IML reconstructs features that are well-known from ML. As a first example of something that cannot directly be expressed in conventional ML, consider first-class polymorphic arguments:

$$f (\text{id} : (a : \mathbf{type}) \Rightarrow a \rightarrow a) = \{x = \text{id int } 5; y = \text{id bool true}\}$$

The function  $f$  takes a function  $\text{id}$  as argument that is (explicitly) polymorphic, and used with different instantiations inside the body. That is not possible with the prenex polymorphism of core ML.

Similarly, existential types are directly expressible:

$$\begin{aligned} \mathbf{type} \text{ SHAPE} &= \{\mathbf{type} \text{ t}; \text{area} : \text{t} \rightarrow \text{float}\} \\ \mathbf{type} \text{ shape} &= \{S : \text{SHAPE}; v : S.t\} \\ \text{totalArea} &= \text{List.foldLeft shape float } (\mathbf{fun} (a : \text{float}) (x : \text{shape}) \Rightarrow a + x.S.\text{area } x.v) \end{aligned}$$

This example uses a two stage construction: SHAPE abstracts an ADT of shapes, which can be implemented e.g. as a circle or rectangle:

```
Circ : SHAPE = {type t = {radius : float}; area c = pi * c.radius * c.radius}
Rect : SHAPE = {type t = {width : float; height : float}; area r = r.width * r.height}
```

The type shape pairs such an ADT with an actual value:

```
aCirc = {S = Circ; v = {radius = 2.0}} :> shape
aRect = {S = Rect; v = {width = 1.8; height = 1.2}} :> shape
```

The magic here is that the “:>” operator corresponds directly to an introduction form for existential types. Because their concrete type is “forgotten” this way, shapes like these can be collected in a heterogeneous list, for which the function totalArea is able to compute the cumulative area:

```
totalArea [aCirc, aRect, aCirc]
```

This works because the computation inside the function is agnostic to the actual type of each shape, and the IML typing rules implicitly “open” the existential package. The upshot is that ADTs and existential types are one and the same thing in IML – driving home Mitchell & Plotkin’s original point (Mitchell & Plotkin, 1988) not just as a theoretical semantic observation, but as very concrete language design choice.

It turns out that the previous examples can still be expressed with packaged modules (Section 1.1), even though they would become fairly verbose.<sup>5</sup> 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 (Wadler & Blott, 1989). It contains two abstract type specifications, which are known as *associated types* in the type class literature (or in the C++ universe). The function entries is parameterised over a corresponding module C – an (explicit) instance of the type class 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 at the core level.<sup>6</sup>

<sup>5</sup> A welcome excuse to omit the boring details.

<sup>6</sup> In OCaml 4, this example can be approximated only 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 = ...
```

Moving to higher kinds, things become even more interesting:

```

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 constructor class, i.e., a type class over a type constructor. As explained in Section 1.1, this kind of polymorphism cannot be expressed in ML, even in dialects 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 `1MLex`:

```
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, the annotation is necessary in general to achieve unique types. But in most cases it can usually be inferred once we add inference to the mix (Section 6).

**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**; or in other words, monomorphic types. Small types thus exclude first-class abstract types, actual functors (functions taking type parameters), and type constructors (which are just functors). For example, all of the following define *large* types:

```

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

```

None of these are expressible as type expressions in conventional ML, and vice versa, all ML type expressions materialise as small types in `1ML`, so nothing is lost in comparison.

The restriction on subtyping affects annotations, parameterisation over types, and the formation of abstract types. For example, for all of the above  $T_i$ , none of the following definitions are well-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, it is the value of  $T_1$  above! But its only type is the transparent (= **type type**), it cannot be abstracted. Preventing a type: type situation is the main challenge with first-class modules, 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**  $T'_1 = (= \text{type int});$   
**type**  $T'_2 = \{\text{type } u = \text{int}\}$

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

**Recursion** The  $\text{IML}_{\text{ex}}$  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.

Recursive functions can be added simply by throwing in a primitive polymorphic fixpoint operator

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

On top of that, it is only a matter of defining suitable syntactic sugar:

**rec**  $X \bar{Y} (Z:T) : U = E \quad := \quad X = \text{fun } \bar{Y} \Rightarrow \text{fix } T U (\text{fun } (X : (Z:T) \rightarrow T') \Rightarrow \text{fun } (Z:T) \Rightarrow E)$

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

**rec**  $\text{stream } (a : \text{type}) : \text{type} = \text{type } \{\text{head} : a; \text{tail} : \text{stream } a\}$

won't type-check, because  $\text{stream}$  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  $\text{System } F_\omega$  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 (Harper & Stone, 2000), is to interpret a recursive datatype like

**datatype**  $t = A \mid B \text{ of } T$

as a module defining a primitive ADT with the signature

**{type**  $t; A : t; B : T \Rightarrow t; \text{expose } a : (\{\} \rightarrow a) \Rightarrow (T \rightarrow a) \Rightarrow t \rightarrow a$

where  $\text{expose}$  is a case-operator accessed by pattern matching compilation. We refer to Harper & Stone (2000) for more details on this approach. There is one caveat, though: datatypes expressed as ADTs require sealing. With the simple system presented in this article alone they hence could not be defined inside applicative functors. This limitation is removed by pure sealing, which, as mentioned before, is an obvious extension from the  $F$ -ing modules paper (Rossberg *et al.*, 2014) that is applicable to IML.

### 3 System $F_\omega$

Before diving deep into the semantics of IML in Section 4, a brief detour is in order. Following the *F-ing modules* approach (Rossberg *et al.*, 2014), IML's semantics is defined by translation into  $\text{System } F_\omega$ , the higher-order polymorphic  $\lambda$ -calculus (Barendregt, 1992).

Readers who know  $F_\omega$  in their sleep should feel free to skip over this sections.

(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)  $\Gamma ::= \cdot \mid \Gamma, \alpha:\kappa \mid \Gamma, x:\tau$

Fig. 3. Syntax of  $F_\omega$ 

**Syntax** Figure 3 gives the syntax of (impredicative) System  $F_\omega$  as we will use it here. It includes simple record types (where we assume that labels are always disjoint), but is otherwise completely standard. In the grammar, and elsewhere, we liberally use the meta-notation  $\overline{A}$  to stand for zero or more repetitions of an object or formula  $A$ , sometimes writing  $\varepsilon$  explicitly for the empty sequence. We also sometimes abuse the notation  $\overline{A}$  to actually denote the unordered set  $\{\overline{A}\}$ .

To differentiate the *external* language (IML) from the *internal* one ( $F_\omega$ ) later in the paper, we consistently use uppercase letters to range over phrases of the former ( $X, T, E$ , etc.), and lowercase for the latter ( $\alpha, x, \tau, e$ ).

**Static Semantics** The full typing rules of  $F_\omega$  are given in Figure 4. The judgements are as usual, with  $\Gamma \vdash \square$  denoting well-formedness of typing environments. Type equivalence is defined as full  $\beta\eta$ -equivalence.

The only point of note is that, unlike in many presentations, environments  $\Gamma$  permit term variables  $x$  (but not type variables  $\alpha$ ) to be shadowed without  $\alpha$ -renaming, which is convenient for translating  $\text{IML}_{\text{ex}}$ -bindings later. Thus, we take the notation  $\Gamma(x)$  to index the rightmost binding of  $x$  in  $\Gamma$ . Finally,  $\text{dom}(\Gamma)$  denotes the set of bound variables in  $\Gamma$ , while  $\text{fv}(\tau)$  and  $\text{fv}(e)$  are the free type variables of  $\tau$  or  $e$ .

**Dynamic Semantics** We assume a standard left-to-right call-by-value operational semantics, which is defined in Figure 5 via small-step reduction rules and evaluation contexts.

There is nothing interesting to see here, we can move on immediately to the meta-theory.

**Properties** The calculus as defined enjoys the standard soundness properties:

*Theorem 3.1 (Preservation)*

If  $\cdot \vdash e : \tau$  and  $e \hookrightarrow e'$ , then  $\cdot \vdash e' : \tau$ .

*Theorem 3.2 (Progress)*

If  $\cdot \vdash e : \tau$  and  $e \neq v$  for any  $v$ , then  $e \hookrightarrow e'$  for some  $e'$ .

The proofs are entirely standard.

The calculus also has the usual technical properties. The most relevant ones for our purposes are the following:

*Lemma 3.3 (Validity)*

1. If  $\Gamma \vdash \tau : \Omega$ , then  $\Gamma \vdash \square$ .
2. If  $\Gamma \vdash e : \tau$ , then  $\Gamma \vdash \tau : \Omega$ .

*Lemma 3.4 (Weakening)*

Let  $\Gamma' \supseteq \Gamma$  with  $\Gamma' \vdash \square$ .



**Environments** $\boxed{\Gamma \vdash \square}$ 

$$\frac{}{\cdot \vdash \square} \quad \frac{\Gamma \vdash \square \quad \alpha \notin \text{dom}(\Gamma)}{\Gamma, \alpha : \kappa \vdash \square} \quad \frac{\Gamma \vdash \tau : \Omega}{\Gamma, x : \tau \vdash \square}$$

**Types** $\boxed{\Gamma \vdash \tau : \kappa}$ 

$$\frac{\Gamma \vdash \tau_1 : \Omega \quad \Gamma \vdash \tau_2 : \Omega}{\Gamma \vdash \tau_1 \rightarrow \tau_2 : \Omega} \quad \frac{\overline{\Gamma \vdash \tau : \Omega} \quad \Gamma \vdash \square}{\Gamma \vdash \{l : \tau\} : \Omega}$$

$$\frac{\Gamma \vdash \square}{\Gamma \vdash \alpha : \Gamma(\alpha)} \quad \frac{\Gamma, \alpha : \kappa \vdash \tau : \Omega}{\Gamma \vdash \forall \alpha : \kappa. \tau : \Omega} \quad \frac{\Gamma, \alpha : \kappa \vdash \tau : \Omega}{\Gamma \vdash \exists \alpha : \kappa. \tau : \Omega}$$

$$\frac{\Gamma, \alpha : \kappa \vdash \tau : \kappa'}{\Gamma \vdash \lambda \alpha : \kappa. \tau : \kappa \rightarrow \kappa'} \quad \frac{\Gamma \vdash \tau_1 : \kappa' \rightarrow \kappa \quad \Gamma \vdash \tau_2 : \kappa'}{\Gamma \vdash \tau_1 \tau_2 : \kappa}$$

**Terms** $\boxed{\Gamma \vdash e : \tau}$ 

$$\frac{\Gamma \vdash \square}{\Gamma \vdash x : \Gamma(x)} \quad \frac{\Gamma \vdash e : \tau' \quad \tau' \equiv \tau \quad \Gamma \vdash \tau : \Omega}{\Gamma \vdash e : \tau}$$

$$\frac{\Gamma, x : \tau \vdash e : \tau'}{\Gamma \vdash \lambda x : \tau. e : \tau \rightarrow \tau'} \quad \frac{\Gamma \vdash e_1 : \tau' \rightarrow \tau \quad \Gamma \vdash e_2 : \tau'}{\Gamma \vdash e_1 e_2 : \tau}$$

$$\frac{\overline{\Gamma \vdash e : \tau} \quad \Gamma \vdash \square}{\Gamma \vdash \{l = e\} : \{l : \tau\}} \quad \frac{\Gamma \vdash e : \{l : \tau, l' : \tau'\}}{\Gamma \vdash e.l : \tau}$$

$$\frac{\Gamma, \alpha : \kappa \vdash e : \tau}{\Gamma \vdash \lambda \alpha : \kappa. e : \forall \alpha : \kappa. \tau} \quad \frac{\Gamma \vdash e : \forall \alpha : \kappa. \tau' \quad \Gamma \vdash \tau : \kappa}{\Gamma \vdash e \tau : \tau'[\tau/\alpha]}$$

$$\frac{\Gamma \vdash \tau : \kappa \quad \Gamma \vdash e : \tau'[\tau/\alpha] \quad \Gamma \vdash \exists \alpha : \kappa. \tau' : \Omega}{\Gamma \vdash \text{pack } \langle \tau, e \rangle_{\exists \alpha : \kappa. \tau'} : \exists \alpha : \kappa. \tau'}$$

$$\frac{\Gamma \vdash e_1 : \exists \alpha : \kappa. \tau' \quad \Gamma, \alpha : \kappa, x : \tau' \vdash e_2 : \tau \quad \Gamma \vdash \tau : \Omega}{\Gamma \vdash \text{unpack } \langle \alpha, x \rangle = e_1 \text{ in } e_2 : \tau}$$

**Type Equivalence** $\boxed{\tau \equiv \tau'}$ 

$$\frac{}{\tau \equiv \tau} \quad \frac{\tau' \equiv \tau}{\tau \equiv \tau'} \quad \frac{\tau \equiv \tau' \quad \tau' \equiv \tau''}{\tau \equiv \tau''}$$

$$\frac{\tau_1 \equiv \tau'_1 \quad \tau_2 \equiv \tau'_2}{\tau_1 \rightarrow \tau_2 \equiv \tau'_1 \rightarrow \tau'_2} \quad \frac{\overline{\tau \equiv \tau'}}{\{l : \tau\} \equiv \{l : \tau'\}}$$

$$\frac{\tau \equiv \tau'}{\forall \alpha : \kappa. \tau \equiv \forall \alpha : \kappa. \tau'} \quad \frac{\tau \equiv \tau'}{\exists \alpha : \kappa. \tau \equiv \exists \alpha : \kappa. \tau'}$$

$$\frac{\tau \equiv \tau'}{\lambda \alpha : \kappa. \tau \equiv \lambda \alpha : \kappa. \tau'} \quad \frac{\tau_1 \equiv \tau'_1 \quad \tau_2 \equiv \tau'_2}{\tau_1 \tau_2 \equiv \tau'_1 \tau'_2}$$

$$\frac{}{(\lambda \alpha : \kappa. \tau_1) \tau_2 \equiv \tau_1[\tau_2/\alpha]} \quad \frac{\alpha \notin \text{fv}(\tau)}{(\lambda \alpha : \kappa. \tau \alpha) \equiv \tau}$$

Fig. 4. Typing rules for  $F_\omega$

**Reduction** $e \hookrightarrow e'$ 

$$\frac{\begin{array}{l} (\lambda x:\tau.e) v \hookrightarrow e[v/x] \\ \{\overline{l_1=v_1}, l=v, \overline{l_2=v_2}\}.l \hookrightarrow v \\ (\lambda \alpha:\kappa.e) \tau \hookrightarrow e[\tau/\alpha] \\ \text{unpack } \langle \alpha, x \rangle = \text{pack } \langle \tau, v \rangle_{\tau'} \text{ in } e \hookrightarrow e[\tau/\alpha][v/x] \end{array}}{\text{unpack } \langle \alpha, x \rangle = \text{pack } \langle \tau, v \rangle_{\tau'} \text{ in } e \hookrightarrow e[\tau/\alpha][v/x]}$$

where:

(values)  $v ::= \lambda x:\tau.e \mid \{\overline{l=v}\} \mid \lambda \alpha:\kappa.e \mid \text{pack } \langle \tau, v \rangle_{\tau}$ (contexts)  $C ::= [] \mid C e \mid v C \mid \{\overline{l_1=v}, l=C, \overline{l_2=e}\} \mid C.l \mid C \tau \mid \text{pack } \langle \tau, C \rangle_{\tau} \mid \text{unpack } \langle \alpha, x \rangle = C \text{ in } e$ Fig. 5. Reduction rules for  $F_{\omega}$ 

1. If  $\Gamma \vdash \tau : \kappa$ , then  $\Gamma' \vdash \tau : \kappa$ .
2. If  $\Gamma \vdash e : \tau$ , then  $\Gamma' \vdash e : \tau$ .

**Lemma 3.5 (Strengthening)**Let  $\Gamma' \subseteq \Gamma$  with  $\Gamma' \vdash \square$  and  $D = \text{dom}(\Gamma) \setminus \text{dom}(\Gamma')$ .

1. If  $\Gamma \vdash \tau : \kappa$  and  $\text{fv}(\tau) \cap D = \emptyset$ , then  $\Gamma' \vdash \tau : \kappa$ .
2. If  $\Gamma \vdash e : \tau$  and  $\text{fv}(e) \cap D = \emptyset$ , then  $\Gamma' \vdash e : \tau$ .

**Theorem 3.6 (Uniqueness of types and kinds)**Assume  $\Gamma \vdash \square$ .

1. If  $\Gamma \vdash \tau : \kappa_1$  and  $\Gamma \vdash \tau : \kappa_2$ , then  $\kappa_1 = \kappa_2$ .
2. If  $\Gamma \vdash e : \tau_1$  and  $\Gamma \vdash e : \tau_2$ , then  $\tau_1 \equiv \tau_2$ .

Finally, all judgments of the  $F_{\omega}$  type system are decidable:**Theorem 3.7 (Decidability)**

1. It is decidable whether  $\Gamma \vdash \square$ .
2. It is decidable whether  $\Gamma \vdash \tau : \kappa$ .
3. It is decidable whether  $\Gamma \vdash e : \tau$ .
4. If  $\Gamma \vdash \tau_1 : \kappa$  and  $\Gamma \vdash \tau_2 : \kappa$ , it is decidable whether  $\tau_1 \equiv \tau_2$ .

**Notational shorthands** The rest of this article often deals with quantification over whole sequences  $\bar{\alpha}$  of type variables, and some other forms of lists of binders. The following (fairly standard) definitions for shortened  $n$ -ary notation will come in handy:

$$\begin{array}{llll} \forall \varepsilon.\tau := \tau & \exists \varepsilon.\tau := \tau & \lambda \varepsilon.\tau := \tau & \tau_0 \varepsilon := \tau_0 \\ \forall \bar{\alpha}.\tau := \forall \alpha_1.\forall \bar{\alpha}'.\tau & \exists \bar{\alpha}.\tau := \exists \alpha_1.\exists \bar{\alpha}'.\tau & \lambda \bar{\alpha}.\tau := \lambda \alpha_1.\lambda \bar{\alpha}'.\tau & \tau_0 \bar{\tau} := \tau_0 \tau_1 \bar{\tau} \\ \\ \lambda \varepsilon.e & := e & & \\ \lambda \bar{\alpha}.e & := \lambda \alpha_1.\lambda \bar{\alpha}'.e & & \\ e \varepsilon & := e & & \\ e \bar{\tau} & := e \tau_1 \bar{\tau} & & \\ \text{pack } \langle \varepsilon, e \rangle_{\exists \varepsilon.\tau_0} & := e & & \\ \text{pack } \langle \bar{\tau}, e \rangle_{\exists \bar{\alpha}.\tau_0} & := \text{pack } \langle \tau_1, \text{pack } \langle \bar{\tau}', e \rangle_{\exists \bar{\alpha}'.\tau_0[\tau_1/\alpha_1]} \rangle_{\exists \bar{\alpha}.\tau_0} & & \\ \text{unpack } \langle \varepsilon, x:\tau \rangle = e_1 \text{ in } e_2 & := \text{let } x:\tau = e_1 \text{ in } e_2 & & \\ \text{unpack } \langle \bar{\alpha}, x:\tau \rangle = e_1 \text{ in } e_2 & := \text{unpack } \langle \alpha_1, x_1 \rangle = e_1 \text{ in unpack } \langle \bar{\alpha}', x:\tau \rangle = x_1 \text{ in } e_2 & & \\ \text{let } \bar{x}:\bar{\tau} \equiv \bar{e}_1 \text{ in } e_2 & := (\lambda \bar{x}:\bar{\tau}.e_2) \bar{e}_1 & & \end{array}$$

(where  $\bar{\alpha} = \alpha_1 \bar{\alpha}'$  and  $\bar{\tau} = \tau_1 \bar{\tau}'$ )

Likewise,  $n$ -ary forms of other constructs (e.g. term abstraction and applications) are defined in the obvious way.

To ease notation, we often drop type annotations from `let`, `pack`, and `unpack` where clear from context.

**Substitution** Our semantics will also make fair use of parallel type substitutions. We write them as  $[\bar{\tau}/\bar{\alpha}]$ , assuming that both vectors have the same arity. Sometimes we use  $\delta$  to range over such substitutions.

The following definitions and lemmas are relevant:

*Definition 3.8 (Typing of Type Substitutions)*

Let  $\delta = [\bar{\tau}/\bar{\alpha}]$ . We write  $\Gamma' \vdash \delta : \Gamma$  if and only if

1.  $\Gamma' \vdash \square$ ,
2.  $\bar{\alpha} \subseteq \text{dom}(\Gamma)$ ,
3. for all  $\alpha \in \text{dom}(\Gamma)$ ,  $\Gamma' \vdash \delta \alpha : \Gamma(\alpha)$ ,
4. for all  $x \in \text{dom}(\Gamma)$ ,  $\Gamma' \vdash x : \delta(\Gamma(x))$ .

*Lemma 3.9 (Type Substitutions)*

Let  $\Gamma' \vdash \delta : \Gamma$ . Then:

1. If  $\Gamma \vdash \tau : \kappa$ , then  $\Gamma' \vdash \delta \tau : \kappa$ .
2. If  $\Gamma \vdash e : \tau$ , then  $\Gamma' \vdash \delta e : \delta \tau$ .

## 4 Type System and Elaboration

The general recipe for  $\text{IML}_{\text{ex}}$  is simple: take the semantics from F-ing modules (Rossberg *et al.*, 2014), 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.

### 4.1 Semantic Types

The entire semantics is defined by *elaborating*  $\text{IML}_{\text{ex}}$  types and terms directly into “equivalent” types and terms of System  $F_{\omega}$  as given in the previous section. However, IML types do not translate to arbitrary  $F_{\omega}$  types – instead, they map to types of very specific shape. The grammar of these *semantic types* is given in Figure 6 (“semantic” as opposed to the “syntactic” types of the surface language which they model). To avoid clutter, we omit all kind annotations on type variable binders. Where needed, we use the notation  $\kappa_{\alpha}$  to refer to the kind implicitly associated with  $\alpha$ .

Semantic types are the essence of the F-ing elaboration. To understand them, it is probably best to walk through a few explanatory examples that elaborate a concrete IML type  $T$  into a System  $F_{\omega}$  type  $\tau$ , a relation we write  $T \rightsquigarrow \tau$ .

- $\{x : \text{bool}, y : \text{int}\} \rightsquigarrow \{x : \text{bool}, y : \text{int}\}$

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

Desugarings into  $F_{\omega}$ :

(types)		(terms)	
$[= \tau]$	::= $\{\text{type} : \tau \rightarrow \{\}\}$	$[\tau]$	::= $\{\text{type} = \lambda x : \tau . \{\}\}$
$\tau_1 \rightarrow_l \tau_2$	::= $\tau_1 \rightarrow \{\bar{l} : \tau_2\}$	$\lambda_l x : \tau . e$	::= $\lambda x : \tau . \{\bar{l} : e\}$

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

Fig. 6. Semantic Types

Elaboration will translate this syntactic type into the semantic type So records map to records, no surprise there. We assume an implicit injection from IML identifiers  $X$  into  $F_{\omega}$  labels  $l$  (and variables  $x$ ), so we can conveniently treat any  $X$  as a label (or variable).

- $\{t : (= \text{type } \text{bool}); u : (= \text{type } t); x : t, y : u\} \rightsquigarrow \{t : [= \text{bool}], u : [= \text{bool}], x : \text{bool}, y : \text{bool}\}$

Structurally, the elaborated type still looks almost like the original, except that all references to the transparent type specifications are replaced with their definitions.

Of course, you may (and should!) wonder what this “[ $= \tau$ ]” is that represents transparent types. Because we keep saying that semantic types are just System F types, but this does not look like one. The answer can be found in Figure 6: it is simply a notational abbreviation that is encoded into bare System F. Its expansion does not actually matter much, it is merely a coding trick. All that matters is that it is a form of type that (1) uniquely determines  $\tau$ , and (2) is inhabited by a term that also uniquely determines  $\tau$ . A function serves the bill.<sup>7</sup> To make the encoding unambiguous, it is wrapped into a record type with a unique label “type” that we assume to be a reserved name disjoint from all valid IML identifiers.

- $\{t : \text{type}; u : (= \text{type } t); x : t; y : u\} \rightsquigarrow \exists \alpha . \{t : [= \alpha], u : [= \alpha], x : \alpha, y : \alpha\}$

This type is almost the same as before except that  $t$  now is an abstract type. The example demonstrates the main magic of elaboration: it inserts appropriate quantifiers to bind abstract types. Following Mitchell & Plotkin (1988), abstract types are represented by existentials. The type variable  $\alpha$  provides a “semantic” name for  $t$ , and all references to  $t$  are replaced by  $\alpha$ . This includes the specification of  $t$  itself, which can now be viewed as being transparently equal to  $\alpha$ . The upshot is that all reified types, whether transparent or opaque, are represented in a uniform manner. It also removes any dependency between fields of the record type – key to elaborating into bare System F without dependent types.

<sup>7</sup> Because all type constructors are separately represented as “functors” in IML, we have no need for types of higher kind in this notation, as was the case in Rossberg *et al.* (2014).

- $\{A : \{t : \mathbf{type}\}; B : \{u : \mathbf{type}; x : A.t; y : B.u\} \rightsquigarrow$   
 $\exists \alpha \beta. \{A : \{t : [= \alpha]\}, B : \{u : [= \beta], x : \alpha, y : \beta\}$

This one shows the handling of nested abstract types: inner quantifiers are collectively hoisted out of records. In particular, this makes them scope over the remainder of any enclosing record, which allows to reference the type variable across subcomponents – such as  $A.t$  and  $B.u$  in the example. A type may have arbitrary many abstract components, which all accumulate in an outermost quantifier.

Accordingly, the grammar of semantic types defines an *abstracted* type  $\Xi = \exists \bar{\alpha}. \Sigma$  that quantifies over all the abstract types (i.e., components of type **type**) from the underlying *concretised* type  $\Sigma$ , by naming them  $\bar{\alpha}$ .  $\Sigma$  itself contains no existential quantifiers – except as part of function types. So let us move to functions then.

- $(t : \mathbf{type}) \rightarrow \{u : \mathbf{type}; x : t, y : u\} \rightsquigarrow$   
 $\forall \alpha. [= \alpha] \rightarrow_{\mathbf{I}} \exists \beta. \{u : [= \beta], x : \alpha, y : \beta\}$

Conceptually, all function types map to polymorphic functions in  $F_{\omega}$ . Being in negative position, the quantifier for the abstract types from the parameter (here  $\alpha$ ) turns into a universal quantifier. It scopes over the *whole* arrow type and thereby allows the codomain to refer to it. Like for nested records, this hoisting avoids the need for dependent types.

For an *impure* function type like the above, the abstract types in the result are encapsulated by a local existential quantifier. It binds all those types that will be “generated” when applying the function. This representation of generative functor types as System F types of the form  $\forall \bar{\alpha}. \Sigma_1 \rightarrow \exists \bar{\beta}. \Sigma_2$  is due to Russo (1999; 2003), and another central ingredient of the F-ing modules semantics.

Arrow types are further annotated by a simple *effect*  $\eta$ , which distinguishes impure ( $\rightarrow_{\mathbf{I}}$ ) from pure ( $\rightarrow_{\mathbf{P}}$ ) function types. Again, this syntax is encoded into bare  $F_{\omega}$  types with a little record trick, see Figure 6 (we assume that record labels used with this form do not overlap with other labels).

- $(t : \mathbf{type}) \Rightarrow \{u : \mathbf{type}; x : t; y : u\} \rightsquigarrow$   
 $\exists \beta. \forall \alpha. [= \alpha] \rightarrow_{\mathbf{P}} \{u : [= \beta \alpha], x : \alpha, y : \beta \alpha\}$

This example shows the *pure* variant of the previous function type. Because it does not generate types at applications time, the existential quantifier changes position in the  $F_{\omega}$  translation: pure function types encode applicative semantics for the abstract types they return by having their existential quantifiers (here for  $\beta$ ) “lifted” over their parameters. Their general form is  $\exists \bar{\beta}. \forall \bar{\alpha}. \Sigma_1 \rightarrow \Sigma_2$ . Conceptually, the abstract types are generated when the function is *defined*, not when it is applied. We impose the syntactic invariant that a pure function type ( $\rightarrow_{\mathbf{P}}$ ) never has an existential quantifier right of the arrow.

However, the representation of an abstract type returned by a function may depend on the function’s parameters. To capture any potential dependencies, the  $\bar{\beta}$  in a pure functions type are skolemised over  $\bar{\alpha}$  (Biswas, 1995; Russo, 2003; Rossberg *et al.*, 2014). That is, the kinds of  $\bar{\beta}$  are of the form  $\bar{\kappa}_{\bar{\alpha}} \rightarrow \kappa$ , which is where higher kinds come into play. In return, all their occurrences are applied to  $\bar{\alpha}$ . In the concrete example, the kind of  $t$ ’s type variable  $\beta$  hence is  $\Omega \rightarrow \Omega$ , and all its references are of the form  $\beta \alpha$ . Once the function is applied to a concrete type, say  $\tau$ , the parameter variable  $\alpha$  is substituted, and all occurrences will have the form  $\beta \tau$ .

- $\{t : \mathbf{type} \Rightarrow \mathbf{type}; x : t \text{ int}; y : t \text{ bool}\} \rightsquigarrow$   
 $\exists \beta. \{t : \forall \alpha. [= \alpha] \rightarrow_P [= \beta \alpha], x : \beta \text{ int}, y : \beta \text{ bool}\}$

Pure functions over type **type** readily subsume abstract type constructors. As we will see in a short while, applying  $(t \text{ int})$  for the above substitutes  $\text{int}$  for  $\alpha$  and results in a value of semantic type  $[= \beta \text{ int}]$  – which is the representation of the syntactic type  $(= \mathbf{type} \ t \ \text{int})$ . Hence, the types of  $x$  and  $y$  in the example have the desired meaning.

We call the general form of abstract type variables applied to parameters *semantic paths*,<sup>8</sup> ranged over by  $\pi$ . Parameters are (only) introduced through pure function abstraction and the aforementioned kind raising that goes with it. In general, an abstract type that is the result of an application of a pure function  $F$  to a value  $M$  is represented by the application of the higher-kinded type variable representing the constructor to the concrete types  $\bar{\sigma}_M$  from the applied argument.

Because we enforce predicativity, these argument types have to be small. Figure 6 defines small types as a subgrammar  $\sigma$  of the general semantic types  $\Sigma$ . The central difference is that small types cannot contain quantifiers. Moreover, small function types are required to be impure, which will simplify type inference (Section 7).

A more systematic summary of the mapping between syntactic types  $T$  and semantic types  $\Xi$  is as follows:

$$\begin{array}{c}
 T \rightsquigarrow \exists \bar{\alpha}. \Sigma \\
 \hline
 (= \mathbf{type} \ T_1) \rightsquigarrow [= \exists \bar{\alpha}_1. \Sigma_1] \\
 \mathbf{type} \rightsquigarrow \exists \alpha. [= \alpha] \\
 \{X_1 : T_1; X_2 : T_2\} \rightsquigarrow \exists \bar{\alpha}_1 \bar{\alpha}_2. \{X_1 : \Sigma_1, X_2 : \Sigma_2\} \\
 (X : T_1) \rightarrow T_2 \rightsquigarrow \forall \bar{\alpha}_1. \Sigma_1 \rightarrow_I \exists \bar{\alpha}_2. \Sigma_2 \\
 (X : T_1) \Rightarrow T_2 \rightsquigarrow \exists \bar{\alpha}_2. \forall \bar{\alpha}_1. \Sigma_1 \rightarrow_P \Sigma_2 \\
 A.t \rightsquigarrow \alpha_{A.t} \\
 F(M) \rightsquigarrow \alpha_{F(-)} \bar{\sigma}_M
 \end{array}$$

Here, we assume that each constituent type  $T_i$  on the left-hand side is recursively mapped to a corresponding  $\exists \bar{\alpha}_i. \Sigma_i$  appearing on the right-hand side.

#### 4.2 Elaboration

The elaboration rules for  $1\text{ML}_{\text{ex}}$  are collected in Figures 7 and 8. There is one judgement for each syntactic class with the following regular structure:

$$\begin{array}{cc}
 \Gamma \vdash T \rightsquigarrow \Xi & \Gamma \vdash E :_{\eta} \Xi \rightsquigarrow e \\
 \Gamma \vdash D \rightsquigarrow \bar{\Xi} & \Gamma \vdash B :_{\eta} \bar{\Xi} \rightsquigarrow e
 \end{array}$$

In addition, there is an additional judgement defining subtyping over semantic types:

$$\Gamma \vdash \Xi_1 \leq_{\bar{\pi}} \Xi_2 \rightsquigarrow \delta; f$$

<sup>8</sup> Because they are the semantic analog to syntactic paths like  $F(M).t$  that appear in some previous type systems for higher-order modules (Leroy, 1995)

**Types**

$$\begin{array}{c}
\boxed{\Gamma \vdash T \rightsquigarrow \Xi} \\
\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} \\
\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_{\mathbf{I}} \exists \bar{\alpha}_2. \Sigma_2} \text{TFUN} \\
\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 [\alpha'_2 \bar{\alpha}_1 / \bar{\alpha}_2]} \text{TPFUN} \\
\frac{\Gamma \vdash E :_{\text{p}} \Sigma \rightsquigarrow e \quad E = E' : > T \vee E = E' : T \vee E = \mathbf{type} T}{\Gamma \vdash (= E) \rightsquigarrow \Sigma} \text{TSING} \\
\frac{\Gamma \vdash T_1 \rightsquigarrow \exists \bar{\alpha}_1. \Sigma_1 \quad \Gamma \vdash T_2 \rightsquigarrow \exists \bar{\alpha}_2. \Sigma_2 \quad \bar{\alpha}_1 = \bar{\alpha}_{11} \uplus \bar{\alpha}_{12} \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**

$$\begin{array}{c}
\boxed{\Gamma \vdash D \rightsquigarrow \Xi} \\
\frac{\Gamma \vdash T \rightsquigarrow \exists \bar{\alpha}. \Sigma}{\Gamma \vdash X : T \rightsquigarrow \exists \bar{\alpha}. \{X : \Sigma\}} \text{DVAR} \quad \frac{\Gamma \vdash T \rightsquigarrow \exists \bar{\alpha}. \{\bar{X} : \Sigma\}}{\Gamma \vdash \mathbf{include} T \rightsquigarrow \exists \bar{\alpha}. \{\bar{X} : \Sigma\}} \text{DINCL} \\
\frac{\Gamma \vdash D_1 \rightsquigarrow \exists \bar{\alpha}_1. \{\bar{X}_1 : \Sigma_1\} \quad \Gamma, \bar{\alpha}_1, \bar{X}_1 : \Sigma_1 \vdash D_2 \rightsquigarrow \exists \bar{\alpha}_2. \{X_2 : \Sigma_2\} \quad \bar{X}_1 \cap \bar{X}_2 = \emptyset}{\Gamma \vdash D_1; D_2 \rightsquigarrow \exists \bar{\alpha}_1 \bar{\alpha}_2. \{X_1 : \Sigma_1, X_2 : \Sigma_2\}} \text{DSEQ} \quad \frac{}{\Gamma \vdash \varepsilon \rightsquigarrow \{\}} \text{DEMPTY}
\end{array}$$

**Subtyping**

$$\begin{array}{c}
\boxed{\Gamma \vdash \Xi' \leq_{\bar{\pi}} \Xi \rightsquigarrow \delta; 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 \{l : \Sigma'\} \leq \{\} \rightsquigarrow \lambda x. \{\}} \text{SEMPY} \\
\frac{\Gamma \vdash \Sigma'_1 \leq_{\bar{\pi}_1} \Sigma_1 \rightsquigarrow \delta_1; f_1 \quad \Gamma \vdash \{l' : \Sigma'\} \leq_{\bar{\pi}_2} \{l : \Sigma\} \rightsquigarrow \delta_2; f_2 \quad \delta_2 \Sigma_1 = \Sigma_1}{\Gamma \vdash \{l_1 : \Sigma'_1, l' : \Sigma'\} \leq_{\bar{\pi}_1 \bar{\pi}_2} \{l_1 : \Sigma_1, l : \Sigma\} \rightsquigarrow \delta_1 \delta_2; \lambda x. \{l_1 = f_1(x.l_1), l = (f_2 x). l\}} \text{SSTR} \\
\frac{\Gamma, \bar{\alpha} \vdash \Sigma \leq_{\bar{\alpha}'} \Sigma' \rightsquigarrow \delta_1; f_1 \quad \eta' \leq \eta \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_{\eta'} \Xi') \leq_{\bar{\pi}} (\forall \bar{\alpha}. \Sigma \rightarrow_{\eta} \Xi) \rightsquigarrow \delta_2; \lambda x. \lambda \bar{\alpha}. \lambda \eta y : \Sigma. f_2((x(\delta_1 \bar{\alpha}')(f_1 y)). \eta')} \text{SFUN} \\
\frac{\Gamma, \bar{\alpha}' \vdash \Sigma' \leq_{\bar{\alpha}} \Sigma \rightsquigarrow \delta; f \quad \bar{\alpha}' \bar{\alpha} \neq \varepsilon}{\Gamma \vdash \exists \bar{\alpha}'. \Sigma' \leq \bar{\alpha}. \Sigma \rightsquigarrow \lambda x. \mathbf{unpack} \langle \bar{\alpha}', y \rangle = x \text{ in } \mathbf{pack} \langle \delta \bar{\alpha}, f y \rangle} \text{SABS}
\end{array}$$

Fig. 7. Elaboration of IML<sub>ex</sub> types

**Expressions**

$$\boxed{\Gamma \vdash E : \eta \ \Xi \rightsquigarrow e}$$

$$\frac{\Gamma(X) = \Sigma}{\Gamma \vdash X : \text{p } \Sigma \rightsquigarrow X} \text{EVAR} \quad \frac{\Gamma \vdash T \rightsquigarrow \Xi}{\Gamma \vdash \mathbf{type } T : \text{p } [= \Xi] \rightsquigarrow [\Xi]} \text{ETYPE}$$

$$\frac{}{\Gamma \vdash \mathbf{true} : \text{p } \text{bool} \rightsquigarrow \text{true}} \text{ETRUE} \quad \frac{}{\Gamma \vdash \mathbf{false} : \text{p } \text{bool} \rightsquigarrow \text{false}} \text{EFALSE}$$

$$\frac{\Gamma \vdash X : \text{p } \text{bool} \rightsquigarrow e \quad \Gamma \vdash E_1 : \eta_1 \ \Xi_1 \rightsquigarrow e_1 \quad \Gamma \vdash \Xi_1 \leq \Xi \rightsquigarrow f_1 \quad \Gamma \vdash T \rightsquigarrow \Xi \quad \Gamma \vdash E_2 : \eta_2 \ \Xi_2 \rightsquigarrow e_2 \quad \Gamma \vdash \Xi_2 \leq \Xi \rightsquigarrow f_2}{\Gamma \vdash \mathbf{if } X \mathbf{ then } E_1 \mathbf{ else } E_2 : T : \eta_1 \vee \eta_2 \vee \eta(\Xi) \ \Xi \rightsquigarrow \text{if } e \text{ then } f_1 \ e_1 \ \text{else } f_2 \ e_2} \text{EIF}$$

$$\frac{\Gamma \vdash B : \eta \ \Xi \rightsquigarrow e}{\Gamma \vdash \{B\} : \eta \ \Xi \rightsquigarrow e} \text{ESTR} \quad \frac{\Gamma \vdash E : \eta \ \exists \bar{\alpha}. \{X' : \Sigma'\} \rightsquigarrow e \quad X : \Sigma \in \overline{X' : \Sigma'}}{\Gamma \vdash E.X : \eta \ \exists \bar{\alpha}. \Sigma \rightsquigarrow \text{unpack } \langle \bar{\alpha}, y \rangle = e \text{ in pack } \langle \bar{\alpha}, y.X \rangle} \text{EDOT}$$

$$\frac{\Gamma \vdash T \rightsquigarrow \exists \bar{\alpha}. \Sigma \quad \Gamma, \bar{\alpha}, X : \Sigma \vdash E : \eta \ \Xi \rightsquigarrow e}{\Gamma \vdash \mathbf{fun} (X : T) \Rightarrow E : \text{p } \forall \bar{\alpha}. \Sigma \rightarrow \eta \ \Xi \rightsquigarrow \lambda \bar{\alpha}. \lambda \eta X : \Sigma. e} \text{EFUN}$$

$$\frac{\Gamma \vdash X_1 : \text{p } \forall \bar{\alpha}. \Sigma_1 \rightarrow \eta \ \Xi \rightsquigarrow e_1 \quad \Gamma \vdash X_2 : \text{p } \Sigma_2 \rightsquigarrow e_2 \quad \Gamma \vdash \Sigma_2 \leq \bar{\alpha} \ \Sigma_1 \rightsquigarrow \delta : f}{\Gamma \vdash X_1 X_2 : \eta \ \delta \Xi \rightsquigarrow (e_1 (\delta \bar{\alpha}) (f e_2)). \eta} \text{EAPP}$$

$$\frac{\Gamma \vdash X : \text{p } \Sigma_1 \rightsquigarrow e \quad \Gamma \vdash T \rightsquigarrow \exists \bar{\alpha}. \Sigma_2 \quad \Gamma \vdash \Sigma_1 \leq \bar{\alpha} \ \Sigma_2 \rightsquigarrow \delta : f}{\Gamma \vdash X : > T : \eta (\exists \bar{\alpha}. \Sigma_2) \ \exists \bar{\alpha}. \Sigma_2 \rightsquigarrow \text{pack } \langle \delta \bar{\alpha}, f e \rangle} \text{ESEAL}$$

**Bindings**

$$\boxed{\Gamma \vdash B : \eta \ \Xi \rightsquigarrow e}$$

$$\frac{\Gamma \vdash E : \eta \ \exists \bar{\alpha}. \Sigma \rightsquigarrow e}{\Gamma \vdash X = E : \eta \ \exists \bar{\alpha}. \{X : \Sigma\} \rightsquigarrow \text{unpack } \langle \bar{\alpha}, x \rangle = e \text{ in pack } \langle \bar{\alpha}, \{X = x\} \rangle} \text{BVAR}$$

$$\frac{\Gamma \vdash E : \eta \ \exists \bar{\alpha}. \{X : \Sigma\} \rightsquigarrow e}{\Gamma \vdash \mathbf{include } E : \eta \ \exists \bar{\alpha}. \{X : \Sigma\} \rightsquigarrow e} \text{BINCL} \quad \frac{}{\Gamma \vdash \varepsilon : \text{p } \{\} \rightsquigarrow \{\}} \text{BEMPTY}$$

$$\frac{\Gamma \vdash B_1 : \eta_1 \ \exists \bar{\alpha}_1. \{\overline{X}_1 : \overline{\Sigma}_1\} \rightsquigarrow e_1 \quad \overline{X}'_1 = \overline{X}_1 - \overline{X}_2 \quad \Gamma, \bar{\alpha}_1, \overline{X}_1 : \overline{\Sigma}_1 \vdash B_2 : \eta_2 \ \exists \bar{\alpha}_2. \{\overline{X}_2 : \overline{\Sigma}_2\} \rightsquigarrow e_2 \quad \overline{X}'_1 : \overline{\Sigma}'_1 \subseteq \overline{X}_1 : \overline{\Sigma}_1}{\Gamma \vdash B_1 ; B_2 : \eta_1 \vee \eta_2 \ \exists \bar{\alpha}_1 \bar{\alpha}_2. \{\overline{X}'_1 : \overline{\Sigma}'_1, \overline{X}_2 : \overline{\Sigma}_2\} \rightsquigarrow \text{unpack } \langle \bar{\alpha}_1, y_1 \rangle = e_1 \text{ in let } \overline{X}_1 = y_1. \overline{X}_1 \text{ in } \text{unpack } \langle \bar{\alpha}_2, y_2 \rangle = e_2 \text{ in } \text{pack } \langle \bar{\alpha}_1 \bar{\alpha}_2, \{\overline{X}'_1 = y_1. \overline{X}'_1, \overline{X}_2 = y_2. \overline{X}_2\} \rangle} \text{BSEQ}$$

Fig. 8. Elaboration of IML<sub>ex</sub> expressions

The rules may seem a bit intimidating at first, especially those for expressions and subtyping. But do not despair: that is only because elaboration defines both static semantics (typing) and dynamic semantics (operational meaning) at the same time. The latter is encapsulated in the greyed parts like “ $\rightsquigarrow e$ ”, and defines what System F term a construct translates to. If you are merely interested in *typing* IML then you can completely ignore those grey parts in all rules – just squint and make them disappear.

**Types and Declarations** The elaboration rules for types implement the translation scheme we laid out in the previous section. Their main job 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 F-ing modules, except that **type** is a free-standing construct instead of being tied to the syntax



of bindings, and IML’s “**where**” construct is slightly more general. Most notable are the following rules.

Using an expression  $E$  as a type path requires that  $E$  is pure and has a type of the form  $[= \Xi]$ . It denotes the very type  $\Xi$  – which may be anything from an abstract type variable to a large type (rule TPATH).

Function types have separate rules for impure (generative) functions (TFUN) and pure (applicative) ones (TPFUN). As explained above, pure function types have the existential quantifier for the  $\bar{\alpha}_2$  from the codomain lifted over the parameters, skolemising them accordingly.

Rule TSING handles “singleton” types. It corresponds to rule S-LIKE in Rossberg *et al.* (2014), except that we restrict it to expression with explicit types, in order to mesh well with type inference later. In return, we can drop the former side condition requiring  $\Sigma$  to be *explicit*. The rule simply infers the (unique) type  $\Sigma$  of the pure expression  $E$ . Note that this type is not allowed to have existential quantifiers, i.e.,  $E$  may not introduce local abstract types. All types  $[= \Xi]$  occurring in  $\Sigma$  thus are transparent.

The most complex rule is TWHERE. Unlike in conventional ML, where the **where** construct simply concretises one (or several) types via substitution, we allow it to specialise an arbitrary subcomponent of  $T_1$  to an arbitrary subtype. The original semantic type of that subcomponent, denoted by  $\Sigma_1.\bar{X}$  in the rule, is replaced by the semantic interpretation of  $T_2$ , written  $\Sigma_2[\bar{X} = \Sigma_2]$  (both these notations are defined in Figure 6). If  $T_2$  defines its own abstract types  $\bar{\alpha}_2$  then they are added to the list. At the same time, it may remove some of the abstract types  $\bar{\alpha}_1$  from the original type. To this end,  $\bar{\alpha}_1$  is partitioned into types  $\bar{\alpha}_{11}$  that remain abstract and  $\bar{\alpha}_{12}$  that are made concrete by the refinement (or just replaced by some of the  $\bar{\alpha}_2$ ). The last premise checks that  $\Sigma_2$  indeed matches  $\Sigma_1.\bar{X}$ , using the subtyping judgement explained below. This yields a substitution  $\delta$  for eliminating the previously abstract types  $\bar{\alpha}_{12}$ . The remaining  $\bar{\alpha}_{11}$  are left abstract, so their quantifiers are kept.

**Expressions and Bindings** The elaboration of expressions closely follows the rules from the first part of Rossberg *et al.* (2014), but adds the tracking of purity as in Section 7 of that paper. However, to keep the current article 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 a relaxed applicative functor semantics. The only other non-editorial change over F-ing modules – apart from the addition of Boolean expression forms – is that as for types, “**type**  $T$ ” is now handled as a first-class value, no longer tied to bindings.

The rules assign effects  $\eta$  and semantic types  $\Xi$  to IML expressions or bindings, and at the same time translate them to corresponding terms  $e$  of System  $F_\omega$  that define the operational semantics of IML. It is an invariant of the judgement that  $e$  has type  $\Xi$  in System  $F_\omega$ .

As we saw before, the main trick in interpreting IML’s module-like types is introducing quantifiers for abstract types. That is reflected in the typing of expressions: an expression that defines new abstract types will have an “abstracted” semantic type  $\Xi$  with existential quantifiers, one quantifier for each new type. When such an expression is nested into a larger expression then the rules have to propagate these quantifiers accordingly. For example, for a projection  $E.x$  to be well-typed,  $E$  obviously needs to have a type of

the form  $\{x : \Sigma, \dots\}$ ; the resulting type would be  $\Sigma$  then. However, if  $E$  creates abstract types locally, then its type will be of the form  $\exists \bar{\alpha}. \{x : \Sigma, \dots\}$  instead. The central idea taken from F-ing modules is to handle such types implicitly by extruding the existential quantifier automatically: that is, the projection  $E.x$  is well-typed and assigned type  $\exists \bar{\alpha}. \Sigma$ , with the same sequence of quantifiers. And so on for other constructs. On the term level, that corresponds to repeated unpacking and repacking of the respective existentials.

Moreover, the sequencing rule BSEQ combines two ( $n$ -ary) existentials into one – and is the only rule that does so, because the kernel syntax allows only variables in most places; the derived forms from Figure 1 all desugar into uses of sequencing.

For the categorically inclined, the existential quantifiers in an expression’s type  $\Xi$  act like an implicit monad of “abstraction effects”. There is no way to “escape” this monad: the sequencing rule BSEQ acts like its bind operator, and the only construct to suspend it is functional abstraction (rule EFUN). The only rules that effectively “run” such a monad are the type formation rules TPATH and TSING, and they require the abstraction effects to be encapsulated in the expression (i.e., no abstract types escape).<sup>9</sup>

It is another invariant of the expression elaboration judgement that  $\eta = \mathbb{I}$  in all cases where  $\Xi$  does *not* degenerate into a concrete type  $\Sigma$  with empty (i.e., no) quantifier – in other words, abstract type “generation” is indeed 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  $\eta(\Xi)$  defined in Figure 6. In that sense, our notion of purity actually corresponds to the stronger property of *valuability* in the parlance of Dreyer (2005), 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** In two places expression typing uses subtyping: for function application (rule EAPP) and for sealing (rule ESEAL). Subtyping is like signature matching in ML modules: it does not just allow for width subtyping on records, it can also substitute abstract types with concrete definitions. It is through that property, that the two occurrences of subtyping subsume universal elimination and existential introduction.

The subtyping judgement is defined on semantic types. It computes a substitution  $\delta$  for the matched abstract types from the right-hand side, and it generates a coercion function  $f$  as computational evidence of the subtyping relation. The domain of that function always is the left-hand type  $\Xi'$ ; to avoid clutter, we omit its explicit annotation from the  $\lambda$ -terms produced by the rules. The rules mostly follow the structure from F-ing modules, merely adding a straightforward rule for abstract type paths  $\pi$  (rule SPATH), because those can now occur as “module types”.

However, we make one structural change: instead of guessing the substitution  $\delta$  non-deterministically in a separate rule (rule U-MATCH in Rossberg *et al.* (2014)), the current

<sup>9</sup> More precisely, we actually have an infinite family of monads – one for each possible sequence of variable kinds – and the bind operator is heterogeneous, yielding the composition of two monads with composition being just nesting. See Rossberg (2016) for a more extensive discussion.

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 with the non-determinism 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  $\Xi$ . The counterparts of those types have to be looked up in the left-hand  $\Xi'$ , which happens one at a time in the rule SFORGET. Those lookups incrementally produce the substitution  $\delta$  whose domain corresponds to the (root variables of) the abstract paths  $\bar{\pi}$ . Substitutions are combined in rule SSTR, where the vector of type paths is partitioned over the record components as needed.

Normally, each path in  $\bar{\pi}$  is just a plain abstract type variable that occurs free in  $\Xi$ . But as we saw 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. For example, matching  $(a : \mathbf{type}) \Rightarrow (= \mathbf{type} \ a)$  against  $(a : \mathbf{type}) \Rightarrow \mathbf{type}$  translates to  $\forall \alpha'. [= \alpha'] \rightarrow_P [= \alpha'] \leq \exists \beta. \forall \alpha. [= \alpha] \rightarrow_P [= \beta \alpha]$  with  $\kappa_\beta = \Omega \rightarrow \Omega$ . Via rule SABS, this would try to match  $\forall \alpha'. [= \alpha'] \rightarrow_P [= \alpha'] \leq_\beta \forall \alpha. [= \alpha] \rightarrow_P [= \beta \alpha]$ , which involves looking up the higher-kinded  $\beta$ . After contra-variantly matching the function parameters, the second premise of rule SFUN then traverses into the functions' codomains with path  $\beta \alpha$ , trying to match  $[= \alpha] \leq_\beta \alpha [= \beta \alpha]$ . And indeed, that is the only way for SFORGET to “find” the path and yield  $\delta = [\lambda \alpha. \alpha / \beta]$ .

No types can actually be looked up from under a(nother) abstracted type  $\Xi$  with non-empty quantifier – because  $\bar{\pi}$  is required to be empty in the conclusion of rule SABS. Hence, for any subtyping derivation  $\Xi' \leq_{\bar{\pi}} \Xi$ , it is either the vector  $\bar{\pi}$  that is empty, or the quantifiers of both  $\Xi$  and  $\Xi'$ . Either case can actually occur when invoking the second premise of rule SFUN, depending on whether the function types are pure or not.

The side conditions of the form  $\delta \Sigma = \Sigma$  in rules SSTR and SFUN enforce that no variables in the domain of  $\delta$  occur in  $\Sigma$  (under the usual freshness assumptions for bound variables). That in turn ensures that all abstract variables are properly scoped, i.e., do not occur before a “definitional” occurrence in a type of the form  $[= \pi]$  that allows to look them up.

Finally, we note that rule SFORGET is the single place in the whole set of rules where the predicativity restriction materialises: the rule only allows a small type  $\sigma$  on the left, so that only small types can end up in the image of the substitution  $\delta$ .

### 4.3 Meta-Theory

It is relatively straightforward to verify by induction that elaboration is correct, i.e., produces only well-formed  $F_\omega$ -objects:

*Proposition 4.1 (Correctness of IML<sub>ex</sub> Elaboration)*

Let  $\Gamma \vdash \square$ .

1. If  $\Gamma \vdash T/D \rightsquigarrow \Xi$ , then  $\Gamma \vdash \Xi : \Omega$ .
2. If  $\Gamma \vdash E/B :_\eta \Xi \rightsquigarrow e$ , then  $\Gamma \vdash e : \Xi$ , and if  $\eta = P$  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 System  $F_\omega$  we can tell that  $1ML_{ex}$  is sound, i.e., a well-typed  $1ML_{ex}$  program will either diverge or terminate with a value of the right type. With the bare definition of  $1ML_{ex}$  we have given divergence cannot even arise, but a more complete ML-like language will surely have it.

*Theorem 4.2 (Soundness of  $1ML_{ex}$ )*

If  $\cdot \vdash E : \Xi \rightsquigarrow e$ , then either  $e \uparrow$  or  $e \hookrightarrow^* v$  such that  $\cdot \vdash v : \Xi$ .

More interestingly, the  $1ML_{ex}$  type system is also decidable:

*Theorem 4.3 (Decidability of  $1ML_{ex}$  Elaboration)*

All  $1ML_{ex}$  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. The one odd case is rule  $TWHERE$ , 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  $\delta$  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. With that, 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*. Because the subtyping rules always remove quantifiers before applying a substitution, the overall weight of types involved decreases in all subtyping rules.

One suitable choice of weight functions is mapping to ordered pairs as follows:

$$\begin{array}{ll}
W[\Xi] & = \langle Q[\Xi], S[\Xi] \rangle \\
Q[\text{bool}] & = 0 & S[\text{bool}] & = 1 \\
Q[\alpha] & = 0 & S[\alpha] & = 1 \\
Q[\pi \sigma] & = 0 & S[\pi \sigma] & = S[\pi] + S[\sigma] + 1 \\
Q[\lambda \bar{\alpha}. \sigma] & = 0 & S[\lambda \bar{\alpha}. \sigma] & = S[\sigma] + 1 \\
Q[=[\Xi]] & = 2 \cdot Q[\Xi] & S[=[\Xi]] & = S[\Xi] + 1 \\
Q[\{\}] & = 0 & S[\{\}] & = 1 \\
Q[\{l_0 : \Sigma_0, \bar{l} : \bar{\Sigma}\}] & = Q[\Sigma_0] + Q[\{\bar{l} : \bar{\Sigma}\}] & S[\{l_0 : \Sigma_0, \bar{l} : \bar{\Sigma}\}] & = S[\Sigma_0] + S[\{\bar{l} : \bar{\Sigma}\}] + 1 \\
Q[\forall \bar{\alpha}. \Sigma \rightarrow_\eta \Xi] & = |\bar{\alpha}| + Q[\Sigma] + Q[\Xi] & S[\forall \bar{\alpha}. \Sigma \rightarrow_\eta \Xi] & = S[\Sigma] + S[\Xi] + 1 \\
Q[\exists \bar{\alpha}. \Sigma] & = |\bar{\alpha}| + Q[\Sigma] & S[\exists \bar{\alpha}. \Sigma] & = S[\Sigma]
\end{array}$$

The definitions assume  $\beta\eta$ -normal form; in the case of paths and structures, they are inductive over the vector of arguments and components, respectively.

Intuitively, the first component of a weight pair,  $Q[\ ]$ , measures the number of quantifiers – or more precisely, the number of quantifiers that need to be eliminated in the subtyping rules, thus the seemingly odd definition for the case  $Q[=[\Xi]]$ . The second component,  $S[\ ]$ , is the actual syntactic size of the type expression. We apply addition point-wise to weights  $W[\ ]$ , and impose a lexicographic ordering on them:

$$\begin{aligned}
\langle q, s \rangle + \langle q', s' \rangle & := \langle q + q', s + s' \rangle \\
\langle q, s \rangle < \langle q', s' \rangle & :\Leftrightarrow q < q' \vee (q = q' \wedge s < s')
\end{aligned}$$

The most important property of this definition of weights is that  $Q[[\sigma]] = 0$  for all small types  $\sigma$ , because small types do not contain quantifiers. Consequently, any large type will have more weight than even the biggest small type. We extend the notion of small type to (higher-order) type constructors and substitutions: a type constructor is small if it is of the form  $\lambda \bar{\alpha}. \sigma$ ; a substitution  $\delta$  is small if  $\delta(\alpha)$  is small for all  $\alpha$ . Then, abbreviating  $|\text{dom}(\delta)|$  to  $|\delta|$ :

*Lemma 4.4 (Weight Reduction under Small Substitution)*

Let  $\delta$  be a small substitution.

1.  $Q[[\delta\Xi]] = Q[[\Xi]]$ .
2.  $W[[\delta\Xi]] < W[[\Xi]] + \langle 1, 0 \rangle$ .
3.  $W[[\delta\Xi]] \leq W[[\Xi]] + \langle |\delta|, 0 \rangle$ .

In more prose, the first property shows that substitution does not change the number of quantifiers; consequently, the second says that eliminating a quantifier is a strictly larger weight reduction than performing a substitution; and the third says that eliminating a sequence of quantifiers via substitution is a weight reduction or at least keeps the weight unchanged. Note how this lemma crucially relies on  $\delta$  being small. The properties are essential in proving the case of rule SFUN for the termination lemma:

*Lemma 4.5 (Termination of Algorithmic Subtyping)*

Let  $\Gamma \vdash \Xi' : \Omega$  and  $\Gamma, \bar{\alpha} \vdash \Xi : \Omega$  and  $\bar{\pi} = \bar{\alpha} \bar{\alpha}'$ . Then  $\Gamma \vdash \Xi' \leq_{\bar{\pi}} \Xi \rightsquigarrow \delta; f$  terminates.

The proof is by case analysis on the algorithm. In each rule, the weight  $W[[\Xi']] + W[[\Xi]] + \langle |\bar{\pi}|, 0 \rangle$  gets strictly smaller for each premise: either its  $Q$  component shrinks, because a quantifier is peeled, or its  $S$ , because the structure is otherwise simplified. In particular,  $Q$  shrinks in all cases where a substitution  $\delta$  is applied in a premise and  $S$  increases.

## 5 Expressiveness

### 5.1 Impredicativity Reloaded

Some readers will think that predicativity is a rather severe restriction. And they are right. So it is worth asking: Can we (re)enable impredicative type abstraction without breaking decidability?

Yes we can. One possibility is the common trick of piggy-backing datatypes: we could allow their data constructors to have large parameters. Because datatypes are *nominal* in ML, impredicativity is “hidden away” and does not interfere with subtyping, thereby avoiding the issue described in Section 1.2.

Here, we suggest a more orthogonal solution, namely introducing *structural* impredicative types. Figure 9 reveals how they can be defined as an easy extension to  $\text{IML}_{\text{ex}}$  and the rules we have previously given.

The trick is that a large type has to be injected into the universe of small types *explicitly*, by way of a special type “**wrap**  $T$ ” (not entirely dissimilar from the bracketted polytypes that occur in the semi-explicit first-class polymorphism of Garrigue & Rémy (1999)). Semantically, this type is represented by a new form “[ $\Xi$ ]” (with the now familiar coding trick of marking it through a special record) that may appear as a small type even when the

30

Andreas Rossberg

**Syntax** (types)  $T ::= \dots \mid \mathbf{wrap} T$   
(expressions)  $E ::= \dots \mid \mathbf{wrap} X:T \mid \mathbf{unwrap} X:T$

Abbreviations:  
 $\mathbf{wrap} E:T := \mathbf{let} X=E \mathbf{in} \mathbf{wrap} X:T$   
 $\mathbf{unwrap} E:T := \mathbf{let} X=E \mathbf{in} \mathbf{unwrap} X:T$

**Semantic Types** (large)  $\Sigma ::= \dots \mid [\Xi]$   
(small)  $\sigma ::= \dots \mid [\Xi]$

Desugarings:  
(types)  $[\Xi] := \{\mathbf{val} : \Xi\}$  (terms)  $[e] := \{\mathbf{val} = e\}$

**Types** 
$$\frac{\Gamma \vdash T \rightsquigarrow \Xi}{\Gamma \vdash \mathbf{wrap} T \rightsquigarrow [\Xi]} \text{TWRAP} \quad \boxed{\Gamma \vdash T \rightsquigarrow \Xi}$$

**Expressions** 
$$\boxed{\Gamma \vdash E :_{\eta} \Xi \rightsquigarrow e}$$

$$\frac{\Gamma \vdash X :_{\mathbf{p}} \Sigma \rightsquigarrow e \quad \Gamma \vdash T \rightsquigarrow [\Xi] \quad \Gamma \vdash \Sigma \leq \Xi \rightsquigarrow f}{\Gamma \vdash \mathbf{wrap} X:T :_{\mathbf{p}} [\Xi] \rightsquigarrow [fe]} \text{EWRAP}$$

$$\frac{\Gamma \vdash X :_{\mathbf{p}} [\Xi'] \rightsquigarrow e \quad \Gamma \vdash T \rightsquigarrow [\Xi] \quad \Gamma \vdash \Xi' \leq \Xi \rightsquigarrow f}{\Gamma \vdash \mathbf{unwrap} X:T :_{\eta([\Xi])} \Xi \rightsquigarrow f(e.\mathbf{val})} \text{EUNWRAP}$$

**Subtyping** 
$$\frac{}{\Gamma \vdash [\Xi] \leq [\Xi] \rightsquigarrow \lambda x : [\Xi]. x} \text{SWRAP} \quad \boxed{\Gamma \vdash \Xi' \leq \Xi \rightsquigarrow \delta; f}$$

Fig. 9. Extending IML<sub>ex</sub> with Impredicativity

inner type is large. The new type comes with respective expression forms for introducing and eliminating it.

The elaboration rules for wrapped types and expressions are straightforward. The crux of this approach is that subtyping does not extend to wrapped types (rule SWRAP): a wrapped type can only be matched by an *equivalent* type, not a subtype. This way, any infinite recursion is avoided; the weight function for the termination proof (Section 4.3) can trivially be extended to the new form of wrapped type as follows:

$$Q[[[\Xi]]] = 0 \quad S[[[\Xi]]] = 1$$

That goes through because SWRAP is simply an axiom, i.e., always a leaf node in any subtyping derivation.

The syntax we chose for wrapped types is reminiscent of packaged modules (Section 1.1), but the construct has far more specialised use cases in IML. In particular, wrapping is never needed if one merely wants to abstract over a *value* of large type (like is the case with packaged modules). It is only needed in the rare case where one wants to *type*-abstract over a large type.

As an example, consider a Church encoding of the option type:

```

type OPT =
{
  type opt a
  none a : opt a
  some a : a → opt a
  caseopt a b : opt a → b → (a → b) → b
}

Opt :> OPT =
{
  type opt a = wrap ((b : type) ⇒ b → (a → b) → b)
  none a = wrap (fun b (n : b) (s : a → b) ⇒ n) : opt a
  some a (x : a) = wrap (fun b (n : b) (s : a → b) ⇒ s x) : opt a
  caseopt a b (o : opt a) = (unwrap o : opt a) b
}

```

The representation of `opt` is a large type, so it has to be wrapped in order to “make it small” and allow matching its abstract declaration in the signature.

## 5.2 Embedding $F_\omega$

The elaboration of  $1ML_{ex}$  embeds the language into System  $F_\omega$ , showing that it is no more expressive than that calculus. We can also show that it is no *less* expressive, by doing the inverse: embedding  $F_\omega$  into  $1ML_{ex}$ .

Doing so is fairly straightforward: Figure 10 gives the canonical translation function  $\llbracket \_ \rrbracket$  for  $F_\omega$  kinds, types, and terms. It assumes a syntactic injection of  $F_\omega$  type and term variables  $\alpha$  and  $x$  into  $1ML$  identifiers  $X_\alpha$  and  $X_x$ , respectively.  $F_\omega$  kinds are translated to suitable  $1ML$  types.  $F_\omega$  types are translated to  $1ML$  expressions; in the case that the type has ground kind  $\Omega$ , the resulting expression will have a type of the form  $\llbracket = \sigma \rrbracket$ , so that it can be used as a  $1ML$  type; otherwise, it will be a function returning such a value. Finally,  $F_\omega$  terms are translated into  $1ML$  expressions, as is to be expected.

Notably, all translated  $F_\omega$  terms become  $1ML$  expressions with a small type. And consequently, ground types translate to small types as well. That is necessary to encompass the impredicativity of  $F_\omega$ . To achieve that for quantified types, the translation wraps the translation of every quantified type, using the extension from the previous Section 5.1. Wrapped expressions require a type annotation; we simply assume that the types  $\tau_e$  and  $\tau_{e_1}$  used in the figure are determined by context – as usual, this could be made more precise via a type-directed translation, but we gloss over this detail for the sake of simple exposition.

To state correctness of the embedding, we also need to define an embedding of  $F_\omega$  environments. It modifies variable bindings in such a way that they match the requirements of the  $1ML$  typing rules:

$$\begin{aligned}
\llbracket \cdot \rrbracket &= \cdot \\
\llbracket \Gamma, x : \tau \rrbracket &= \llbracket \Gamma \rrbracket, X_x : \sigma \quad \text{if } \llbracket \Gamma \rrbracket \vdash \llbracket \tau \rrbracket : \llbracket = \sigma \rrbracket \\
\llbracket \Gamma, \alpha : \kappa \rrbracket &= \llbracket \Gamma \rrbracket, \alpha : \kappa, X_\alpha : \Sigma \quad \text{if } \llbracket \Gamma \rrbracket \vdash \llbracket \kappa \rrbracket \rightsquigarrow \exists \alpha : \kappa. \Sigma
\end{aligned}$$

With this, we can show that the embedding is sound, i.e., produces well-formed  $1ML$  programs from well-formed  $F_\omega$  terms:

*Theorem 5.1 (Soundness of Embedding)*

(kinds)	
$\llbracket \Omega \rrbracket$	$= \mathbf{type}$
$\llbracket \kappa_1 \rightarrow \kappa_2 \rrbracket$	$= \llbracket \kappa_1 \rrbracket \Rightarrow \llbracket \kappa_2 \rrbracket$
(types)	
$\llbracket \alpha \rrbracket$	$= X_\alpha$
$\llbracket \tau_1 \rightarrow \tau_2 \rrbracket$	$= \mathbf{type} (\llbracket \tau_1 \rrbracket \rightarrow \llbracket \tau_2 \rrbracket)$
$\llbracket \tau_1 \times \tau_2 \rrbracket$	$= \mathbf{type} (\llbracket \tau_1 \rrbracket \times \llbracket \tau_2 \rrbracket)$
$\llbracket \forall \alpha : \kappa. \tau \rrbracket$	$= \mathbf{type} (\mathbf{wrap} ((X_\alpha : \llbracket \kappa \rrbracket) \Rightarrow \llbracket \tau \rrbracket))$
$\llbracket \exists \alpha : \kappa. \tau \rrbracket$	$= \mathbf{type} (\mathbf{wrap} ((X_\alpha = \llbracket \kappa \rrbracket) \times \llbracket \tau \rrbracket))$
$\llbracket \lambda \alpha : \kappa. \tau \rrbracket$	$= \mathbf{fun} (X_\alpha : \llbracket \kappa \rrbracket) \Rightarrow \llbracket \tau \rrbracket$
$\llbracket \tau_1 \tau_2 \rrbracket$	$= \llbracket \tau_1 \rrbracket \llbracket \tau_2 \rrbracket$
(expressions)	
$\llbracket x \rrbracket$	$= X_x$
$\llbracket \lambda x : \tau. e \rrbracket$	$= \mathbf{fun} (X_x : \llbracket \tau \rrbracket) \Rightarrow \llbracket e \rrbracket$
$\llbracket e_1 e_2 \rrbracket$	$= \llbracket e_1 \rrbracket \llbracket e_2 \rrbracket$
$\llbracket \langle e_1, e_2 \rangle \rrbracket$	$= (\llbracket e_1 \rrbracket, \llbracket e_2 \rrbracket)$
$\llbracket e.i \rrbracket$	$= \llbracket e \rrbracket.i$
$\llbracket \lambda \alpha : \kappa. e \rrbracket$	$= \mathbf{wrap} (\mathbf{fun} (X_\alpha : \llbracket \kappa \rrbracket) \Rightarrow \llbracket e \rrbracket) : (X_\alpha : \llbracket \kappa \rrbracket) \Rightarrow \llbracket \tau_e \rrbracket$
$\llbracket e_1 \tau_2 \rrbracket$	$= (\mathbf{unwrap} \llbracket e_1 \rrbracket : \llbracket \tau_{e_1} \rrbracket) \llbracket \tau_2 \rrbracket$
$\llbracket \langle \tau_1, e_2 \rangle \tau \rrbracket$	$= \mathbf{wrap} (\llbracket \tau_1 \rrbracket, \llbracket e_2 \rrbracket) : \llbracket \tau \rrbracket$
$\llbracket \mathbf{unpack} \langle \alpha, x \rangle = e_1 \text{ in } e_2 \rrbracket$	$= \mathbf{let} (X_\alpha, X_x) = \mathbf{unwrap} \llbracket e_1 \rrbracket : \llbracket \tau_{e_1} \rrbracket \text{ in } \llbracket e_2 \rrbracket$
where:	
$(X = T_1) \times T_2$	$:= \{ \_1 : T_1; \_2 : \mathbf{let} X = \_1 \text{ in } T_2 \}$
$(E_1, E_2)$	$:= \{ \_1 = E_1; \_2 = E_2 \}$
$E.i$	$:= E.\_i$
$\mathbf{let} (X_1, X_2) = E_1 \text{ in } E_2$	$:= \mathbf{let} X = E_1; X_1 = X.\_1; X_2 = X.\_2 \text{ in } E_2$

Fig. 10. Embedding of  $F_\omega$  in 1ML

1. For all  $\kappa, \cdot \vdash \llbracket \kappa \rrbracket \rightsquigarrow \exists \alpha. \Sigma$  with  $\kappa_\alpha = \kappa$  and  $\cdot \vdash \exists \alpha. \Sigma : \Omega$ .
2. For all  $\Gamma$ , if  $\Gamma \vdash \square$ , then  $\llbracket \Gamma \rrbracket \vdash \square$ .
3. For all  $\tau$ , if  $\Gamma \vdash \tau : \kappa$ , then  $\llbracket \Gamma \rrbracket \vdash \llbracket \tau \rrbracket : \Sigma$  and  $\llbracket \Gamma \rrbracket \vdash \llbracket \kappa \rrbracket \rightsquigarrow \Xi$  and  $\llbracket \Gamma \rrbracket \vdash \Sigma \leq \Xi$   
(the latter implies that if  $\kappa = \Omega$  then  $\Sigma = [= \sigma]$ ).
4. For all  $e$ , if  $\Gamma \vdash e : \tau$ , then  $\llbracket \Gamma \rrbracket \vdash \llbracket e \rrbracket : \sigma$  and  $\llbracket \Gamma \rrbracket \vdash \llbracket \tau \rrbracket : [= \sigma]$ .

The proof of (1) is by induction on the structure of  $\kappa$ , and the rest by (simultaneous) induction on the first derivations, respectively.

For a pedantic proof of correctness we would also need to show that the embedding is adequate, i.e., the produced  $1ML_{\text{ex}}$  programs are computationally equivalent to the original  $F_\omega$  terms. That e.g. is possible through a suitable logical relation, but would be rather involved, requiring simultaneous reasoning about the embedding just defined *and* the 1ML elaboration back into  $F_\omega$ . Given the informal “obviousness” of this result we take the liberty to punt.

Although we will not show it here either, it is also worth noting that a *predicative* version of System  $F_\omega$  can be embedded into plain  $1ML_{\text{ex}}$ . All uses of **wrap** and **unwrap** in the embedding of types and terms could simply be dropped – except in the case of existential formation, where it is replaced by sealing, i.e.,  $\llbracket \langle \tau_1, e_2 \rangle \tau \rrbracket = (\llbracket \tau_1 \rrbracket, \llbracket e_2 \rrbracket) :> \llbracket \tau \rrbracket$ .



### 5.3 Compositionality

The takeaway from the last two sections is that  $\text{IML}_{\text{ex}}$  is no more and no less powerful than System  $F_{\omega}$ . But if it is all just a matter of a fairly direct translation, why is it so much less compelling to program in  $F_{\omega}$  directly? What is the fundamental advantage that a language like IML provides, apart from just nifty type inference and a lot of syntax sugar?

The answer is *modular compositionality*. The heart of modularity is the ability to group types and values together into reusable units, a.k.a. modules and signatures. But also, that from these, larger units can be formed in a compositional manner.

Consider the following toy signatures for demonstration:

```
type A = {type t; v : t}
type B = {type u; w : u}
```

They can be freely composed into new types:

```
type S = {a : A; b : B; c : {v : a.t; w : b.u}}
type F = S → S
```

Similarly, values of these types can be constructed and composed, for example:

```
m1 :> A = {type t = int; v = 1}
m2 :> B = {type u = bool; w = true}
m3 : S = {a = m1; b = m2; c = {v = a.v; w = b.w}}
```

More interestingly, compound types can be abstracted over as they are:

```
g = fun (x : A) (y : B) (f : F) ⇒ f {a = x; b = y}
h = fun (p : S) (f : F) ⇒ (f p).a
```

With raw  $F_{\omega}$  such composition is not directly possible. It is not possible to group types and values *together* into composable and reusable units like the signatures and structures above. Type components must instead be represented by quantified types, whose form, placement and scope is dependent on the circumstances.

For example, a naive attempt of translating A and B to  $F_{\omega}$  would be to define them as the corresponding existential types, in direct adherence to Mitchell & Plotkin (1988):

$$A = \exists t. \{v : t\}$$

$$B = \exists u. \{w : u\}$$

But that's no good! How can we define S in terms of these types? How functions corresponding to g and h? It is not possible.

The closest one can come to defining reusable named signatures in  $F_{\omega}$  is by  $\lambda$ -abstracting the contained abstract types (Russo, 2003):

$$A = \lambda t. \{v : t\}$$

$$B = \lambda u. \{w : u\}$$

$$S = \lambda t u. \{a : A t, b : B u, c : \{v : t, w : u\}\}$$

That way, they can be associated with the appropriate binders as needed, for example:

$$F = \forall t u. S t u \rightarrow \exists t' u'. S t' u'$$

$$g = \lambda t u. \lambda x : A t. \lambda y : B u. \lambda f : F. f t u \{a = x, b = y\}$$

$$h = \lambda t u. \lambda p : S t u. \lambda f : F. \text{unpack } \langle t' u', q \rangle = f t u p \text{ in pack } \langle t', q.a \rangle_{\exists t'. A t'}$$

**Syntax**                    (types)  $T ::= \dots \mid \_ \mid '(X : \mathbf{type}) \Rightarrow T$

Abbreviations:

(expressions)	$\mathbf{if} E_1 \mathbf{then} E_2 \mathbf{else} E_3$	$:=$	$\mathbf{if} E_1 \mathbf{then} E_2 \mathbf{else} E_3 : \_$
	$\mathbf{fun} X \Rightarrow E$	$:=$	$\mathbf{fun} (X : \_) \Rightarrow E$
(types)	$'X \Rightarrow T$	$:=$	$'(X : \mathbf{type}) \Rightarrow T$
(declarations)	$X \overline{Y} : T$	$:=$	$X : '(Y : \mathbf{type}) \Rightarrow T$

Fig. 11. Extension to Full IML

In general, composition in  $F_\omega$  requires systematically separating all “static” components of a type (i.e., the types it contains) from its “dynamic” components (i.e., proper values), a transformation known as *phase-splitting* in the module literature (Dreyer *et al.*, 2003). And if the latter definitions look awfully like the result of IML elaboration then that is no coincidence! Behind the syntax distraction, IML elaboration is primarily a means of automatically and systematically phase-splitting the input program.

Of course, the programmer could try to do that manually, following the sketch above. However, from those simple examples it is already easy to extrapolate that this approach gets very laborious very quickly. Let alone the managing of corresponding existentials on the term level. Harper & Pierce (2005) in fact argue that modular programming based on parameterised signatures (which is what the above is, essentially) leads to a combinatorial explosion of parameters accumulating with each layer of abstraction.

That problem carries over to other proposals that suggest adding or emulating modular structure via manual management of polymorphism or separating type components. The fact that under such approaches it is simply not possible to really build fully composable program units means that such approaches are not really modular. IML avoids this problem because it does not treat modular composition as an afterthought, but instead takes it as a starting point at the centre of its design.

## 6 Full IML

A language without type inference is unworthy being named ML. Because that is so, this section extends  $\text{IML}_{\text{ex}}$  to full IML, which includes implicit typing and recovers ML-style implicit polymorphism. Figure 11 shows the minimal syntactic extension necessary.

### 6.1 Extensions

Syntactically, full IML merely adds two new forms of type expression to the kernel: *wildcards* “ $\_$ ” for types inferred from context, and *implicit* function types, distinguished by a leading tick ‘ (a choice that will become clear in a moment). Let us walk through a number of examples, and also explain some of the new syntax sugar on the way.

**Inferred Types** The wildcard “ $\_$ ” stands for a type that is to be inferred from context – similar to newer versions of OCaml. 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.

For example, with this syntax, the precise argument type of a function or conditional can be left out:

or  $(x : \_)(y : \_) = \text{if } x \text{ then true else } y : \_$

The type system will infer that all wildcards in this function are placeholders for `bool`. As sugar on top of this new piece of kernel syntax we allow a type annotation “: `_`” on a function parameter or conditional to be omitted entirely, thereby recovering the implicitly typed expression syntax familiar from ML. So the above can be condensed to just

or  $x y = \text{if } x \text{ then true else } y$

At this point we drop the  $\text{IML}_{\text{ex}}$  interpretation of an unannotated parameter as a **type**; we only keep that interpretation in those declaration or binding shorthands that are headed by the **type** keyword, that is:

**type**  $t \ a \ b = a \rightarrow b$

still defaults  $a$  and  $b$  to have type **type**. But in all other places parameter types default to “`_`” in full IML, that is, they can be any small type.

Wildcards can also be used to give partial annotations:

$\text{projPlus1 } b \ (p : \{x : \_ ; y : \_ \}) = (\text{if } b \text{ then } p.x \text{ else } p.y) + 1$

In this example, the annotation specifies the outer shape of the parameter  $p$  (so that it avoids running into inference problems with records, see below). But the component types are inferred to be `int`.

Large types cannot be inferred. Consequently, type annotations still have to be given when a parameter type is large (as, e.g., for polymorphic parameters) or when a conditional’s type is large (as in the “computed modules” example given in Section 2) – that is, in all cases that the core type system of conventional ML cannot express. However, the ability to give partial annotations can still come in handy:

$f \ (M : \{\text{type } t ; \text{inj} : \_ ; \text{proj} : \_ \}) \ (x : M.t) \ y = M.\text{proj } x + M.\text{proj } (M.\text{inj } y) + y$

Due to its type component  $t$ , the parameter  $M$  has a large type, but the monomorphic types of its other components can still be inferred as  $\text{inj} : \text{int} \rightarrow t$  and  $\text{proj} : t \rightarrow \text{int}$  (note that they involve  $t!$ ).

**Polymorphic Inference** The previous examples all where monomorphic. But of course, the claim to fame of ML’s type inference (Damas & Milner, 1982) is that it can assign *polymorphic* types to let-bindings. And of course, that ought to work in IML as well:

$\text{pick } b \ x \ y = \text{if } b \text{ then } x \text{ else } y$

Here, the instantiation of the (omitted) wildcards on the parameters  $x$  and  $y$  is arbitrary and the type of `pick` thus polymorphic. This polymorphic type can be expressed syntactically with the other new form of type expression, *implicit function types*:

$\text{pick} : '(a : \text{type}) \Rightarrow \text{bool} \Rightarrow a \Rightarrow a \Rightarrow a$

This corresponds to an ML-style polymorphic type, but makes the binder explicit. For obvious reasons, an implicit function has to be pure. 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. Because the type annotation is fixed to **type** anyway, we allow to omit it (per Figure 11) and just write the above type as:

```
pick : 'a ⇒ bool ⇒ a ⇒ a ⇒ a
```

Or even moving the implicit parameter to the left-hand side, just like we already allow for explicit parameters (cf. Figure 2):

```
pick 'a : bool ⇒ a ⇒ a ⇒ a
```

You see it now: the tick becomes a pun on ML's type variable syntax, but without relying on brittle implicit scoping rules.

As the name would suggest, there are no explicit introduction or elimination forms for implicit functions. Instead, implicit parameters are introduced and eliminated implicitly. So, as one would rightly expect from an ML, the previous function can be invoked without worrying about its type argument:

```
x = pick false 2 3
y = pick true (pick false) (pick true)
```

Because implicit parameters always have type **type**, again only small types can be inferred as arguments. So `(pick (size > threshold) HashMap TreeMap)` would *not* be well-typed, but demands an explicitly polymorphic version of the pick function.

Implicit type generalisation applies to expressions, not just bindings. That is true in some formulations of ML-style polymorphism as well, but the difference usually is not observable. In 1ML it is, because we can have contexts that require a polymorphic value:

```
makeAB (g : '(a : type) ⇒ a ⇒ a ⇒ a) = {a = g 1 2; b = g true false}
```

Given this, the argument function in the invocation

```
makeAB (fun x y ⇒ x)
```

is inferred to have polymorphic type. In fact, we could even apply `makeAB (pick true)`, and here is why.

**Purity restriction** 1ML has no need for ML's syntactic *value restriction* (Wright, 1995). Instead, it (once more) replaces a blunt syntactic restriction with a more precise semantic one. And we can (once more) reuse a concept that we already have introduced: purity. Any *pure* expression can have its type generalised. For example:

```
pickFst = pick true
```

has polymorphic type

```
pickFst : '(a : type) ⇒ a ⇒ a ⇒ a
```

because the fat arrow in (the first explicit argument of) the type of pick tells that its application has no side effects. Compare that to

```
pickImp1 b x y = if !b then x else y
pickImp2 b = if !b then fun x y ⇒ x else fun x y ⇒ y
```

Assuming  $(!) : '(a : \mathbf{type}) \Rightarrow \text{ref } a \rightarrow a$ , as the implicitly polymorphic version of the `rd` operator from Section 2, the previous functions have types as follows:

```
pickImp1 : '(a : type) ⇒ ref bool ⇒ a ⇒ a → a
pickImp2 : '(a : type) ⇒ ref bool → a ⇒ a ⇒ a
```

Respectively,  $(\text{pickImp}_1 \text{ true})$  could still be generalised, but  $(\text{pickImp}_2 \text{ true})$  could not.

**Subtyping** The definition and use of expressions is just one way to introduce or eliminate implicit functions. The other is through subtyping: an implicit function type is a subtype of any possible application of this function. If you imagine the two types:

```
type T1 = int → bool
type T2 = 'a ⇒ 'b ⇒ a → b
```

then  $T_2$  is a subtype of  $T_1$ . This is another generalisation of a notion that already exists in traditional ML: when matching signatures, a value component can be matched by a more polymorphic value component. Here, it becomes an independent notion.

**Example** With these few extensions, the `Map` functor from Section 2 can now be written in IML very much like you would have it in traditional ML:

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

## 6.2 Semantics

The semantics for the extension of  $\text{IML}_{\text{ex}}$  to Full IML can be seen in Figure 12. As is typical for type systems with inference, the formulation of the semantics is purely *declarative*: the typing rules are non-deterministic, they simply “guess” the right types and quantifiers in various places. We defer the discussion of an actual *algorithm* that implements type inference for this system to Section 7).

**Semantic Types** Wildcards do not induce any extension to semantic types, since they just stand for a suitable small type  $\sigma$ . But obviously, we have to extend our semantic types with a representation for implicit function types.

**Semantic Types** (large signatures)  $\Sigma ::= \dots \mid \forall \bar{\alpha}. \{ \} \rightarrow_A \Sigma$

**Types**

$$\frac{\Gamma \vdash \sigma : \Omega}{\Gamma \vdash \_ \rightsquigarrow \sigma} \text{TINFER} \quad \frac{\Gamma, \alpha, X : [= \alpha] \vdash T \rightsquigarrow \Sigma \quad \kappa_\alpha = \Omega}{\Gamma \vdash '(X : \mathbf{type}) \Rightarrow T \rightsquigarrow \forall \alpha. \{ \} \rightarrow_A \Sigma} \text{TIMPL}$$

$$\boxed{\Gamma \vdash T \rightsquigarrow \Xi}$$

**Expressions**

$$\frac{\Gamma, \bar{\alpha} \vdash E :_P \Sigma \rightsquigarrow e \quad \overline{\kappa_\alpha = \Omega}}{\Gamma \vdash E :_P \forall \bar{\alpha}. \{ \} \rightarrow_A \Sigma \rightsquigarrow \lambda \bar{\alpha}. \lambda_A x. \{ \}. e} \text{EGEN}$$

$$\frac{\Gamma \vdash E :_\eta \exists \bar{\alpha}. \forall \bar{\alpha}'. \{ \} \rightarrow_A \Sigma \rightsquigarrow e \quad \overline{\Gamma, \bar{\alpha} \vdash \sigma : \kappa_{\alpha'}}}{\Gamma \vdash E :_\eta \exists \bar{\alpha}. \Sigma[\bar{\sigma}/\bar{\alpha}'] \rightsquigarrow \text{unpack } \langle \bar{\alpha}, x \rangle = e \text{ in pack } \langle \bar{\alpha}, (x \bar{\sigma} \{ \}) \rangle. A} \text{EINST}$$

$$\boxed{\Gamma \vdash E :_\eta \Xi \rightsquigarrow e}$$

**Subtyping**

$$\frac{\overline{\Gamma \vdash \sigma : \kappa_{\alpha'}} \quad \Gamma \vdash \Sigma'[\bar{\sigma}/\bar{\alpha}'] \leq_{\bar{\pi}} \Sigma \rightsquigarrow \delta; f}{\Gamma \vdash \forall \bar{\alpha}'. \{ \} \rightarrow_A \Sigma' \leq_{\bar{\pi}} \Sigma \rightsquigarrow \delta; \lambda x. f((x \bar{\sigma} \{ \}) \cdot A)} \text{SIMPLL}$$

$$\frac{\Gamma, \bar{\alpha} \vdash \Sigma' \leq_{\bar{\pi}} \Sigma \rightsquigarrow \delta; f \quad \overline{\text{fv}(\delta \bar{\pi}) \not\uparrow \bar{\alpha}}}{\Gamma \vdash \Sigma' \leq_{\bar{\pi}} \forall \bar{\alpha}. \{ \} \rightarrow_A \Sigma \rightsquigarrow \delta; \lambda x. \lambda \bar{\alpha}. \lambda_A y. \{ \}. f x} \text{SIMPLR}$$

$$\boxed{\Gamma \vdash \Xi' \leq_{\bar{\pi}} \Xi \rightsquigarrow \delta; f}$$

Fig. 12. New elaboration rules for full IML

The semantic representation of the syntactic type  $'(a : \mathbf{type}) \Rightarrow \{ \dots \}$  is  $\forall \alpha. \{ \} \rightarrow_A \{ \dots \}$ . We write these types with an arrow  $\rightarrow_A$  (“A” for automatic), in order to reuse the notational encoding we already introduced for effects. It is a distinct form from  $\rightarrow_\eta$ , however, i.e., A is not an effect, and not included in  $\eta$ . All implicit functions are pure.

The semantic denotation of an explicitly written, syntactic implicit function type always has exactly one quantified type variable. However, *inferred* implicit functions can quantify multiple type variables at once. For example, the polymorphic function

$$f \times y = \{a = x; b = y\}$$

can conveniently be typed  $\forall \alpha \beta. \{ \} \rightarrow_A \alpha \rightarrow_P \beta \rightarrow_P \{a : \alpha, b : \beta\}$ . As we will see, this type is interchangeable (via subtyping in both directions) with  $\forall \alpha. \{ \} \rightarrow_A \forall \beta. \{ \} \rightarrow_A \alpha \rightarrow_P \beta \rightarrow_P \{a : \alpha, b : \beta\}$ , so this is simply a technical convenience.

The term argument to an implicit function is always  $\{ \}$ , a.k.a. unit. That is so because the function’s parameters are simply a sequence of types  $\bar{\alpha}$ , and there is little value in reifying them as terms  $\overline{[= \alpha]}$ , since they are implicit anyway.<sup>10</sup>

**Typing Rules** With regard to elaboration, we first need the formation rules for the new pieces of type syntax. The rule for wildcards simply guesses a (well-formed) small type non-deterministically. The one for implicit functions mirrors the one for regular function

<sup>10</sup> We could have dropped the term argument entirely, representing implicit function types simply as  $\forall_A \bar{\alpha}. \Sigma$ , but we opted for keeping the shape of all function types uniform and reusing the  $\rightarrow_x$  notation. In part, because we foresee nice generalisations of implicit functions.

types, but with a fixed parameter type. Moreover, implicit functions may not produce abstract types – a restriction that wouldn’t strictly be necessary, but precludes questionable examples like the following, where each use of weird could denote a completely different type:<sup>11</sup>

$f$  (weird : 'a  $\Rightarrow$  **type**) ( $x$  : weird) ( $y$  : weird) = ...

Next, the new typing rules for introducing and eliminating implicit functions, EGEN and EINST, match common formulations of ML-style polymorphism (Damas & Milner, 1982). The former rule simply introduces (non-deterministically) a sequence of fresh type variables into the context that can be used to type an expression – which must play nicely by being pure, however. In the conclusion, the rule turns these variables into implicit function parameters. Conversely, the latter rule can take any expression with implicit function type and replaces its type parameters with arbitrary (small) types. The only difference to plain Damas/Milner is that the rule must allow for abstract types  $\bar{\alpha}$  being generated by the (possibly impure) expression; as usual, their quantifiers are just hoisted.

Finally, subtyping allows the implicit introduction and elimination of implicit functions as well, via instantiation on the left, or skolemisation on the right (rules SIMPLL and SIMPLR). As mentioned earlier, this closely corresponds to ML’s signature matching rules. In combination, the two rules also implement reflexivity for implicit function types. For example,  $\forall \alpha. \{ \} \rightarrow_A [= \alpha]$  can be derived to be a subtype of itself by first (reversely) applying SIMPLR, moving the right-hand side  $\alpha$  from the quantifier to the context, and then applying SIMPLL to derive  $\forall \alpha'. \{ \} \rightarrow_A [= \alpha'] \leq [= \alpha]$  via instantiation  $[\alpha/\alpha']$ .

### 6.3 Meta-Theory Revisited

The additional semantics of 1ML does not introduce any difficulty for the soundness proof from Section 4.3. The following properties extend easily to the new rules:

*Proposition 6.1 (Correctness of Full 1ML Elaboration)*

Let  $\Gamma \vdash \square$ .

1. If  $\Gamma \vdash T/D \rightsquigarrow \Xi$ , then  $\Gamma \vdash \Xi : \Omega$ .
2. If  $\Gamma \vdash E/B :_{\eta} \Xi \rightsquigarrow e$ , then  $\Gamma \vdash e : \Xi$ , and if  $\eta = P$  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  $\Gamma \vdash \delta : \Gamma, \bar{\alpha} \vdash f : \Xi' \rightarrow \delta \Xi$ .

Like before, soundness is still a direct corollary:

*Theorem 6.2 (Soundness of Full 1ML)*

If  $\cdot \vdash E : \Xi \rightsquigarrow e$ , then either  $e \uparrow$  or  $e \hookrightarrow^* v$  such that  $\cdot \vdash v : \Xi$ .

What breaks, though, is our previous proof for decidability, because several of the new rules are “guessing”. The next section is dedicated to coping with that problem.

<sup>11</sup> Our type system still allows *definitions* like “**type** weird =  $\_$ ” to be polymorphic with type  $\forall \alpha. \{ \} \rightarrow_A [= \alpha]$ . Those could be ruled out by removing  $[= \sigma]$  from the syntax of small types. However, that is a more substantial restriction. More practical experience with the language is probably needed to figure out the most pragmatic trade-off.

## 7 Type Inference

With the additions from Figure 12 we have turned the deterministic typing and elaboration judgements of  $1ML_{ex}$  non-deterministic. They have to guess types (in rules  $TINFER$ ,  $EINST$ ,  $SIMPLL$ ) and quantifiers (in rule  $EGEN$ ). Moreover, they have to decide when to apply rules  $EGEN$  and  $EINST$ . Clearly, an algorithm is needed.

Fortunately, what’s going on is not fundamentally different from conventional (core) ML. Where core ML would require type equivalence (and type inference would use unification), the 1ML rules require subtyping.

Subtyping sneaking into inference may seem scary at first. But closer inspection of the subtyping rules reveals that, when applied to small types, they almost degenerate to type equivalence! The only exception is width subtyping on records, which we will get back to in Section 7.2. The 1ML 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* – holding place for small types – and a variation of unification.

### 7.1 Algorithm

Figures 13-15 show how the elaboration rules from the previous sections (including the rules for impredicative types) can be turned into an inference algorithm. Because we are just interested in typing, witness *terms* (i.e., the greyed out bits) are omitted from the rules this time.

The basic idea of these algorithmic rules is to introduce a (free) inference variable  $v$  wherever the original declarative rules have to guess a (small) type – for example, in rule  $ITINFER$ . Furthermore, the rules are augmented by outputting a substitution  $\theta$  for *resolved* inference variables: all judgements have the form  $\Gamma \vdash_{\theta} J$ , which, roughly, implies the respective  $\bar{v}, \theta \Gamma \vdash \theta J$  declaratively, where  $\bar{v}$  in the context binds the unresolved inference variables that still appear free in  $\theta \Gamma$  or  $\theta J$  (as if they were regular type variables).

In rules with multiple premises, the output substitution of one premise needs to be applied to the inputs of the next, and all output substitutions must be composed. We avoid too much notational awkwardness with abbreviations of the form

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

for all inference judgement forms  $J$ , where  $\theta J$  is meant to apply  $\theta$  to  $J$ ’s “inputs”, i.e., everything that is not right of a colon “:”, squiggly arrow “ $\rightsquigarrow$ ”, or instantiation operator “ $\preceq$ ” (in the case of the auxiliary Instantiation judgement, see below). This notation is used to consistently thread and compose substitutions through all rules that have multiple premises.

Other than the systematic threading of  $\theta$  there are fairly few changes relative to the declarative rules. In Figures 13-15, all those other differences are **highlighted** for easier reference.

The basics of the inference algorithm are not that unusual, focussing around generalisation, instantiation, and resolution of inference variables.

**Generalisation** The declarative typing rules allow generalising the type of an expression at any point. However, there are only two situations under which generalisation actually



## Types

$$\boxed{\Gamma \vdash_{\theta} T \rightsquigarrow \Xi}$$

$$\frac{\Gamma \vdash_{\theta}^! E :_{\mathbb{P}} [= \Xi]}{\Gamma \vdash_{\theta} E \rightsquigarrow \Xi} \text{ITPATH} \quad \frac{v \text{ fresh} \quad \Delta_v = \text{dom}(\Gamma)}{\Gamma \vdash_{\square} - \rightsquigarrow v} \text{ITINFER}$$

$$\frac{}{\Gamma \vdash_{\square} \mathbf{type} \rightsquigarrow \exists \alpha. [= \alpha]} \text{ITTYPE} \quad \frac{}{\Gamma \vdash_{\square} \mathbf{bool} \rightsquigarrow \text{bool}} \text{ITBOOL}$$

$$\frac{\Gamma \vdash_{\theta} D \rightsquigarrow \Xi}{\Gamma \vdash_{\theta} \{D\} \rightsquigarrow \Xi} \text{ITSTR} \quad \frac{\Gamma \vdash_{\theta_1} T_1 \rightsquigarrow \exists \bar{\alpha}_1. \Sigma_1 \quad \Gamma; \bar{\alpha}_1, X : \Sigma_1 \vdash_{\theta_1} \theta_2 T_2 \rightsquigarrow \exists \bar{\alpha}_2. \Sigma_2}{\Gamma \vdash_{\theta_2} (X : T_1) \rightarrow T_2 \rightsquigarrow \forall \bar{\alpha}_1. \Sigma_1 \rightarrow_1 \exists \bar{\alpha}_2. \Sigma_2} \text{ITFUN}$$

$$\frac{\Gamma \vdash_{\theta_1} T_1 \rightsquigarrow \exists \bar{\alpha}_1. \Sigma_1 \quad \Gamma; \bar{\alpha}_1, X : \Sigma_1 \vdash_{\theta_1} \theta_2 T_2 \rightsquigarrow \exists \bar{\alpha}_2. \Sigma_2 \quad \overline{\kappa_{\alpha'_2} = \bar{\kappa}_{\alpha_1} \rightarrow \kappa_{\alpha_2}}}{\Gamma \vdash_{\theta_2} (X : T_1) \Rightarrow T_2 \rightsquigarrow \exists \bar{\alpha}'_2. \forall \bar{\alpha}_1. \Sigma_1 \rightarrow_{\mathbb{P}} \Sigma_2[\bar{\alpha}'_2 \bar{\alpha}_1 / \bar{\alpha}_2]} \text{ITPFUN}$$

$$\frac{\Gamma; \alpha, X : [= \alpha] \vdash_{\theta} T \rightsquigarrow \Sigma \quad \kappa_{\alpha} = \Omega}{\Gamma \vdash_{\theta} (X : \mathbf{type}) \Rightarrow T \rightsquigarrow \forall \alpha. \{\} \rightarrow_{\mathbb{A}} \Sigma} \text{ITIMPL}$$

$$\frac{\Gamma \vdash_{\theta} T \rightsquigarrow \Xi}{\Gamma \vdash_{\theta} \mathbf{wrap} T \rightsquigarrow [\Xi]} \text{ITWRAP} \quad \frac{\Gamma \vdash_{\theta} E :_{\mathbb{P}} \Sigma \quad E = E' :> T \vee E = E' : T \vee E = \mathbf{type} T}{\Gamma \vdash_{\theta} (= E) \rightsquigarrow \Sigma} \text{ITSING}$$

$$\frac{\Gamma \vdash_{\theta_1} T_1 \rightsquigarrow \exists \bar{\alpha}_1. \Sigma_1 \quad \Gamma \vdash_{\theta_1} \theta_2 T_2 \rightsquigarrow \exists \bar{\alpha}_2. \Sigma_2 \quad \Gamma; \bar{\alpha}_{11}, \bar{\alpha}_2 \vdash_{\theta} \Sigma_2 \leq_{\bar{\alpha}_{12}} \Sigma_1. \bar{X} \rightsquigarrow \delta \quad \bar{\alpha}_1 = \bar{\alpha}_{11} \uplus \bar{\alpha}_{12}}{\Gamma \vdash_{\theta} T_1 \mathbf{where} (\bar{X} : T_2) \rightsquigarrow \exists \bar{\alpha}_{11} \bar{\alpha}_2. \delta \Sigma'_1[\bar{X} = \Sigma_2]} \text{ITWHERE}$$

## Declarations

$$\boxed{\Gamma \vdash_{\theta} D \rightsquigarrow \Xi}$$

$$\frac{\Gamma \vdash_{\theta} T \rightsquigarrow \exists \bar{\alpha}. \Sigma}{\Gamma \vdash_{\theta} X : T \rightsquigarrow \exists \bar{\alpha}. \{X : \Sigma\}} \text{IDVAR} \quad \frac{\Gamma \vdash_{\theta} T \rightsquigarrow \exists \bar{\alpha}. \{\bar{X} : \bar{\Sigma}\}}{\Gamma \vdash_{\theta} \mathbf{include} T \rightsquigarrow \exists \bar{\alpha}. \{\bar{X} : \bar{\Sigma}\}} \text{IDINCL}$$

$$\frac{\Gamma \vdash_{\theta_1} D_1 \rightsquigarrow \exists \bar{\alpha}_1. \{\bar{X}_1 : \bar{\Sigma}_1\} \quad \Gamma; \bar{\alpha}_1, \bar{X}_1 : \bar{\Sigma}_1 \vdash_{\theta_1} \theta_2 D_2 \rightsquigarrow \exists \bar{\alpha}_2. \{\bar{X}_2 : \bar{\Sigma}_2\} \quad \bar{X}_1 \cap \bar{X}_2 = \emptyset}{\Gamma \vdash_{\theta_2} D_1 ; D_2 \rightsquigarrow \exists \bar{\alpha}_1 \bar{\alpha}_2. \{\bar{X}_1 : \bar{\Sigma}_1, \bar{X}_2 : \bar{\Sigma}_2\}} \text{IDSEQ} \quad \frac{}{\Gamma \vdash_{\square} \varepsilon \rightsquigarrow \{\}} \text{IDEMPTY}$$

Fig. 13. Inference for Types

needs to happen: (1) when a value is bound to a variable, because then it could potentially be used at multiple different types, and (2) when an expression's type needs to match an explicit annotation (directly or indirectly), because that type may demand polymorphism.

The first case corresponds to the classical let-rule of Damas/Milner, and occurs in the rule IBVAR: this invokes the auxiliary Generalisation judgement (Figure 16). For pure expressions, it finds all inference variables  $\bar{v}$  in the expression's type that do not also occur in the context (we write  $\text{undet}(A)$  to denote the free inference variables of a semantic object  $A$ ) and quantifies over them – in our case, by introducing an implicit function abstraction.

**Expressions**

$$\boxed{\Gamma \vdash_{\theta} E :_{\eta} \Xi}$$

$$\frac{\Gamma(X) = \Sigma}{\Gamma \vdash_{\emptyset} X :_{\mathbb{P}} \Sigma} \text{IEVAR} \quad \frac{\Gamma \vdash_{\theta} T \rightsquigarrow \Xi}{\Gamma \vdash_{\theta} \text{type } T :_{\mathbb{P}} [\Xi]} \text{IETYPE}$$

$$\frac{}{\Gamma \vdash_{\emptyset} \text{true} :_{\mathbb{P}} \text{bool}} \text{IETRUE} \quad \frac{}{\Gamma \vdash_{\emptyset} \text{false} :_{\mathbb{P}} \text{bool}} \text{IEFALSE}$$

$$\frac{\Gamma \vdash_{\theta_0}^! X :_{\mathbb{P}} \text{bool} \quad \Gamma \theta_0 \vdash_{\theta_1}^{\vee} E_1 :_{\eta_1} \Xi_1 \quad \Gamma \theta_3 \vdash_{\theta_4} \Xi_1 \leq \Xi}{\Gamma \theta_2 \vdash_{\theta_3} T \rightsquigarrow \Xi \quad \Gamma \theta_1 \vdash_{\theta_2}^{\vee} E_2 :_{\eta_2} \Xi_2 \quad \Gamma \theta_4 \vdash_{\theta_5} \Xi_2 \leq \Xi} \text{IEIF}$$

$$\frac{}{\Gamma \vdash_{\theta_5} \text{if } X \text{ then } E_1 \text{ else } E_2 :_{\eta_1 \vee \eta_2 \vee \eta(\Xi)} \Xi} \text{IEIF}$$

$$\frac{\Gamma \vdash_{\theta} B :_{\eta} \Xi}{\Gamma \vdash_{\theta} \{B\} :_{\eta} \Xi} \text{IESTR} \quad \frac{\Gamma \vdash_{\theta}^! E :_{\eta} \exists \bar{\alpha}. \{X : \Sigma, \bar{X}' : \bar{\Sigma}'\}}{\Gamma \vdash_{\theta} E.X :_{\eta} \exists \bar{\alpha}. \Sigma} \text{IEDOT}$$

$$\frac{\Gamma \vdash_{\theta_1} T \rightsquigarrow \exists \bar{\alpha}. \Sigma \quad \Gamma; \bar{\alpha}, X : \Sigma \theta_1 \vdash_{\theta_2} E :_{\eta} \Xi}{\Gamma \vdash_{\theta_2} \text{fun } (X : T) \Rightarrow E :_{\mathbb{P}} \forall \bar{\alpha}. \Sigma \rightarrow_{\eta} \Xi} \text{IEFUN}$$

$$\frac{\Gamma \vdash_{\theta_1}^! X_1 :_{\mathbb{P}} \forall \bar{\alpha}. \Sigma_1 \rightarrow_{\eta} \Xi \quad \Gamma \theta_2 \vdash_{\theta_3} \Sigma_2 \leq_{\bar{\alpha}} \Sigma_1 \rightsquigarrow \delta}{\Gamma \vdash_{\theta_3} X_1 X_2 :_{\eta} \delta \Xi} \text{IEAPP}$$

$$\frac{\Gamma \vdash_{\theta_1} X :_{\mathbb{P}} \Sigma_1 \quad \Gamma \theta_1 \vdash_{\theta_2} T \rightsquigarrow \exists \bar{\alpha}. \Sigma_2 \quad \Gamma \theta_2 \vdash_{\theta_3} \Sigma_1 \leq_{\bar{\alpha}} \Sigma_2 \rightsquigarrow \delta}{\Gamma \vdash_{\theta_3} X :_{>T} :_{\eta(\exists \bar{\alpha}. \Sigma_2)} \exists \bar{\alpha}. \Sigma_2} \text{IESEAL}$$

$$\frac{\Gamma \vdash_{\theta_1} X :_{\mathbb{P}} \Sigma \quad \Gamma \theta_1 \vdash_{\theta_2} T \rightsquigarrow [\Xi] \quad \Gamma \theta_2 \vdash_{\theta_3} \Sigma \leq \Xi}{\Gamma \vdash_{\theta_3} \text{wrap } X : T :_{\mathbb{P}} [\Xi]} \text{IEWRAP} \quad \frac{\Gamma \vdash_{\theta_1}^! X :_{\mathbb{P}} [\Xi'] \quad \Gamma \theta_1 \vdash_{\theta_2} T \rightsquigarrow [\Xi] \quad \Gamma \theta_2 \vdash_{\theta_3} \Xi' \leq \Xi}{\Gamma \vdash_{\theta_3} \text{unwrap } X : T :_{\mathbb{P}} \Xi} \text{IEUNWRAP}$$

**Bindings**

$$\boxed{\Gamma \vdash_{\theta} B :_{\eta} \Xi}$$

$$\frac{\Gamma \vdash_{\theta}^{\vee} E :_{\eta} \exists \bar{\alpha}. \Sigma}{\Gamma \vdash_{\theta} X = E :_{\eta} \exists \bar{\alpha}. \{X : \Sigma\}} \text{IBVAR} \quad \frac{\Gamma \vdash_{\theta}^! E :_{\eta} \exists \bar{\alpha}. \{X : \Sigma\}}{\Gamma \vdash_{\theta} \text{include } E :_{\eta} \exists \bar{\alpha}. \{X : \Sigma\}} \text{IBINCL}$$

$$\frac{\Gamma \vdash_{\theta_1} B_1 :_{\eta_1} \exists \bar{\alpha}_1. \{X_1 : \Sigma_1\} \quad \bar{X}'_1 = \bar{X}_1 - \bar{X}_2 \quad \Gamma; \bar{\alpha}_1, \bar{X}_1 : \Sigma_1 \theta_1 \vdash_{\theta_2} B_2 :_{\eta_2} \exists \bar{\alpha}_2. \{X_2 : \Sigma_2\} \quad \bar{X}'_1 : \Sigma'_1 \subseteq \bar{X}_1 : \Sigma_1}{\Gamma \vdash_{\theta_2} B_1 ; B_2 :_{\eta_1 \vee \eta_2} \exists \bar{\alpha}_1 \bar{\alpha}_2. \{X'_1 : \Sigma'_1, X_2 : \Sigma_2\}} \text{IBSEQ} \quad \frac{}{\Gamma \vdash_{\emptyset} \varepsilon :_{\mathbb{P}} \{\}} \text{IEEMPTY}$$

Fig. 14. Inference for Expressions

In fact, this rule already covers almost all cases of the second kind as well, because our kernel grammar is a sort of normal form which requires named variables in most places (Figure 1). In particular, this is the case for function arguments. Consider this example:

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

Desugaring rewrites this application into an expression that has an explicit binding for the argument (**fun** x  $\Rightarrow$  x), such that rule IBVAR kicks in. The same observations applies to other places where an expression's type must match an annotation: rules IEsEal, IEwraP

$$\begin{array}{c}
\text{Subtyping} \quad \Gamma \vdash_{\theta} \Xi \leq \Xi' \quad := \quad \Gamma \vdash_{\theta} \Xi \leq_{\varepsilon} \Xi' \rightsquigarrow \text{id} \quad \boxed{\Gamma \vdash_{\theta} \Xi' \leq_{\bar{\pi}} \Xi \rightsquigarrow \delta} \\
\\
\frac{}{\Gamma \vdash_{\square} v \leq v} \text{ISREFL} \\
\\
\frac{\Gamma \vdash_{\theta}^! v \approx \Sigma \quad \Gamma \theta \vdash_{\theta'} v \leq \Sigma}{\Gamma \vdash_{\theta'} v \leq \Sigma} \text{ISRESL} \quad \frac{\Gamma \vdash_{\theta}^! v \approx \Sigma' \quad \Gamma \theta \vdash_{\theta'} \Sigma' \leq v}{\Gamma \vdash_{\theta'} \Sigma' \leq v} \text{ISRESR} \\
\\
\frac{\Gamma \vdash_{\theta} \pi' = \pi}{\Gamma \vdash_{\theta} \pi' \leq \pi} \text{ISPATH} \quad \frac{}{\Gamma \vdash_{\square} \text{bool} \leq \text{bool}} \text{ISBOOL} \\
\\
\frac{\Gamma \vdash_{\theta} \Xi' \leq \Xi \quad \Gamma \theta \vdash_{\theta'} \Xi \leq \Xi'}{\Gamma \vdash_{\theta'} [= \Xi'] \leq [= \Xi] \rightsquigarrow \square} \text{ISTYPE} \quad \frac{}{\Gamma \vdash_{\square} [= \sigma] \leq_{\alpha_0 \bar{\alpha}} [= \alpha_0 \bar{\alpha}] \rightsquigarrow [\lambda \bar{\alpha}. \sigma / \alpha_0]} \text{ISFORGET} \\
\\
\frac{\Gamma \vdash_{\theta} \Xi' = \Xi}{\Gamma \vdash_{\theta} [\Xi'] \leq [\Xi]} \text{ISWRAP} \quad \frac{}{\Gamma \vdash_{\square} \{l: \Sigma'\} \leq \{\}} \text{ISEMPTY} \\
\\
\frac{\Gamma \vdash_{\theta_1} \Sigma'_1 \leq_{\bar{\pi}_1} \Sigma_1 \rightsquigarrow \delta_1 \quad \Gamma \theta_1 \vdash_{\theta_2} \{l': \Sigma'\} \leq_{\bar{\pi}_2} \{l: \delta_1 \Sigma\} \rightsquigarrow \delta_2 \quad \theta_2 \delta_2 \Sigma_1 = \theta_2 \Sigma_1}{\Gamma \vdash_{\theta_2} \{l_1: \Sigma'_1, l': \Sigma'\} \leq_{\bar{\pi}_1 \bar{\pi}_2} \{l_1: \Sigma_1, l: \Sigma\} \rightsquigarrow \delta_1 \delta_2} \text{ISSTR} \\
\\
\frac{\Gamma; \bar{\alpha} \vdash_{\theta_1} \Sigma \leq_{\bar{\alpha}'} \Sigma' \rightsquigarrow \delta_1 \quad \eta' \leq \eta \quad \Gamma; \bar{\alpha} \theta_1 \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_{\eta'} \Xi') \leq_{\bar{\pi}} (\forall \bar{\alpha}. \Sigma \rightarrow_{\eta} \Xi) \rightsquigarrow \delta_2} \text{ISFUN} \\
\\
\frac{\Gamma; \bar{\alpha} \vdash_{\theta} \Sigma' \leq_{\bar{\alpha}} \Sigma \rightsquigarrow \delta \quad \bar{\alpha}' \bar{\alpha} \neq \varepsilon}{\Gamma \vdash_{\theta} \exists \bar{\alpha}'. \Sigma' \leq \exists \bar{\alpha}. \Sigma} \text{ISABS} \\
\\
\frac{\overline{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{SIMPLL} \\
\\
\frac{\Gamma; \bar{\alpha} \vdash_{\theta} \Sigma' \leq_{\bar{\pi}} \Sigma \rightsquigarrow \delta \quad \overline{\text{fv}(\theta \delta \pi) \not\uparrow \bar{\alpha}}}{\Gamma \vdash_{\theta} \Sigma' \leq_{\bar{\pi}} \forall \bar{\alpha}. \{ \} \rightarrow_{\mathbb{A}} \Sigma \rightsquigarrow \delta} \text{SIMPLR}
\end{array}$$

Fig. 15. Inference for Subtyping

and  $\text{IEunwrap}$ . Similarly, it extends to sugar such as transparent type ascriptions  $E:T$  as defined in Figure 2.

The only exception are conditionals, because their arms are not variables, yet the type annotation may require them to be polymorphic. Rule  $\text{IEIF}$  hence also has to invoke the Generalisation judgement.

**Instantiation** Instantiation takes place in a slightly different location than with Damas/Milner. Since implicit functions are first-class in IML, it is not just variables that can have “polymorphic” type. For example, the result of a projection  $\text{m.f}$  or of an application  $\text{g} \times \text{m}$  might have a polymorphic type, too, that potentially requires instantiation.

**Generalisation**

$$\boxed{\Gamma \vdash_{\theta}^{\forall} E :_{\eta} \Xi}$$

$$\frac{\Gamma \vdash_{\theta} E :_{\mathbb{I}} \Xi \quad \text{IGIMPURE} \quad \Gamma \vdash_{\theta} E :_{\mathbb{P}} \Sigma \quad \bar{v} = \text{undet}(\theta\Sigma) - \text{undet}(\theta\Gamma) \quad \overline{\kappa\alpha = \Omega}}{\Gamma \vdash_{\theta}^{\forall} E :_{\mathbb{I}} \Xi \quad \Gamma \vdash_{\theta}^{\forall} E :_{\mathbb{P}} \forall\bar{\alpha}.\{\} \rightarrow_{\mathbb{A}} \Sigma[\bar{\alpha}/\bar{v}] \quad \text{IGPURE}}$$

**Instantiation**

$$\boxed{\Gamma \vdash_{\theta}^{\exists} E :_{\eta} \Xi}$$

$$\frac{\Gamma \vdash_{\theta} E :_{\eta} \Xi \quad \text{INREFL} \quad \Gamma \vdash_{\theta}^{\exists} E :_{\eta} \exists\bar{\alpha}'.\forall\bar{\alpha}.\{\} \rightarrow_{\mathbb{A}} \Sigma \quad \bar{v} \text{ fresh} \quad \overline{\Delta_v = \text{dom}(\Gamma)}}{\Gamma \vdash_{\theta}^{\exists} E :_{\eta} \Xi \quad \Gamma \vdash_{\theta}^{\exists} E :_{\eta} \exists\bar{\alpha}'.\Sigma[\bar{v}/\bar{\alpha}'] \quad \text{INIMPL}}$$

$$\frac{\Gamma \vdash_{\theta} E :_{\eta} \exists\bar{\alpha}.v \quad \Gamma \vdash_{\theta'} v \approx \theta'v}{\Gamma \vdash_{\theta'} E :_{\eta} \exists\bar{\alpha}.v} \text{INRES}$$

**Resolution**

$$\Gamma \vdash_{\theta}^{\exists} v \approx \Sigma \quad := \quad v \notin \text{undet}(\Sigma) \wedge \Gamma \vdash_{\theta} v \approx \Sigma$$

$$\boxed{\Gamma \vdash_{\theta} v \approx \Sigma}$$

$$\frac{v'' \text{ fresh} \quad \Delta_{v''} = \Delta_v \cap \Delta_{v'}}{\Gamma \vdash_{[v''/v, v''/v']} v \approx v'} \text{IRINFER} \quad \frac{\alpha \in \Delta_v \quad \bar{v}' \text{ fresh} \quad \overline{\Delta_{v'} = \Delta_v}}{\Gamma \vdash_{[\alpha\bar{v}'/v]} v \approx \alpha\bar{\sigma}} \text{IRPATH}$$

$$\frac{}{\Gamma \vdash_{[\text{bool}/v]} v \approx \text{bool}} \text{IRBOOL} \quad \frac{v' \text{ fresh} \quad \Delta_{v'} = \Delta_v}{\Gamma \vdash_{[v'/v]} v \approx [v']} \text{IRTYPE}$$

$$\frac{v' \text{ fresh} \quad \Delta_{v'} = \Delta_v}{\Gamma \vdash_{[v'/v]} v \approx [v']} \text{IRWRAP} \quad \frac{v_1, v_2 \text{ fresh} \quad \Delta_{v_1} = \Delta_{v_2} = \Delta_v}{\Gamma \vdash_{[(v_1 \rightarrow_{\mathbb{I}} v_2)/v]} v \approx \forall\bar{\alpha}.\Sigma \rightarrow_{\eta} \Xi} \text{IRFUN}$$

**Unification**

$$\boxed{\Gamma \vdash_{\theta} \Xi = \Xi'}$$

$$\frac{}{\Gamma \vdash_{\emptyset} v = v} \text{IUREFL} \quad \frac{\Gamma \vdash_{\theta} \Xi' = \Xi}{\Gamma \vdash_{\theta} \Xi = \Xi'} \text{IUSYMM}$$

$$\frac{\bar{v}' = \text{undet}(\sigma) \quad v \notin \bar{v}' \quad \Delta_v \supseteq \text{fv}(\sigma) \quad \bar{v}'' \text{ fresh} \quad \overline{\Delta_{v''} = \Delta_{v'} \cap \Delta_v}}{\Gamma \vdash_{[\bar{v}''/\bar{v}'] \circ [\sigma/v]} v = \sigma} \text{IUBIND}$$

$$\frac{\Gamma \vdash_{\theta} \Xi = \Xi'}{\Gamma \vdash_{\theta'} [v] = [v']} \text{IUTYPE} \quad \frac{\Gamma \vdash_{\theta} \Xi = \Xi'}{\Gamma \vdash_{\theta} [v] = [v']} \text{IUWRAP} \quad \frac{\Gamma \vdash_{\theta} \bar{\sigma} = \bar{\sigma}'}{\Gamma \vdash_{\theta} \alpha\bar{\sigma} = \alpha\bar{\sigma}'} \text{IUPATH}$$

$$\frac{\Gamma \vdash_{\theta} \bar{\Sigma} = \bar{\Sigma}'}{\Gamma \vdash_{\theta} \{\bar{l}:\bar{\Sigma}\} = \{\bar{l}:\bar{\Sigma}'\}} \text{IUSTR} \quad \frac{\Gamma; \bar{\alpha} \vdash_{\theta_1} \Sigma = \Sigma' \quad \Gamma; \bar{\alpha} \vdash_{\theta_2} \Xi = \Xi'}{\Gamma \vdash_{\theta_2} \forall\bar{\alpha}.\Sigma \rightarrow_{\eta} \Xi = \forall\bar{\alpha}.\Sigma' \rightarrow_{\eta} \Xi'} \text{IUFUN}$$

$$\frac{\Gamma; \bar{\alpha} \vdash_{\theta} \Sigma = \Sigma'}{\Gamma \vdash_{\theta} \exists\bar{\alpha}.\Sigma = \exists\bar{\alpha}.\Sigma'} \text{IUABS} \quad \frac{\Gamma; \bar{\alpha} \vdash_{\theta} \Sigma = \Sigma'}{\Gamma \vdash_{\theta} \forall\bar{\alpha}.\{\} \rightarrow_{\mathbb{A}} \Sigma = \forall\bar{\alpha}.\{\} \rightarrow_{\mathbb{A}} \Sigma'} \text{IUIPML}$$

Fig. 16. Generalisation, Instantiation, Resolution, and Unification

For that reason, the rules delay instantiation as much as possible, and only perform it in the typing rules of *elimination forms* (e.g. rules IEIF, IEDOT, IEAPP, but also ITPATH). They do so with the help of an auxiliary *instantiation* judgement, shown in Figure 16.

The instantiation judgement  $\Xi' \preceq \Xi$ , may seem non-deterministic because of rule INRES, which allows fabricating a seemingly random type. However, all invoking rules require an output type of a certain shape, so that it is always enforced by context that all implicit functions are actually eliminated, i.e., rule INIMPL is always applied exhaustively, and a specific resolution rule is determined. In other words, the “expected” outermost constructor of the right-hand  $\Xi$  should conceptually be thought of as an input to the rule, too.

Instantiation can also happen implicitly as part of subtyping (rule ISIMPLL), which covers the case where a polymorphic value is matched against a monomorphic (or other polymorphic) parameter. For example,  $\forall \alpha_1 \alpha_2. \{ \} \rightarrow_A \alpha_1 \rightarrow_I \alpha_2 \leq \forall \beta. \{ \} \rightarrow_A \beta \rightarrow_I \beta$  will be checked by first applying ISIMPLR, turning the right type monomorphic, and then instantiating the left with ISIMPLL, so that the check is down to  $v_1 \rightarrow_I v_2 \leq \beta \rightarrow_I \beta$ , resolving easily.

**Resolution** There are two main complications to the previous scheme, both due to the fact that, unlike good old ML, our IML allows small types to be intermixed with large ones.

For one, it is necessary to be able to infer small types from large ones 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  $\sigma$ . Through contravariance, similar situations can arise with an inference variable on the left, e.g.:

$$v \leq (\forall \alpha. [= \alpha] \rightarrow_P [= \alpha]) \rightarrow_I \text{bool}$$

which can be solved with  $v = ([= \sigma] \rightarrow_I [= \sigma]) \rightarrow_I \text{bool}$ .

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  $\Sigma$  may or may not be small) it invokes the auxiliary *Resolution* judgement  $\Gamma \vdash_{\theta} v \approx \Sigma$  (Figure 16), which only resolves  $v$  so far as to match the outermost constructor of  $\Sigma$  and inserts fresh inference variables for its subcomponents – like with the Instantiation judgement, the “expected” constructor of the right-hand  $\Sigma$  should be thought of as an input to the resolution algorithm. After refining  $v$  that way, subtyping “tries again” in the second premise of ISRESL or ISRESR, and is guaranteed to proceed at least one step further. The occurs check encoded in the shorthand  $\Gamma \vdash_{\theta}^! v \approx \Sigma$  prevents this process from looping.

Resolution is actually invoked in two situations: during subtyping, in the rules ISRESL and ISRESR, which deal with subtyping constraints as just described, but also in rule INRES of the instantiation judgement that we have glossed over earlier. The latter rule comes into play when an elimination form is type-checked but the type of the eliminated expression is merely an unresolved inference variable. In that case, delegating to the resolution judgement will also refine the shape of the yet unknown type one step.

**Variable Scoping** The second complication to IML inference is that an inference variable  $v$  can be introduced within the scope of arbitrary abstract types (i.e., regular type variables

$\alpha$ ). It would be incorrect to resolve  $v$  to a type containing type variables that are *not* in scope for *all* occurrences of  $v$  in a derivation.

To prevent that, each  $v$  is associated with a set  $\Delta_v$  of type variables that are known to be in scope for all uses of  $v$ . The set is verified when resolving  $v$  (see resolution rule IRPATH in particular). The set also is propagated to any other  $v'$  that the original  $v$  is unified with, by intersecting  $\Delta_{v'}$  with  $\Delta_v$  – 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  $\Delta_v$  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 regular type variables go out of scope. That is achieved by employing the following notation (note the semicolon!) in all rules that locally extend  $\Gamma$  with type variables:

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

where  $\bar{v} = \text{undet}(\theta'' J)$  and  $\bar{v}'$  fresh with  $\overline{\Delta_{v'}} = \overline{\Delta_v \cap \text{dom}(\Gamma)}$

The net effect is that all local  $\alpha$ 's from  $\Gamma'$  are stripped from all  $\Delta$ -sets of inference variable that remain after executing  $\Gamma, \Gamma' \vdash J$ , thereby disallowing any resolution with a type that would contain those variables. We omit  $\theta$  in this notation when it is the identity.

**Unification** Finally, our algorithm also includes a separate, old-fashioned unification algorithm (Figure 16). However, unification is only needed by two rules: the rule ISPATH for paths and the rule ISWRAP for the wrapped types of the impredicative extension – because those are the only rule where subtyping does not recurse into subcomponents. In their declarative formulation (Figure 12), the respective rules simply require the types on both sides of the relation to be syntactically equivalent (modulo  $\alpha\beta$ -conversion, as noted). The presence of inference variables makes that formulation insufficient, of course, and necessitates traversing the structure of the two types to unify possible inference variables embedded in them.

## 7.2 Incompleteness

Unfortunately, there are a couple of sources of incompleteness in this algorithm:

**Width subtyping** Subtyping constraints like  $v \leq \{\bar{l} : \bar{\sigma}\}$  do not determine the shape of the record type  $v$ : the set of its labels can still vary over all possible subsets of  $\bar{l}$ . Consequently, the Resolution judgement in Figure 16 has no rule for structures – instead a structure type must be determined by the previous context.

This is, in fact, similar to Standard ML (Milner *et al.*, 1997), where record types cannot be inferred either, and require type annotations. 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 1ML, 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 in this article and leave out this possible improvement.

The real solution of course would be to incorporate not just row inference but *row polymorphism* (Rémy, 1989), so that width subtyping on structures can be recast as universal and existential quantification over row variables. We leave investigating such an extension for future work. We only note here that it would not easily solve the problem completely, because of **include**: that extends the current scope with bindings that are determined by a structure type; to handle the case where the structure type is row-polymorphic we somehow would have to allow *environments*  $\Gamma$  to be row-polymorphic, too, or restrict **include** in a suitable manner.

**Type Scoping** Tracking of the sets  $\Delta_v$  is conservative: after leaving the scope of a type variable  $\alpha$ , we exclude any solution for  $v$  that would still involve  $\alpha$ , even if  $v$  only appears inside a type binder for  $\alpha$ . Consider the following example, adapted from Example (c) in Dreyer & Blume (2007):

```
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);
```

Assume  $\text{id} : 'a \Rightarrow a \rightarrow a$ . Because  $\text{id}$  is impure, the definition of  $f$  is impure, and its type cannot be generalised; consequently,  $G$  is impure, too, and hence generative. The algorithm will infer  $G$ 's type as

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

with  $\beta \notin \Delta_v$  – because  $\beta$  goes out of scope the moment we bind it with a local quantifier, while all  $v$  are treated as “global”. It then generalises to

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

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

Although not well-known, this very problem is already present in good old ML, as Dreyer & Blume (2007) pointed out: all existing type inference implementations for ML 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 (as for  $G$  above). However, generalisation happens outside the abstraction, which is more rigid than what they propose (but which is not expressible in System  $F_{\omega}$ ). Consequently, IML can type some examples from their paper – such as Example (b):

```
F (t : type) = {f = id id};
A = F int; B = H bool;
x = A.f 10; y = B.f "dude";
```

This works because  $F$  is polymorphic in IML (with type  $\forall \alpha. \{ \} \rightarrow_{\text{A}} \forall \beta. [= \beta] \rightarrow_{\text{P}} \{f : \alpha \rightarrow_{\text{I}} \alpha\}$ ), but that is not enough to handle other examples, like  $G$  above.

In other words, our algorithm can infer the types of some programs that regular ML inference unexpectedly can't handle. At the same time, the liberal mixture of modules and core programming in IML introduces more of these cases, and some remain incomplete.

**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 I. But arguably, that is cheating a little by side-stepping the issue. In particular, it makes an extension of the notion of (im)purity to other effects, as suggested in Section 2, somewhat inconvenient, because pure function types could not be inferred in parameter positions.

Again, the solution would be more polymorphism, in this case a simple form of effect polymorphism (Talpin & Jouvelot, 1992): if effects could be quantified over, they could also be inferred in a complete manner. We explored this solution in a separate paper that extends IML with effect polymorphism (Rossberg, 2016), which – interestingly – induces a novel notion of *generativity polymorphism* for functors.

Despite these limitations, we found IML inference quite usable. In practice, MLs have long given up on complete type inference: various limitations exist in both SML and OCaml (and the extended language family including Haskell), necessitating type annotations or declarations. In our limited experience with a prototype, IML is not substantially worse, at least not when used in the same manner as traditional ML.

In fact, we conjecture that any SML program not using features omitted from IML – but including both modules and Damas/Milner polymorphism – can be directly transliterated into valid IML without adding a single type annotation.

### 7.3 Correctness

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

*Theorem 7.1 (Soundness of IML Inference)*

Let  $\bar{v}, \Gamma$  be a well-formed  $F_\omega$  environment.

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

In order to prove this result, we need to refine the statement a little bit. In particular, instead of just bundling together  $\bar{v}, \Gamma$  as an environment, we want to construct an environment that respects the scoping constraints on the free inference variables  $\bar{v}$  relative to the regular type variables in their  $\Delta$ -sets.

To that end, we define a little auxiliary judgement  $\bar{v} \vdash \Gamma \rightsquigarrow \Gamma'$  that intersperses  $\Gamma$  with  $\bar{v}$  such that for each  $v$  the type variables  $\Delta_v$  are in scope:

$$\frac{}{\varepsilon \vdash \cdot \rightsquigarrow \cdot} \quad \frac{v \in \bar{v} \quad \bar{v} - v \vdash \Gamma \rightsquigarrow \Gamma' \quad \Delta_v \subseteq \text{dom}(\Gamma)}{\bar{v} \vdash \Gamma \rightsquigarrow \Gamma', v}$$

$$\frac{\bar{v} \vdash \Gamma \rightsquigarrow \Gamma' \quad \alpha \notin \text{dom}(\Gamma')}{\bar{v} \vdash \Gamma, \alpha \rightsquigarrow \Gamma', \alpha} \quad \frac{\bar{v} \vdash \Gamma \rightsquigarrow \Gamma' \quad \Gamma' \vdash \Sigma : \Omega}{\bar{v} \vdash \Gamma, X:\Sigma \rightsquigarrow \Gamma', X:\Sigma}$$



This assumes that  $\text{dom}(\Gamma) \not\uparrow \bar{v}$  initially.  $\Gamma'$  is an extension of  $\Gamma$  that includes bindings for  $\bar{v}$ , treated as ordinary type variables, and placed such that they adhere to the scoping constraints encoded in  $\Delta$ . Note that this is a relation: there may be many possible  $\Gamma'$  for a given pair of  $\Gamma$  and  $\bar{v}$  (although they are all equivalent in the sense that they all admit the same set of  $F_\omega$  typing judgments).

With that, for any  $F_\omega$  or IML typing judgement  $J$ , we define

$$\bar{v}; \Gamma \vdash J \quad :\Leftrightarrow \quad \bar{v} \vdash \Gamma \rightsquigarrow \Gamma' \wedge \Gamma' \vdash J \text{ (for some } \Gamma')$$

Furthermore, define well-formedness of resolution substitutions:

$$\begin{aligned} \bar{v}'; \Gamma' \vdash \theta : \bar{v}; \Gamma \quad :\Leftrightarrow \quad & \bar{v} \vdash \Gamma \rightsquigarrow \Gamma'' \wedge \bar{v}'; \Gamma' \vdash \theta : \Gamma'' \wedge \\ & \theta \text{ small} \wedge \text{dom}(\theta) \not\uparrow \bar{v}' \wedge \\ & \forall v \in \bar{v}, \forall v'' \in \text{undet}(\theta v), \Delta_{v''} \subseteq \Delta_v \end{aligned}$$

Well-typed substitutions can be composed:

*Lemma 7.2 (Composition of Substitutions)*

If  $\bar{v}'; \Gamma' \vdash \theta : \bar{v}; \Gamma$  and  $\bar{v}''; \Gamma'' \vdash \theta' : \bar{v}'; \Gamma'$  and  $\text{dom}(\theta) \not\uparrow \bar{v}''$ , then  $\bar{v}''; \Gamma'' \vdash \theta' \circ \theta : \bar{v}; \Gamma$ .

The refined soundness theorem in its full beauty is then stated as follows:

*Theorem 7.3 (Soundness of IML Inference Refined)*

Let  $\text{dom}(\Gamma) \not\uparrow \bar{v}$  and  $\bar{v} \vdash \Gamma \rightsquigarrow \Gamma'$ . Let  $J$  range over the IML inference judgements.

1.  $\Gamma' \vdash \square$  and  $\bar{v} \subseteq \text{dom}(\Gamma')$ .
2. If  $\theta \Gamma \vdash_{\theta'} J$ , then  $\theta \Gamma \vdash_{\theta'} \theta J$ .
3. If  $\Gamma \vdash_{\theta} J$ , then  $\bar{v}'; \theta \Gamma \vdash \theta : \bar{v}; \Gamma$  with  $\bar{v}' - \bar{v}$  fresh.
4. If  $\Gamma \vdash_{\theta'} J$  and  $\bar{v}'; \theta \Gamma \vdash \theta : \bar{v}; \Gamma$ , then  $\bar{v}''; \theta' \Gamma \vdash \theta' : \bar{v}; \Gamma$  with  $\bar{v}'' - \bar{v}' - \bar{v}$  fresh.
5. If  $\Gamma_1; \Gamma_2 \vdash_{\theta'} J$  and  $\bar{v}'; \theta \Gamma \vdash \theta : \bar{v}; \Gamma$  with  $\Gamma = \Gamma_1, \Gamma_2$ , then  $\bar{v}''; \theta' \Gamma \vdash \theta' : \bar{v}; \Gamma$  with  $\bar{v}'' - \bar{v}' - \bar{v}$  fresh and  $\Delta_{\bar{v}''} \subseteq \text{dom}(\Gamma_1)$ .
6. If  $\Gamma \vdash_{\theta} T/D \rightsquigarrow \Xi$ , then  $\bar{v}'; \theta \Gamma \vdash T/D \rightsquigarrow \theta \Xi$ .
7. If  $\Gamma \vdash_{\theta} E/B :_{\eta} \Xi$ , then  $\bar{v}'; \theta \Gamma \vdash E/B :_{\eta} \theta \Xi \rightsquigarrow \theta e$ .
8. If  $\Gamma \vdash_{\theta} \Xi \leq_{\bar{\pi}} \Xi' \rightsquigarrow \delta$ , and  $\bar{v}; \Gamma \vdash \Xi : \Omega$  and  $\bar{v}; \Gamma \vdash \Xi' : \Omega$ , then  $\bar{v}'; \theta \Gamma \vdash \theta \Xi \leq_{\bar{\pi}} \theta \Xi' \rightsquigarrow \theta \delta; \theta f$ .
9. If  $\Gamma \vdash_{\theta}^{\forall} E :_{\eta} \Xi$ , then  $\bar{v}'; \theta \Gamma \vdash E :_{\eta} \theta \Xi \rightsquigarrow \theta e$ .
10. If  $\Gamma \vdash_{\theta}^{\exists} E :_{\eta} \Xi$ , then  $\bar{v}'; \theta \Gamma \vdash E :_{\eta} \theta \Xi \rightsquigarrow \theta e$ .
11. If  $\Gamma \vdash_{\theta} v \approx \Sigma$ , then  $\theta v$  has the same outermost constructor as  $\Sigma$ .
12. If  $\Gamma \vdash_{\theta} \Xi = \Xi'$ , and  $\bar{v}; \Gamma \vdash \Xi : \Omega$  and  $\bar{v}; \Gamma \vdash \Xi' : \Omega$ , then  $\theta \Xi = \theta \Xi'$ .

The proof of the first part is by induction on the derivation of  $\bar{v} \vdash \Gamma \rightsquigarrow \Gamma'$ , for the others by simultaneous induction on the (first) derivation.

**Termination** We can also show that – despite its incompleteness – the inference algorithm at least does not go astray by diverging:

*Theorem 7.4 (Termination of IML Inference)*

All IML type inference judgements terminate.

As for the  $\text{IML}_{\text{ex}}$  typing rules, termination is obvious for most the  $\text{IML}$  inference judgements: the main ones are inductive on the syntactic structure of the language; for the auxiliary Unification judgement, termination can be proved in the usual manner, Instantiation is inductive on the right-hand side type, and Resolution is not even recursive. Again, it only remains subtyping as the problem child.

We can prove its termination by extending the prove from Section 4.3. First, we trivially extend the weight function to handle inference variables:

$$Q[[v]] = 0 \quad S[[v]] = 1$$

The Weight Reduction lemma still holds for  $\delta$ -substitutions on type variables, but we can formulate it analogously for  $\theta$ -substitutions on inference variables, which are always small as well:

*Lemma 7.5 (Weight Reduction under Small Resolution)*

1.  $Q[[\theta\Xi]] = Q[[\Xi]]$ .
2.  $W[[\theta\Xi]] < W[[\Xi]] + \langle 1, 0 \rangle$ .
3.  $W[[\theta\Xi]] \leq W[[\Xi]] + \langle |\theta|, 0 \rangle$ .

The following key lemma postulates that a subtyping derivation either resolves more inference variables than it introduces, or the amount of excess variables (which will come from rule  $\text{ISIMPLL}$ ) is bounded by the number of quantifiers in the types involved.

*Lemma 7.6 (Resolution Progress in Subtyping Inference)*

Let  $\Gamma \vdash_{\theta} J$  an inference derivation and  $Y$  the set of fresh inference variables it generates.

1. If  $J$  is  $\Xi' = \Xi$ , then either  $|Y| = |\theta| = 0$  or  $|Y| < |\theta|$ .
2. If  $J$  is  $\Xi' \leq_{\bar{\pi}} \Xi \rightsquigarrow \delta; f$ , then either  $|Y| = |\theta| = 0$  or  $|Y| < |\theta|$  or  $|Y| - |\theta| \leq Q[[\Xi']] + Q[[\Xi]]$ . Furthermore, if  $\Xi \neq \Xi'$  and either  $\Xi = v$  or  $\Xi' = v$ , then  $v \in \text{dom}(\theta)$ . Else if  $\Xi = v \neq \Xi' = v'$ , then  $|\{v, v'\} \cap \text{dom}(\theta)| = 1$ .

With this, the termination proof can proceed similarly to before, except that we have to take  $\theta$  into account.

*Lemma 7.7 (Termination of Subtyping Inference)*

Let  $\bar{v} \vdash \Gamma \rightsquigarrow \Gamma'$  and  $\Gamma' \vdash \Xi' : \Omega$  and  $\Gamma', \bar{\alpha} \vdash \Xi : \Omega$  and  $\bar{\pi} = \overline{\alpha \bar{\alpha}'}$ . Then  $\Gamma \vdash_{\theta} \Xi' \leq_{\bar{\pi}} \Xi \rightsquigarrow \delta; f$  terminates.

The proof is again by case analysis on the algorithm, mostly as before, but this time using the weight  $W[[\Xi']] + W[[\Xi]] + \langle |\bar{\pi}| + |\bar{v}'|, 0 \rangle$ , with  $\bar{v}' = \text{undet}(\Xi') \cup \text{undet}(\Xi)$ . In each rule, this weight gets smaller for all premises. The only exceptions are the rules  $\text{ISRESL}$  and  $\text{ISRESR}$ , which invoke the Resolution judgement that may introduce additional temporary inference variables. That locally increases the weight of the judgement, but we can show by a simple case analysis of the  $\Sigma$  in these rules that the temporary variables are resolved immediately and weight will decrease again in the recursive invocation of the subtyping judgement.

## 8 Related Work

**Packaged Modules** The first concrete proposal for extending ML with packaged modules was by Russo (2000), and is implemented in Moscow ML. Later work on type systems for

modules routinely included them (Dreyer *et al.*, 2003; Dreyer, 2005; Rossberg & Dreyer, 2013; Rossberg *et al.*, 2014), and variations have been implemented in other ML dialects, such as Alice ML (Rossberg, 2006) and OCaml (Garrigue & Frisch, 2010).

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, which forbid 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 (Harper & Mitchell, 1993). It is a dependent type theory modeling modules as terms of Martin-Löf-style  $\Sigma$  and  $\Pi$  types, closely following MacQueen’s original ideas (MacQueen, 1986). 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. XML allows both  $U_1 : U_2$  and  $U_1 \subseteq U_2$ , which is similar in 1ML. However, XML lacks any account of either sealing or translucency, which makes it fall short as a foundation for contemporary ML modules.

That gap was closed by Harper & Lillibridge’s calculus of *translucent sums* (Harper & Lillibridge, 1994; Lillibridge, 1997), 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 deconstructed into the notion of *singleton types* (Stone & Harper, 2006); they formed the foundation of Dreyer *et al.*’s type theory for higher-order modules (Dreyer *et al.*, 2003). However, to avoid undecidability, this system went back to second-class modules (though it provides packaged modules).

One central 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.* (1990) investigate phase separation for the XML calculus. The beauty of the F-ing approach underlying 1ML is that it enforces phase separation by construction, since it avoids dependent types altogether.

Shields & Peyton Jones (Shields & Peyton Jones, 2002) proposed first-class modules for Haskell. However, these do not provide modular compositionality in the sense discussed in Section 5.3, because they cannot contain type components (only module *types* can), there is no sealing, and type abstraction is manual through raw existential types.

**Applicative Functors** Leroy proposed applicative semantics for functors (Leroy, 1995), as implemented in OCaml. Russo later combined both generative and applicative functors in one language (Russo, 2003) and implemented them in Moscow ML; others followed (Shao, 1999; Dreyer *et al.*, 2003; Dreyer, 2005; Rossberg *et al.*, 2014).

A system like Leroy’s, where *all* functors are applicative, would be incompatible with first-class modules even in a pure language, 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 all generative functors into applicative ones. Like Dreyer (Dreyer, 2005) or F-ing modules (Rossberg *et al.*,

2014), 1ML 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 of the conference version of this article (Rossberg, 2015) sketches a similar extension to 1ML.

In the basic language introduced in the present article, an applicative functor can only be created by sealing a fully transparent functor with pure function type, very much like in Shao's system (Shao, 1999).

**Type Inference** There has been little work that has considered type inference for modules. Russo examined the interplay between core-level inference and modules (Russo, 2003), elegantly dealing with variable scoping via unification under a mixed prefix (which might also form a more elegant basis for 1ML inference than our  $\Delta$ -sets do).

Dreyer & Blume (2007) investigated how functors interfere with the value restriction. They show that – contrary to popular belief – Damas/Milner type inference (Damas & Milner, 1982) with a value restriction (Wright, 1995), as employed in core ML, is no longer complete when combined with ML-style functors. They propose a way to fix this problem by allowing functors to become (implicitly) polymorphic as well. However, their solution requires a *destination-passing style* semantics for type abstraction – inspired by work on recursive modules (Dreyer, 2007) – that is well beyond the realms of what System  $F_\omega$  can express. Consequently, although 1ML can implicitly generalise functors as well and thus is able to handle some of the examples from Dreyer & Blume's paper, it cannot deal with the more interesting ones.

At the same time, there have been ambitious extensions of ML-style type inference with higher-rank or impredicative types (Garrigue & Rémy, 1999; Le Botlan & Rémy, 2003; Vytiniotis *et al.*, 2008; Russo & Vytiniotis, 2009). Unlike those systems, 1ML never tries to infer a polymorphic type annotation: all guessed types are monomorphic; polymorphic parameters always require annotations. The impredicative extension from the Section 5.1 is similar to Garrigue & Rémy's semi-explicit polymorphism (Garrigue & Rémy, 1999), but also allows unboxed polytypes.

On the other hand, 1ML 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 always remain implicit and are effectively inferred as part of the elaboration process.

## 9 Conclusion and Future Work

1ML is a coherent, unified redesign of ML in a single language that subsumes both the ML core and the ML module language. It is no more and no less powerful than System  $F_\omega$  (the higher-order polymorphic lambda calculus), but unlike the latter, supports type inference and modular composition.

But 1ML, as shown here, is but a first step. There are many possible improvements and extensions that are worth investigating.

**Implementation** We have implemented a simple prototype interpreter for IML ([mpi-sws.org/~rossberg/1ml/](http://mpi-sws.org/~rossberg/1ml/)) that allows running smallish examples. But it would be great to gather more experience with a “real” implementation, especially regarding the viability and user-friendliness of IML type inference.

**Applicative Functors** To keep the present paper simple, it adopts a rather basic notion of applicative functor. It is not hard to extend it to *pure sealing* à la F-ing modules – as mentioned earlier, a sketch can be found in the Technical Appendix of the conference version of this article (Rossberg, 2015). A greater challenge would be to go a step further and make it properly *abstraction-safe*, by tracking value identities as described in Section 8 of Rossberg *et al.* (2014).

**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* (Wadler & Blott, 1989) or Coq-style *canonical structures* (Mahboubi & Tassi, 2013). In particular, generalised implicit functions could be a semantic foundation for modular implicits (White *et al.*, 2014).

**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 with some of the more advanced approaches to inference described in literature, some of which we mentioned in Section 8. However, it is non-obvious how well these methods will interact with the implicit quantifier pushing in IML’s elaboration.

**Row Polymorphism** Replacing subtyping with a more potent form of polymorphism might increase expressiveness and lead to better inference: in particular, *row polymorphism* (Rémy, 1989) could express width subtyping without compromising completeness of inference. The main issue with such an extension probably is its effect on the small/large type universe distinction, since many current uses of subtyping on small types will require large types once they involve quantification of row variables.

**Effect Polymorphism** Similarly, the utility of IMLs effect system and its pure function types could be greatly increased if higher-order functions could vary their effect by employing *effect polymorphism* (Talpin & Jouvelot, 1992). We have recently explored this direction in a separate paper (Rossberg, 2016). Its most interesting implication is that it enables “functors” to be polymorphic over their generativity, i.e., whether they behave applicatively or generatively. Incidentally, this ability finally allows expressing the motivating examples of MacQueen’s “truly higher-order modules” (MacQueen & Tofte, 1994; Kuan & MacQueen, 2009), which previously have eluded a nice type-theoretic explanation, and were not supported by any other module type system.

**Recursive Modules** Recursive modules would add a whole new dimension of expressiveness – and complexity! – to IML. With MixML (Rossberg & Dreyer, 2013) we gave a fully general design for recursive modules, elaborating into an extension of System F. It would be interesting (but complicated) to redo it IML-style, in order to achieve a more uniform

treatment of recursion than IML has. The result of such a cross-breed would presumably share some similarities with Scala (Odersky & Zenger, 2005), but with a more expressive semantics for type equivalences.

**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 (a.k.a. existential types). Incorporating them, along with the quantifier pushing of our elaboration, into a dependent type system poses a thrilling challenge.

### References

- Barendregt, Henk. (1992). Lambda calculi with types. *Chap. 2, pages 117–309 of: Abramsky, Samson, Gabbay, Dov, & Maibaum, T.S.E. (eds), Handbook of logic in computer science, vol. 2.* Oxford University Press.
- Biswas, Sandip K. (1995). Higher-order functors with transparent signatures. *ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL).*
- Damas, Luis, & Milner, Robin. (1982). Principal type-schemes for functional programs. *ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL).*
- Dreyer, Derek. (2005). *Understanding and Evolving the ML Module System.* Ph.D. thesis, Carnegie Mellon University.
- Dreyer, Derek. (2007). Recursive type generativity. *Journal of Functional Programming (JFP)*, **17**(4&5), 433–471.
- Dreyer, Derek, & Blume, Matthias. (2007). Principal type schemes for modular programs. *European Symposium on Programming (ESOP).*
- Dreyer, Derek, & Rossberg, Andreas. (2008). Mixin' up the ML module system. *ACM SIGPLAN International Conference on Functional Programming (ICFP).*
- Dreyer, Derek, Crary, Karl, & Harper, Robert. (2003). A type system for higher-order modules. *ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL).*
- Garrigue, Jacques, & Frisch, Alain. (2010). First-class modules and composable signatures in Objective Caml 3.12. *ACM SIGPLAN Workshop on ML.*
- Garrigue, Jacques, & Rémy, Didier. (1999). Semi-explicit first-class polymorphism for ML. *Information and computation*, **155**(1-2).
- Harper, Robert, & Lillibridge, Mark. (1994). A type-theoretic approach to higher-order modules with sharing. *ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL).*
- Harper, Robert, & Mitchell, John C. (1993). On the type structure of Standard ML. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, **15**(2), 211–252.
- Harper, Robert, & Pierce, Benjamin C. (2005). Design considerations for ML-style module systems. *Chap. 8 of: Pierce, Benjamin C. (ed), Advanced Topics in Types and Programming Languages.* MIT Press.
- Harper, Robert, & Stone, Chris. (2000). A type-theoretic interpretation of Standard ML. *Proof, Language, and Interaction: Essays in Honor of Robin Milner.* MIT Press.
- Harper, Robert, Mitchell, John C., & Moggi, Eugenio. (1990). Higher-order modules and the phase distinction. *ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL).*
- Kuan, George, & MacQueen, David. (2009). Engineering higher-order modules in SML/NJ. *International Symposium on the Implementation and Application of Functional Languages (IFL).* LNCS, vol. 6041.

- Le Botlan, Didier, & Rémy, Didier. (2003). MLF: Raising ML to the power of System F. *ACM SIGPLAN International Conference on Functional Programming (ICFP)*.
- Leroy, Xavier. (1995). Applicative functors and fully transparent higher-order modules. *ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL)*.
- Lillibridge, Mark. (1997). *Translucent sums: A foundation for higher-order module systems*. Ph.D. thesis, Carnegie Mellon University.
- MacQueen, David B. (1986). Using dependent types to express modular structure. *ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL)*.
- MacQueen, David B., & Tofte, Mads. (1994). A semantics for higher-order functors. *European Symposium on Programming (ESOP)*.
- Mahboubi, Assia, & Tassi, Enrico. (2013). Canonical structures for the working Coq user. *Pages 19–34 of: Blazy, Sandrine, Paulin-Mohring, Christine, & Pichardie, David (eds), Interactive theorem proving*. LNCS, no. 7998. Springer-Verlag.
- Milner, Robin. (1978). A theory of type polymorphism in programming languages. *Journal of computer and system sciences*, **17**, 348–375.
- Milner, Robin, Tofte, Mads, Harper, Robert, & MacQueen, David. (1997). *The definition of Standard ML (revised)*. MIT Press.
- Mitchell, John C., & Plotkin, Gordon D. (1988). Abstract types have existential type. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, **10**(3), 470–502.
- Odersky, Martin, & Zenger, Matthias. (2005). Scalable component abstractions. *Object-Oriented Programming, Systems, Languages and Applications (OOPSLA)*. ACM Press.
- Rémy, Didier. (1989). Records and variants as a natural extension of ML. *ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL)*.
- Rossberg, Andreas. (2006). The Missing Link – Dynamic components for ML. *ACM SIGPLAN International Conference on Functional Programming (ICFP)*.
- Rossberg, Andreas. (2015). *IML – Core and modules united (Technical Appendix)*. [mpi-sws.org/~rossberg/1ml/](http://mpi-sws.org/~rossberg/1ml/).
- Rossberg, Andreas. (2016). IML with special effects. *Pages 336–355 of: A List of Successes That Can Change the World – WadlerFest 2016*. LNCS, no. 9600. Springer-Verlag.
- Rossberg, Andreas, & Dreyer, Derek. (2013). Mixin’ up the ML module system. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, **35**(1).
- Rossberg, Andreas, Russo, Claudio, & Dreyer, Derek. (2014). F-ing modules. *Journal of Functional Programming (JFP)*, **24**(5), 529–607.
- Russo, Claudio, & Vytiniotis, Dimitris. (2009). QML: Explicit first-class polymorphism for ML. *ACM SIGPLAN Workshop on ML*.
- Russo, Claudio V. (1999). Non-dependent types for Standard ML modules. *International Conference on Principles and Practice of Declarative Programming (PPDP)*.
- Russo, Claudio V. (2000). First-class structures for Standard ML. *Nordic Journal of Computing*, **7**(4), 348–374.
- Russo, Claudio V. (2003). Types for Modules. *Electronic Notes in Theoretical Computer Science (ENTCS)*, **60**.
- Shao, Zhong. (1999). Transparent modules with fully syntactic signatures. *ACM SIGPLAN International Conference on Functional Programming (ICFP)*.
- Shields, Mark, & Peyton Jones, Simon. (2002). First-class modules for Haskell. *Pages 28–40 of: International Workshop on Foundations of Object-Oriented Languages (FOOL)*.
- Stone, Christopher A., & Harper, Robert. (2006). Extensional equivalence and singleton types. *ACM Transactions on Computational Logic (TOCL)*, **7**(4), 676–722.

- Talpin, Jean-Pierre, & Jouvelot, Pierre. (1992). Polymorphic type, region and effect inference. *Journal of Functional Programming (JFP)*, 2(3), 245-271.
- Vytiniotis, Dimitrios, Weirich, Stephanie, & Peyton Jones, Simon. (2008). FPH: First-class polymorphism for Haskell. *ACM SIGPLAN International Conference on Functional Programming (ICFP)*.
- Wadler, Philip, & Blott, Stephen. (1989). How to make ad-hoc polymorphism less ad hoc. *ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL)*.
- White, Leo, Bour, Frédéric, & Yallop, Jeremy. (2014). Modular implicits. *Pages 22–63 of: Kiselyov, Oleg, & Garrigue, Jacques (eds), ACM SIGPLAN ML Family / OCaml Workshops*. EPTCS, no. 198.
- Wright, Andrew. (1995). Simple imperative polymorphism. *LISP and Symbolic Computation (LASC)*, 343–356.