

Polymonads

Extended version of paper submitted to POPL'13

Nataliya Guts[†] Michael Hicks[†] Nikhil Swamy^{*} Daan Leijen^{*} Gavin Bierman^{*}

[†]University of Maryland, College Park

^{*}Microsoft Research

Abstract

From their semantic origins to their use in structuring effectful computations, monads are now also used as a programming pattern to structure code in a number of important scenarios, including functional reactivity, information flow tracking and probabilistic computation. However, whilst these examples are inspired by monads they are not strictly speaking monadic but rather something more general. The first contribution of this paper is the definition of a new structure, the polymonad, which subsumes monads and encompasses the monad-like programming patterns that we have observed. A concern is that given such a general setting, a program would quickly become polluted with polymonadic coercions, making it hard to read and maintain. The second contribution of this paper is to build on previous work to define a polymorphic type inference algorithm that supports programming with polymonads using a direct style, e.g., as if computations of type $M \tau$ were expressions of type τ . During type inference the program is rewritten to insert the necessary polymonadic coercions, a process that we prove is coherent—all sound rewritings produce programs with the same semantics. The resulting programming style is powerful and lightweight.

1. Monads and more

For most programmers, a monad is an abstract datatype represented by a unary type constructor m and two operations, *bind* and *unit* (a.k.a., *return*), with the following signature:

$$\begin{aligned} \text{bind} &: \forall \alpha, \beta. m \alpha \rightarrow (\alpha \rightarrow m \beta) \rightarrow m \beta \\ \text{unit} &: \forall \alpha. \alpha \rightarrow m \alpha \end{aligned}$$

Implementations of these operations are expected to obey the *monad laws*, which enable reasoning about transformations of monadic programs, both by programmers and by tools such as optimizing compilers [19].

Since the time that Moggi first connected monads to effectful computation [18], monads have proven to be a surprisingly versatile computational structure. Perhaps best known as the foundation of Haskell's support for state, I/O, and other effects, monads have also been used to structure APIs for libraries that implement a wide range of programming tasks, including parsers [13], probabilistic computations [22], functional reactivity [8, 4], and information flow tracking [23].

While conceptually simple, programming directly against a monadic API is impractical. Programmers must insert calls to *bind* and *unit* pervasively, and when composing multiple monads, *morphisms* between the monads must also be inserted. However, effective type inference algorithms have been devised to reduce this burden. In the context of Haskell, monadic type inference relies on the mechanisms of typeclasses and a specialized syntax (the *do* notation) to infer the placement of *binds* and *units*. For ML, our

own prior work [24] has shown that the existing *let*-structure of a call-by-value program can be used to infer the placement of the morphisms in addition to the *binds* and *units*.

The monadic programming pattern is sufficiently appealing that many researchers have developed subtle variations to adapt and apply monads to new problem domains. Examples include Wadler and Thiemann's [27] indexed monad for typing effectful computations; Atkey's parameterized monad [2], which has been used to encode disciplines like regions [15] and session types [21]; Devriese and Piessens' [7] monad-like encodings for information flow controls; Danielsson's [6] counting monad for computational complexity; and many others. Oftentimes these extensions are used to prove stronger properties about computations than would be possible with monads, or to prevent undesirable behavior (such as illegal information flows, memory errors, etc.).

We observe that in each of these cases a family of abstract datatypes $\{m_1, \dots, m_n\}$ with *bind* and *unit*-like operations is provided. But, unlike for traditional monads, the *binds* have signatures of the form

$$\forall \alpha, \beta. m_1 \alpha \rightarrow (\alpha \rightarrow m_2 \beta) \rightarrow m_3 \beta$$

We call binds of this form *non-uniform binds* and refer to the collection of m_i and the binds among them as a *polymonad*.

This paper explores the idea of polymonads, a generalization of monads and monad morphisms. Section 2 defines polymonads precisely, including laws that are analogs of the monad and morphism laws. While the notion of a non-uniform bind has been considered previously [16], the laws underlying the behavior of polymonads have never been articulated. As with monads, these laws are important for reasoning about intuitive program transformations. The laws are also important for type inference, as we explain shortly. We show that every set of monads and monad morphisms can be encoded using polymonads while obeying the polymonad laws.

Next, in Section 3, we present a variety of examples and show them to be polymonads. Our examples include several from the literature mentioned previously, as well as two new polymonadic constructions. The first is an encoding of information flow controls in the presence of side effects, and the second is an encoding of contextual effects [20]. We show that each of our examples obey the polymonad laws.

To make polymonads easier to program with, in Section 4 we develop a novel type inference and rewriting algorithm for an ML-like, call-by-value programming language. Rather than write non-uniform binds in programs directly, the programmer can use expressions of type $m \tau$ as if they were of type τ , thus programming in a direct style. Our algorithm will automatically infer which non-uniform binds are needed and where they should be placed. While in Haskell the *do* notation must be used to identify where *binds* and *units* should be placed (and type class inference determines which binds/units to insert), for ML no special syntax is needed. General-

constructors	\mathcal{M}	::=	$Id \mid M, \mathcal{M}$
signature	$\Sigma_{\mathcal{M}}$::=	$\mathbf{b}_{Id, Id, Id} \mid \mathbf{b}:s, \Sigma_{\mathcal{M}}$
bind type	s	::=	$\forall \alpha \beta. m_1 \alpha \rightarrow (\alpha \rightarrow m_2 \beta) \rightarrow m_3 \beta$

Figure 1. Syntax of a polymonadic signature

izing our prior work [24], normal let bindings and applications naturally identify where binds should be placed when programming with (poly)monadic computations, thanks to ML’s call-by-value, fixed evaluation order. We prove that our algorithm enjoys principal types.

One might wonder why a special algorithm is needed: could not a standard type inference algorithm based on OML [14], the framework of qualified types on which Haskell’s typeclass inference algorithm is based, be applied to polymonads? Indeed, this has been done for the typeclass-based encoding of polymonads in Kmett’s `Control.Monad.Parameterized` package for Haskell [16]. But this approach has two important problems.

First, because vanilla type class inference is unaware of the polymonadic laws many simple programs are rejected. For example, we cannot use `Control.Monad.Parameterized` to check that the program $\backslash x. \text{id } (\text{id } x)$ has the type $\alpha \rightarrow m \alpha$ for some polymonad m . Haskell infers type class constraints for this program that could allow ambiguous instantiations: different instantiations of certain type variables could result in different program semantics. In Section 5, we prove a coherence result for our type inference algorithm for polymonadic programs, extending our prior result for monadic programs. We show that by appealing to the polymonad laws, instantiations that OML would conservatively deem as ambiguous can in fact be accepted, as there is no possibility for ambiguity of a program’s semantics. Thus our algorithm permits many useful functions that would otherwise be rejected.

Second, principal types turn out to be exceedingly hard to read. Even simple terms have a large number of quantified variables and constraints. However, as we show in Section 6, by applying constraint simplifications justified by the polymonad laws we can *improve* types to make them much simpler without negatively impacting their generality.

Section 7 describes a simple prototype implementation we have built to check the examples in the paper.

Concurrently with our work, Tate [25] developed a semantic construction he calls *productors*. After some interaction, we discovered that polymonads match a specialization of productors applied to what Tate calls *semi-strict* effect systems. Tate proves that such effect systems are quite general, and thus polymonads are equally general. As polymonads are focused on practical programming, while productors are focused on semantics, our two papers provide complementary views of an exciting, underlying generalization of monads. We discuss more related work in Section 8.

2. Polymonads

We begin by defining polymonads from the perspective of a programmer, who may think of them as a set of abstract types with a collection of operations which, when taken together, must respect a specific set of equations. To help provide intuition as to the implications of these equations, we prove that every collection of monad and monad morphisms also induces a polymonad. The increased expressive power of polymonads is put on display in the next section, which incorporates them into a simple programming language and presents a series of examples.

2.1 Syntax

Figure 1 defines the signature $\Sigma_{\mathcal{M}}$ of a set \mathcal{M} of polymonadic abstract types. The set \mathcal{M} contains unary type constructors M , including a distinguished type constructor Id . We expect all the type constructors occurring in $\Sigma_{\mathcal{M}}$ to be in \mathcal{M} . We write m as a metavariable ranging over monadic type constructors. When the set of constructors is irrelevant or clear from the context, we refer to the signature as Σ , rather than $\Sigma_{\mathcal{M}}$.

The signature represents a map from names \mathbf{b} to types s where each type has the shape $\forall \alpha \beta. m_1 \alpha \rightarrow (\alpha \rightarrow m_2 \beta) \rightarrow m_3 \beta$. We refer to each \mathbf{b} as a *bind*, for its relation to the monadic bind operation [26], although, unlike a monadic bind, a polymonadic bind involves three abstract types. We also refer to the signature Σ as the *bind set*. Syntactically, we require the bind set Σ to contain an element $\mathbf{b}_{Id, Id, Id}$ which implicitly has the type $\forall \alpha \beta. Id \alpha \rightarrow (\alpha \rightarrow Id \beta) \rightarrow Id \beta$. One may think of $Id \tau$ as a synonym for τ , in which case (adhering to the laws below) $\mathbf{b}_{Id, Id, Id}$ is reverse apply.

We use the following shorthands:

- $(m_1, m_2) \triangleright m$ is the type $\forall \alpha \beta. m_1 \alpha \rightarrow (\alpha \rightarrow m_2 \beta) \rightarrow m \beta$.
- $(m_1, m_2) \triangleright m \in \Sigma$ means $\exists \mathbf{b}. \mathbf{b} : (m_1, m_2) \triangleright m \in \Sigma$. This notation is unambiguous as we assume that binds \mathbf{b} with the same type s have the same semantics.
- We write $\mathbf{b}_{m_1, m_2, m} \in \Sigma$ to mean $\mathbf{b} : (m_1, m_2) \triangleright m \in \Sigma$.
- We define unit_m , with type $\forall \alpha. \alpha \rightarrow m \alpha$, to be $\lambda x. \mathbf{b}_{Id, Id, m} x (\lambda z. z)$, and write $\text{unit}_m \in \Sigma$ to mean $(Id, Id) \triangleright m \in \Sigma$.
- We write $m_1 \succ m_2$ to mean $(m_1, Id) \triangleright m_2$ or $(Id, m_1) \triangleright m_2$.

2.2 Polymonad laws

We impose several requirements on Σ for it to define a valid polymonad. First, we require each bind operation to be given an *interpretation* in an underlying typed lambda calculus. We choose System F [11] as the underlying calculus in this paper, although for clarity, we omit explicit type abstraction and application in System F terms. Thus, we write $\mathbf{b}_{m_1, m_2, m_3} e_1 e_2$ to mean the application in System F of the interpretation of $\mathbf{b}_{m_1, m_2, m_3}$ to two terms $e_1 : m_1 \tau$ and $e_2 : \tau \rightarrow m_2 \tau'$. Equivalent formulations in other calculi are also feasible.

Any valid interpretation of a bind set Σ is expected to respect certain equations between well-typed terms, where we interpret equality as β, η -equivalence. We also require a bind set to be closed, in the sense that the presence of certain binds in Σ mandate the existence in Σ of other, related binds. These requirements are expressed by the six polymonad laws below (which we explain in detail shortly); we write $\models \Sigma$ when a signature satisfies these laws.

Functorial law: For all $m \in \mathcal{M}$, we require $(m, Id) \triangleright m \in \Sigma$. That is, every m comes equipped with a *map* function.

Left identity: For all m_1, m_2, e, k , if $\mathbf{b}_{m_1, m_2, m_2} \in \Sigma$ and $\text{unit}_{m_1} \in \Sigma$, then $\mathbf{b}_{m_1, m_2, m_2} (\text{unit}_{m_1} e) k = k e$.

Right identity: For all m_1, m_2, e , if $\mathbf{b}_{m_1, m_2, m_1} \in \Sigma$ and $\text{unit}_{m_2} \in \Sigma$, then $\mathbf{b}_{m_1, m_2, m_1} e \text{unit}_{m_2} = e$.

Associativity:

- (1) For all m_1, m_2, m_3, m_{123} ,
 $\exists m_{12}. \{ \mathbf{b}_{m_1, m_2, m_{12}}, \mathbf{b}_{m_{12}, m_3, m_{123}} \} \subseteq \Sigma$
if and only if
 $\exists m_{23}. \{ \mathbf{b}_{m_1, m_{23}, m_{123}}, \mathbf{b}_{m_2, m_3, m_{23}} \} \subseteq \Sigma$.
- (2) For all $m_1, m_2, m_3, m_{12}, m_{23}, m_{123}, e_1, k_2, k_3$,
if $\{ \mathbf{b}_{m_1, m_2, m_{12}}, \mathbf{b}_{m_{12}, m_3, m_{123}}, \mathbf{b}_{m_1, m_{23}, m_{123}}, \mathbf{b}_{m_2, m_3, m_{23}} \} \subseteq \Sigma$,
then $\mathbf{b}_{m_{12}, m_3, m_{123}} (\mathbf{b}_{m_1, m_2, m_{12}} e_1 k_2) k_3 = \mathbf{b}_{m_1, m_{23}, m_{123}} e_1 (\lambda x. \mathbf{b}_{m_2, m_3, m_{23}} (k_2 x) k_3)$.

monad sig. $\mathcal{S} ::= (D_{Id}, D_1, \dots, D_n, f_{M_1, M'_1}, \dots, f_{M_k, M'_k})$
monad $D ::= (M, \text{unit}_M, \text{bind}_M), f_{M, M}$
morphism $f_{M, N} : M \triangleright N$

$$\frac{M \in \mathcal{S}}{S \models M \triangleright M} \quad \frac{M \in \mathcal{S}}{S \models Id \triangleright M} \quad \frac{f : M_1 \triangleright M_2 \in \mathcal{S}}{S \models M_1 \triangleright M_2}$$

$$\frac{S \models M_1 \triangleright M_2 \quad S \models M_2 \triangleright M_3}{S \models M_1 \triangleright M_3}$$

$$\begin{aligned} \text{bind}_M (\text{unit}_M e) k &= k e & (i) \\ \text{bind}_M e \text{unit}_M &= e & (ii) \\ \text{bind}_M (\text{bind}_M e k_1) k_2 &= & \\ \text{bind}_M e (\lambda x. \text{bind}_M (k_1 x) k_2) & & (iii) \\ f_{M_1, M_2} (\text{unit}_{M_1} e) &= \text{unit}_{M_2} e & (iv) \\ f_{M_1, M_2} (\text{bind}_{M_1} e k) &= & \\ s\text{bind}_{M_2} (f_{M_1, M_2} e) (f_{M_1, M_2} \circ k) & & (v) \\ f_{M_2, M_3} \circ f_{M_1, M_2} &= f_{M_1, M_3} & (vi) \\ f_{M, M} e &= e & (vii) \end{aligned}$$

Figure 2. Monads

Paired morphisms law: For all $m_1, m_2 \in \mathcal{M}$, if $(m_1, Id) \triangleright m_2 \in \Sigma$ then $(Id, m_1) \triangleright m_2 \in \Sigma$, and vice versa.

Composition closure law: For all $m_1, m_2, m_3, m'_1, m'_2, m'_3$, if $(m_1, m_2) \triangleright m_3 \in \Sigma$ and if $\{m'_1 \succ m_1, m'_2 \succ m_2, m_3 \succ m'_3\} \subseteq \Sigma$, then $(m'_1, m'_2) \triangleright m'_3 \in \Sigma$.

We choose these laws for two reasons. First, we aim for poly-monads to be a generalization of monads, yielding reasoning principles with which programmers are already familiar, e.g., associativity of bind operations, which provides justification for many simple program transformations. This motivates our choice of the first four laws.

Second, we aim to exploit the polymonad laws to devise an coherent constraint solving procedure for type inference. Thus, when certain binds can be defined unambiguously in terms of other binds, we require these to be present in the bind set. This is the motivation behind the last two laws.

For example, consider the paired morphism law. We think of both $(m_1, Id) \triangleright m_2$ and $(Id, m_1) \triangleright m_2$ as *morphisms* between m_1 and m_2 and write $m_1 \succ m_2$ for a morphism between m_1 and m_2 . Given a $\mathbf{b}_{m_1, Id, m_2} \in \Sigma$ one can easily define $\mathbf{b}_{Id, m_1, m_2}$ as $\lambda x. \lambda f. \mathbf{b}_{m_1, Id, m_2}(fx)(\lambda x. x)$. It is easy to check that this is the only interpretation of $\mathbf{b}_{Id, m_1, m_2}$ consistent with the other laws. A similar construction yields $\mathbf{b}_{m_1, Id, m_2}$ from $\mathbf{b}_{Id, m_1, m_2}$.

The composition closure law allows morphisms to be composed with binds. For example, given $\mathbf{b}_1: m'_1 \succ m_1$, $\mathbf{b}_2: m'_2 \succ m_2$, $\mathbf{b}_3: m_3 \succ m'_3$, and $\mathbf{b}: (m_1, m_2) \triangleright m_3$, one can define $\mathbf{b}_{m'_1, m'_2, m'_3}$ as $\lambda x. \lambda f. \mathbf{b}_3 (\mathbf{b} (\mathbf{b}_1 x (\lambda y. y)) (\lambda x. \mathbf{b}_2 (fx) (\lambda y. y))) (\lambda y. y)$, where we apply \mathbf{b}_1 and \mathbf{b}_2 contravariantly, and \mathbf{b}_3 covariantly. As with the paired morphism law, it is easy to check that this construction is the only way to compose the binds and morphisms in a manner compatible with the rest of the laws, so we simply require these binds to be present.

2.3 Comparison to conventional monads

Here we formalize the relationship between conventional monads and polymonads.

Monad families. We define a family of monads using a *monadic signature* \mathcal{S} , defined in Figure 2. A signature consists of a series of monad definitions D and a series of morphisms $f_{M, N}$ between

monads M and N . The relation $S \models M \triangleright N$ indicates that according to signature S we can convert a monad M into a monad N , either reflexively, via a morphism, or via morphism composition. We also assume the existence of a monad Id whose semantics is the identity for unit_{Id} and reverse apply for bind_{Id} ; as such unit_M for all M can be viewed as a morphism from Id to M .

A monadic signature is well-formed, denoted $\models S$, if morphisms in S only refer to monads also defined in S , and if the morphisms and monadic operators satisfy the laws given at the bottom of Figure 2. Observe that the polymonad laws are structurally similar the conventional monad laws. In particular, making all indices m_1, \dots, m_6 equal to a single M maps the polymonad laws *left identity*, *right identity*, and *associativity* to the monad laws (i)–(iii), respectively. In fact, we can prove that all monads are polymonads.

Translating signatures. We write $\langle S \rangle$ to denote the polymonad $\mathcal{M}, \Sigma_{\mathcal{M}}$ equivalent to monadic signature S , where $\mathcal{M} = \{M \mid (M, \text{unit}_M, \text{bind}_M) \in S\}$ and $\Sigma_{\mathcal{M}} = \text{Clos}(S)$, defined as follows:

$$\text{Clos}(S) = \{ \mathbf{b}_{M_1 M_2 M} : (M_1, M_2) \triangleright M \mid S \models M_1 \triangleright M \wedge S \models M_2 \triangleright M \}$$

Given the definitions of its binds and morphisms, we can define each of the polymonadic binds $\mathbf{b}_{M_1 M_2 M} \in \text{Clos}(S)$ as follows:

$$\mathbf{b}_{M_1 M_2 M} x g = \text{bind}_M (f_{M_1, M} x) (\lambda y. f_{M_2, M} (g y))$$

That the morphisms $f_{M_1, M}$ and $f_{M_2, M}$ are present in S follows directly from the definition of $\text{Clos}(S)$.

We can prove that polymonad and monad laws coincide for polymonadic representation of monads.

Lemma 1. $\models S$ iff $\models \langle S \rangle$ for all S .

Proof. See Appendix A. □

Example As an example of this translation, consider the signature $S = (D_{Id}, (M, \text{bind}_M, \text{unit}_M), f_{M, M})$. The corresponding polymonad is $\mathcal{M} = M, Id$ and

$$\begin{aligned} \Sigma_{\mathcal{M}} &= \mathbf{b}_{MMM} : (M, M) \triangleright M, \mathbf{b}_{Id, Id, M} : (Id, Id) \triangleright M, \\ &\mathbf{b}_{Id, M, M} : (Id, M) \triangleright M, \mathbf{b}_{M, Id, M} : (M, Id) \triangleright M, \\ &\mathbf{b}_{Id, Id, Id} : (Id, Id) \triangleright Id \end{aligned}$$

We can define the polymonad binds using bind_M and unit_M as follows (after simplifying terms produced by the definition above):

$$\begin{aligned} \mathbf{b}_{M, M, M} &= \text{bind}_M \\ \mathbf{b}_{Id, Id, M} x f &= \text{unit}_M (f x) \\ \mathbf{b}_{Id, M, M} x f &= \text{bind}_M (\text{unit}_M x) f = f x \\ \mathbf{b}_{M, Id, M} x f &= \text{bind}_M x (\lambda y. \text{unit}_M (f y)) \end{aligned}$$

2.4 Categorical foundations

Whilst the focus of our work is on the programmatic aspects of polymonads, we have developed some categorical analysis of polymonads. Our categorical model consists of a collection of functors (modeling the type constructors) and over this a collection of natural transformations of the form $T_1(T_2(A)) \rightarrow T_3(A)$ where the T_i are taken from the collection of functors. The collection of natural transformations must satisfy a number of conditions that are the categorical analogs of the polymonad laws from Section 2.2. Some details of our categorical model appear in Appendix D. Interestingly, almost identical categorical constructions have been independently proposed by Tate as models of his generalized effects framework [25]. (Indeed we are grateful to Tate for spotting some shortcomings in our initial model.)

3. Programming with monads

The previous section defines a core, abstract definition for monads. In this section, we define λPM , a lambda calculus with support for monadic programming. λPM integrates monads with constructs familiar from the polymorphic lambda calculus. Notably, rather than insisting on having only unary type constants in \mathcal{M} , we will permit monadic type constructors to have additional parameters, and to carry constraints on these parameters within types.

We develop a series of λPM examples demonstrating the usefulness of monads. With our type inference algorithm, detailed in the next section, programmers write λPM programs in direct style, our algorithm infers monadic types and then elaborates the typed source programs into System F with explicit applications of the monadic bind operations.

3.1 λPM : A language for monadic programming

Figure 3 presents the syntax of λPM , a call-by-value, polymorphic lambda calculus. Its term language is standard—we have variables x , constants c , undecorated λ -abstractions, function application and let-bindings. We expect λPM programs to be written in direct style, as computations of type $m \tau$ are simply expressions of type τ . We then infer where and which bind operations must be inserted to typecheck a λPM program against a monadic interface.

To make inference feasible, we rely crucially on λPM 's call-by-value structure. Following our prior work on monadic programming for ML, we restrict the shape of types assignable to a λPM program by separating value types τ from the types of monadic computations $m \tau$. The co-domain of every function is required to be a computation type $m \tau$, although pure functions can be typed $\tau \rightarrow \tau'$, which is a synonym for $\tau \rightarrow \text{Id } \tau'$. We also include types $T \bar{\tau}$ for fully applied type constructors, e.g., list *int*.

Programs can also be given type schemes σ that are polymorphic in their monads, e.g., $\forall \alpha \mu \beta. (\alpha \rightarrow \mu \beta) \rightarrow \alpha \rightarrow \mu \beta$. Here, the variable α ranges over all value types τ , while μ ranges over computation types m . Type schemes may also be qualified by a set of bind constraints, P . For example, $\forall \alpha \mu. (\mu, \text{Id}) \triangleright M \Rightarrow (\text{int} \rightarrow \mu \text{int}) \rightarrow M \text{int}$ is the type of a function that abstracts over a morphism $\mu \succ M$.

λPM is parameterized by a set of monad type constructors \mathcal{M} , where each constructor $M/k \in \mathcal{M}$ is a $(k + 1)$ -ary type constructor (unlike in Section 2 where we only considered unary constructors). For example, we may write monadic types like $\text{ST } h \text{ int}$, indexing the state monad $\text{ST}/1$ with a phantom type h for a heap variable, as is common in a language like Haskell. We often omit the arity k for brevity. Note, our metavariable m for monadic types now includes both monad variables μ , as well as monadic constants $M \bar{\tau}$ applied to a sequence of type indices. Our intention is that type indices are *phantom*, meaning that they are used as a type-level representation of some property of the monad's current state, but a monadic bind's implementation does not depend on them. For example, we would expect that binds would treat objects of type $\text{ST } h \tau$ uniformly, for all h ; different values of h would be used to statically prevent unsafe operations like double-frees or dangling pointer dereferences. If an object has different states that would affect the semantics of binds, the programmer can use different constructors M for each state (rather than different type indices and the same constructor).

As before, a bind set $\Sigma_{\mathcal{M}}$ is a map from bind names \mathbf{b} to their types s . However, unlike in Section 2, where we required Σ to be specified intensionally as a set, here, we allow an extensional definition of Σ using a language of *theory constraints* Φ . λPM 's type system is parametric in the choice of theory constraints Φ . This generality allows us to encode a variety of prior monad-like systems with λPM .

values	v	$::=$	$x \mid c \mid \lambda x. e$
expressions	e	$::=$	$v \mid e_1 e_2 \mid \text{let } x = e_1 \text{ in } e_2$
types	τ	$::=$	$() \mid \alpha \mid T \bar{\tau} \mid \tau_1 \rightarrow m \tau_2$
type schemes	σ	$::=$	$\forall \bar{v}. P \Rightarrow \tau$
monadic types	m	$::=$	$M \mid \mu$
ground monads	M	$::=$	$M \bar{\tau}$
type variables	ν	$::=$	$\alpha \mid \mu$
bind constraint	π	$::=$	$(m_1, m_2) \triangleright m$
bind constraints	P	$::=$	$\cdot \mid \pi, P$
substitutions	θ	$::=$	$\cdot \mid \mu \mapsto m \mid \alpha \mapsto \tau \mid \theta, \theta$
environment	Γ	$::=$	$\cdot \mid \Gamma, c : \sigma \mid \Gamma, x : \sigma$

Parameterized by:

k -ary constructors	\mathcal{M}	$::=$	$\text{Id}/0 \mid M/k, \mathcal{M}$
bind set	$\Sigma_{\mathcal{M}}$	$::=$	$\mathbf{b}_{\text{Id}, \text{Id}, \text{Id}} \mid \mathbf{b} : s, \Sigma_{\mathcal{M}}$
bind types	s	$::=$	$\forall \bar{\alpha}. \Phi \Rightarrow (M_1, M_2) \triangleright M_3$
theory constraints	Φ		
theory entailment	\models	$::=$	$2^{\Sigma \times P \times \theta \times \bar{\mathbf{b}}}$

Figure 3. Syntax of λPM

For example, to model Wadler and Thiemann's [27] indexed monad, which represents a type and effect system, we can introduce a monadic constructor $W/1$, and use $W \epsilon \tau$ to represent a computation that produces a τ -result after exhibiting effects contained within the set ϵ . To specify a monadic bind for W , we can use λPM 's bind type $\forall \alpha, \beta, \epsilon_1, \epsilon_2, \epsilon_3. (\epsilon_3 = \epsilon_1 \cup \epsilon_2) \Rightarrow (W \epsilon_1, W \epsilon_2) \triangleright W \epsilon_3$. Here, we have instantiated the theory Φ to include equality and set operators like \cup , and the bind type indicates (informally) that when composing two computations, the effects are additive.

To interpret the theory constraints, λPM requires a theory entailment relation \models , where elements of this relation are written $\Sigma \models \theta P \rightsquigarrow \bar{\mathbf{b}}$. This states that for each $\pi_i \in P$, the bind $\mathbf{b}_i : \theta \pi_i$ is provided by Σ , for some substitution θ of the free variables of P . We write $\Sigma \models \theta P$ when we do not care about the elaborated binds. The solving of inferred bind constraints in λPM is complete modulo the completeness of a decision procedure for the entailment relation \models . Note, however, that the type schemes σ for a λPM program are entirely independent of the choice of the theory— Φ constraints never appear in a type scheme σ .

Of course, we require the entailment relation to still define a monad, i.e., \models is admissible if and only if the set $\{\mathbf{b}_{m_1, m_2, m} \mid \Sigma \models \cdot (m_1, m_2) \triangleright m \rightsquigarrow \mathbf{b}\}$ satisfies the monad laws.

3.2 Parameterized monads

Our first example shows that Atkey's parameterized monad [2] is a monad and illustrates how it can be used to program safe communication protocols. Atkey proposes an abstract data type A , a ternary type constructor with two operations, *unitA* and *bindA*, with the signature shown below.

$$\begin{aligned} \text{unitA} &: \forall \alpha \phi. \quad \alpha \rightarrow A \phi \phi \alpha \\ \text{bindA} &: \forall \alpha \beta \phi \gamma \psi. \quad A \phi \gamma \alpha \rightarrow (\alpha \rightarrow A \gamma \psi \beta) \rightarrow A \phi \psi \beta \end{aligned}$$

The type constructor $A p q \tau$ can be thought of (informally) as the type of a computation producing a τ -typed result, with a pre-condition p and a post-condition q . The *bindA* operator matches the post-condition parameter of the first computation with the pre-condition parameter of the composing function, producing a computation having the pre-condition of the former and the post-condition of the latter. The *unitA* operator lifts a pure computation into a parameterized monad with the same pre- and post-condition.

Notice that $A p q$ is not a monad, for all indexes p and q —a unit is only available when the indexes are the same, and the type indices vary in the bind operator. However, Atkey's construction

can be seen as a monad with the signature given below.

$$\begin{aligned} \mathcal{M} &= Id, A/2 \\ \Sigma &= \mathbf{b}_{Id, Id, Id}, \\ &\text{mapA} : \forall \phi, \psi. (A \phi \psi, Id) \triangleright A \phi \psi, \\ &\text{appA} : \forall \phi, \psi. (Id, A \phi \psi) \triangleright A \phi \psi, \\ &\text{unitA} : \forall \phi. (Id, Id) \triangleright A \phi \phi, \\ &\text{bindA} : \forall \phi, \gamma. (A \phi \gamma, A \gamma \psi) \triangleright A \phi \psi \end{aligned}$$

The bind set Σ includes mapA , the functorial map over A (as required by the functorial law). The paired morphism law requires us to include appA (since it is dual to mapA). We include unitA , the analog of Atkey’s *unitA*—note that the monad laws do not require us to provide a unit for every instance of A . Finally, we have bindA , the analog of *bindA*—of course, the varying type constructors are natural with monads.

The composition closure law requires us to close the bind set Σ under the composition of a morphism and a bind. Here, we have one non-trivial morphism, i.e., $\text{unitA} : \forall \phi. Id \succ A \phi \phi$. Composing unitA with bindA we get binds of the form $\forall \phi, \psi. (A \phi \psi, Id) \triangleright A \phi \psi$ and $\forall \phi, \psi. (Id, A \phi \psi) \triangleright A \phi \psi$. These have exactly the same form as mapA and appA , so Σ is already closed. Note, by instantiating the indexes, one can see Σ as an infinite set of binds. Also, observe that in this example the theory constraints Φ are empty; however, as we will see shortly, we still require an entailment relation for bind constraints over the empty theory.

Session types. To check the remaining monad laws, we need to instantiate the abstract type A , provide interpretations for each bind, and check that they satisfy the necessary equations. As an example, we choose Pucella and Tov’s encoding of session types [21]. This provides a way to safely program two-party communication protocols. In what follows, we write *Sess* as a synonym for A , to emphasize the connection to session types.

The type $Sess \phi \gamma \alpha$ represents a computation involved in a two-party session which starts in protocol state ϕ and completes in state γ , returning a value of type α . We begin by instantiating the language of phantom type indices to describe protocol states and transitions. Concretely, we use the type index $send \gamma \alpha$ to denote a protocol state that requires a message of type α to be sent, and then transitions to γ . Similarly, the type index $recv \psi \beta$ denotes the protocol state in which once a message of type β is received, the protocol transitions to ψ . We also use the index end to denote the protocol end state. The signatures of two primitive operations for sending and receiving messages captures this behavior.

$$\begin{aligned} \text{send} &: \forall \alpha, \gamma. \alpha \rightarrow Sess (send \gamma \alpha) \gamma () \\ \text{recv} &: \forall \alpha, \gamma. () \rightarrow Sess (recv \gamma \alpha) \gamma \alpha \end{aligned}$$

The concrete representation of the abstract type $Sess \phi \gamma \alpha$ is simple. Assuming that the underlying language provides support for primitive effects like *I/O*, $Sess \phi \gamma \alpha$ is just a synonym for α . Under this interpretation, all five bind operations correspond to reverse function application—it is easy to check that this interpretation satisfies the left- and right-identities and the associativity law.

Using these definitions, consider the following λ PM program that implements one side of a simple protocol that sends a message x , waits for an integer reply y , and returns $y+1$.

```
let go =  $\lambda x. \mathbf{let} \_ = \text{send } x \mathbf{in} \text{incr} (\text{recv } ())$ 
```

Type inference in λ PM is closely related to the algorithm of Jones’ OML [14], a variant of which is also implemented by Haskell. In Section 4, we show how λ PM programs can be embedded in OML, and how, based on OML’s principal types property, we can compute a principal type of λ PM programs. Using this strategy, we compute the following principal type for go :

$$\begin{aligned} \forall \alpha, \beta, \gamma, \mu_{10}, \mu_2, \mu_8, \mu_9, \mu_5, \mu_3, \mu_6, \mu_7, \mu_6, \mu_{11}, \\ (\mu_2, \mu_8) \triangleright \mu_{10}. \end{aligned}$$

$$\begin{aligned} (\mu_5, Id) \triangleright \mu_9, \\ (Id, Sess (send \alpha \beta) \alpha) \triangleright \mu_3, \\ (Id, \mu_9) \triangleright \mu_8, \\ (Id, \mu_3) \triangleright \mu_2, \\ (Id, Sess (recv \gamma \text{int}) \gamma) \triangleright \mu_6, \\ (Id, \mu_6) \triangleright \mu_5, \\ \Rightarrow (\beta \rightarrow \mu_{10} \text{int}) \end{aligned}$$

Inference also produces a rewritten term that contains, or abstracts, the needed binds. For this example, the term starts with a sequence of λ abstractions, one for each of the seven bind constraints. If we imagine these are numbered $b_1 \dots b_7$ then the main body looks as follows

```
 $\lambda x. b_1 (b_5 \text{send } (\lambda y. b_3 \times y))$   
 $(\lambda \_ . b_4 \text{incr}$   
 $(\lambda z. b_2 (b_7 \text{recv } (\lambda w. b_6 () w)) z))$ 
```

While maximally general, the principal type of go is also unreadable! Worse yet, OML (and Haskell) reject this type as ambiguous. The reason is that the constraint $\text{recv } \beta \rightarrow \mu_{10} \text{int}$. A general-purpose solver for such constraints (not being aware of the monad laws) assumes that the particular instantiation of these variables could influence the semantics of the program, and so requires a programmer to explicitly instantiate each of these variables.

Thankfully, by exploiting the monad laws, λ PM can do much better. In Section 5, we show that for a given result type, all possible solutions to a set of monad constraints are coherent, i.e., they have the same semantics. Next, in Section 6, we show how the monad laws allow us to aggressively improve types, eliminating constraints that we show cannot affect typability. Using our improvement procedure, the type of go becomes slightly more readable (as does the rewritten term, which we elide):

$$\begin{aligned} \forall \alpha, \beta, \gamma, \mu_{10}. \\ (Sess (send \alpha \beta) \alpha, Sess (recv \gamma \text{int}) \gamma) \triangleright \mu_{10}, \\ (Sess (recv \gamma \text{int}) \gamma, Id) \triangleright Sess (recv \gamma \text{int}) \gamma, \\ (Id, Sess (recv \gamma \text{int}) \gamma) \triangleright Sess (recv \gamma \text{int}) \gamma, \\ (Id, Sess (send \alpha \beta) \alpha) \triangleright Sess (send \alpha \beta) \alpha \\ \Rightarrow (\beta \rightarrow \mu_{10} \text{int}) \end{aligned}$$

This type is still rather unwieldy. But, note that all but the first constraint are tautologies. By the functorial and paired morphism laws, we know that for every instantiation of α, β, γ , there is guaranteed to be a bind in Σ of the form $(m, Id) \triangleright m$ and $(Id, m) \triangleright m$. Thus, we can further simplify the type to the one shown below.

$$\begin{aligned} \forall \alpha, \beta, \gamma, \mu_{10}. \\ (Sess (send \alpha \beta) \alpha), (Sess (recv \gamma \text{int}) \gamma) \triangleright \mu_{10} \\ \Rightarrow (\beta \rightarrow \mu_{10} \text{int}) \end{aligned}$$

Finally, when at the top-level a programmer calls $go\ 0$, and instantiates the result type, say, to $Sess\ end\ end$, we obtain a constraint $\pi = (Sess (send \alpha \beta) \alpha), (Sess (recv \gamma \text{int}) \gamma) \triangleright Sess\ end\ end$. To complete typing $go\ 0$, we make use of the decision procedure \models to solve this constraint. For this particular example, without any theory constraints, the relation \models is a simple unification-based procedure that can compute the substitution $\theta = \alpha \mapsto (recv\ end\ \text{int}), \gamma \mapsto end$, such that $\text{bindA} \in \Sigma$ can be instantiated to $\theta\pi$.

3.3 Polymorphic information flow controls

Several researchers have proposed type systems or libraries with a monad-like structure to implement information flow controls [7, 23, 17, 5, 1]. These controls allow a programmer to indicate that some program inputs are secrets, and that some outputs are public. The goal is to ensure that the public outputs are independent of the secret inputs—a property called noninterference [12]. In this section, we develop a monad, *IST*, suitable for enforcing information flow controls in stateful programs.

Our encoding has several elements. First, as usual, we introduce a lattice of security labels $l \in \{L, H\}$, each a nullary type constructor. We assume $L < H$ in a lattice ordering, indicating that data labeled H is more secret than data labeled L .

Next, we introduce a type of integer references,¹ $\text{intref } l$, a reference to an integer value labeled l . Finally, we have a polymonadic type constructor IST which takes two phantom type indices. In particular, $IST \text{ pc } l \tau$ classifies a computation that potentially writes to references labeled l and returns a τ -result that is labeled l . For instance, $IST \text{ H } L$ is the type of a computation that may write to secret storage cells while computing a public value.

To make this notion precise, we introduce a lattice theory Φ and use it in the definition of the bind set below.

$$\begin{aligned} \mathcal{M} &= IST/2 \\ \Phi &::= \tau \leq \tau \mid \Phi, \Phi \\ \Sigma &= \mathbf{b}_{Id, Id, Id}, \\ &\text{app}IST : \forall pc, l. (Id, IST \text{ pc } l) \triangleright IST \text{ pc } l, \\ &\text{map}IST : \forall pc, l. (IST \text{ pc } l, Id) \triangleright IST \text{ pc } l, \\ &\text{b}IST : \forall pc_1, l_1, pc_2, l_2, pc_3, l_3. \\ &\quad l_1 \leq pc_2, l_1 \leq l_3, l_2 \leq l_3, pc_3 \leq pc_1, pc_3 \leq pc_2 \\ &\quad \Rightarrow (IST \text{ pc}_1 l_1, IST \text{ pc}_2 l_2) \triangleright IST \text{ pc}_3 l_3 \end{aligned}$$

When composing a computation $IST \text{ pc}_1 l_1 \alpha$ with a function $\alpha \rightarrow IST \text{ pc}_2 l_2 \beta$, the type of $\text{b}IST$ requires $l_1 \leq pc_2$ to prevent the second computation from leaking information about its l_1 -secure α -typed argument into a reference cell that is less than l_1 -secure. Dually, the next two constraints ensures that the β -typed result of the composed computation is at least as secure as the results of each component. The last two constraints ensure that the effect lower bound of the composed computation is a lower bound of the effects of each component.

To interpret the theory constraints, we instantiate \models to a standard theory of lattice constraints. Specifying this decision procedure is orthogonal to the purposes of this paper. However, for a finite lattice of labels, an easy decision procedure simply constructs the finite enumeration of binds obtained by instantiating $\text{b}IST$ in all ways consistent with the lattice constraints Φ . Deciding whether a particular bind constraint π is derivable then simply amounts to testing the membership of π in the finite enumeration, for some substitution θ of the free variables in π .

To implement IST , we must give it a representation in System F and provide an interpretation for $\text{b}IST$ and check that it satisfies the polymonad laws. One choice is for IST to be the state monad augmented with phantom label indexes, as shown below.

```
heap = list (int * int) (* an assoc. list for a heap *)
intref l = int (* intref l is a synonym of int *)
IST  $\alpha\beta\gamma$  = heap  $\rightarrow$  ( $\gamma$  * heap) (* state monad with phantom indexes *)
bIST =  $\lambda c \text{ f h} = \mathbf{let} (x, h) = c \text{ h in f } x$ 
```

We can then give polymonadic signatures to functions for allocating, reading, and writing references, as shown below.

```
alloc   :  $\forall pc. \text{int} \rightarrow IST \text{ pc } pc$  (intref pc)
read    :  $\forall l. \text{intref } l \rightarrow IST \text{ H } l \text{ int}$ 
write   :  $\forall l. \text{intref } l \rightarrow \text{int} \rightarrow IST \text{ l } L$  ()
```

With these definitions in hand, we can write the following λ PM program, infer a principal type for it (very verbose, as with sessions), and improve its type using the polymonad laws to the one

¹ A generalization to polymorphic references is feasible, but we use integer references here to keep our heap model simple.

shown below.

```
 $\lambda$ PM source   let add_interest =  $\lambda$ savings.  $\lambda$ interest.
                let current = read savings in
                write savings (current + interest)
improved type  $\forall l, \mu. (IST \text{ H } l, IST \text{ l } L) \triangleright \mu$ 
                 $\Rightarrow$  intref l  $\rightarrow$  int  $\rightarrow \mu$  ()
```

Consider applying `add_interest secret_ac 100`, where `secret_ac` : `intref l` is a reference to a bank account at a secrecy level l , for some l . We now have a constraint $(IST \text{ H } l, IST \text{ l } L) \triangleright \mu$ to solve. If we instantiate the variable μ to $IST \text{ l } l$ which produces type $\forall l. \text{intref } l \rightarrow \text{int} \rightarrow IST \text{ l } l$ (), we must check that $(IST \text{ H } l, IST \text{ l } L) \triangleright IST \text{ l } l$ exists, and it does: the theory constraints on $\text{b}IST$ are satisfied since the resulting pc label $l \leq l$ and $l \leq H$; and, likewise, the resulting confidentiality label $l \geq L$ and $l \geq l$. Intuitively, the type of `add_interest` tells us that it is a function that writes to its argument reference and returns a value that is at least as secret as the contents of that reference. This type is precise as far as its effect lower bound goes, but it is a little imprecise in that its result value is (), hence uninformative, and so could be given the label L . Overall, polymonadic information flow tracking provides a convenient syntactic way of eliminating information leaks in a program, but, being syntactic, it is necessarily imprecise.

3.4 Contextual type and effect systems

We have already sketched an encoding of Wadler and Thiemann's [27] type and effect systems as polymonad. As our final example, we show that a recent type and effect system for *contextual effects* [20] (which subsumes traditional type and effects) is a polymonad. We define a polymonad $CE \epsilon \alpha \omega \tau$, for the type of a computation which itself has effect ϵ and produces a value of type τ , and appears in a context in which the prior computations have effect α and whose subsequent computations have effect ω .

As with our information flow encoding, we start by describing a language of type indices to describe effect sets. Indices are sets of atomic effects $a_1 \dots a_n$, with \emptyset the empty effect, \top the effect set that includes all other effects, and \cup the union of two effects. We also introduce theory constraints for subset relations and extensional equality on sets, with the obvious interpretation.

```
types           $\tau ::= \dots \mid a_1 \dots a_n \mid \emptyset \mid \top \mid \tau_1 \cup \tau_2$ 
theory constraints  $\Phi ::= \tau \subseteq \tau' \mid \tau = \tau' \mid \Phi, \Phi$ 
```

The following binds capture the tracking of contextual effects:

$$\begin{aligned} \mathcal{M} &= Id, CE/3 \\ \Sigma &= \mathbf{b}_{Id, Id, Id}, \\ &\text{unitce} : (Id, Id) \triangleright CE \top \emptyset \top \\ &\text{bindce}_{Id} : \forall \alpha_1, \alpha_2, \epsilon_1, \epsilon_2, \omega_1, \omega_2. \\ &\quad (\alpha_2 \subseteq \alpha_1, \epsilon_1 \subseteq \epsilon_2, \omega_2 \subseteq \omega_1) \Rightarrow \\ &\quad (Id, CE \alpha_1 \epsilon_1 \omega_1) \triangleright CE \alpha_2 \epsilon_2 \omega_2 \\ &\text{bindce} : \forall \alpha_1, \epsilon_1, \omega_1, \alpha_2 \epsilon_2, \omega_2, \epsilon_3. \\ &\quad (\epsilon_2 \cup \omega_2 = \omega_1, \epsilon_1 \cup \alpha_1 = \alpha_2, \epsilon_1 \cup \epsilon_2 = \epsilon_3) \Rightarrow \\ &\quad (CE \alpha_1 \epsilon_1 \omega_1, CE \alpha_2 \epsilon_2 \omega_2) \triangleright CE \alpha_1 \epsilon_3 \omega_2 \end{aligned}$$

The bind `unitce` lifts a computation into a contextual effect monad with empty effect and any prior or future effects. The bind `bindceId` expresses that it is safe to consider an additional effect for the current computation (the ϵ s are covariant), and fewer effects for the prior and future computations (α s and ω s are contravariant). Finally, `bindce` composes two computations such that the future effect of the first computation includes the effect of the second one, provided that the prior effect of the second computation includes the first computation; the effect of the composition includes both effects, while the prior effect is the same as before the first

computation, and the future effect is the same as after the second computation.

As with Atkey’s monad, it is easy to check that Σ is closed under the composition with the single non-trivial morphism bind_{Id} . However, unlike our other examples, we do not provide an example of programming with contextual effects. Typing any non-trivial example requires an implementation of a decision procedure for the set theory we have introduced—something that our current prototype lacks. In the future, we aim to extend our work to the setting of an SMT-based dependently typed programming language, where we hope to make use of the theory support of the underlying solver to efficiently decide polymonadic theory constraints.

4. Type inference for λPM

This section presents a formalization of type inference for λPM . We defer to the next section how type inference is (also) drives the conversion of the direct-style source program to an elaborated program containing the needed polymonadic binds. We prove that our type system enjoys principal types via a sound and complete translation to Jones’ OML [14] (the theoretical foundation of Haskell type classes), that type inference produces principal types. However, this is not the end of the story—the next two sections show that appealing to the polymonadic laws makes it possible to safely accept what might otherwise be deemed ambiguous programs, and to simplify principal types without reducing their generality.

4.1 Type rules

Figure 4 gives the syntax-directed type rules of our language. As shown in earlier work [24] it is straightforward to give declarative rules too. There are two forms of rules, one for values, $P \mid \Gamma \vdash v : \tau$ and one for expressions $P \mid \Gamma \vdash e : m \tau$ (which have a polymonadic type). The rules state that under an environment Γ , and predicates P , the value v or expression e have type τ or $m \tau$ respectively.

The rules are all straightforward: Rules (TS-Var) and (TS-Const) look up a variable or constant in the environment and return an instantiated type scheme. Note how the rule (TS-Inst) ensures that the constraints in the type scheme are entailed by the assumed predicates P . The entailment relation is defined as:

$$\frac{\pi \in P}{P \models \pi} \quad \frac{\Sigma \models \cdot \pi}{P \models \pi} \quad \frac{P \models \pi_1 \quad \dots \quad P \models \pi_n}{P \models \pi_1, \dots, \pi_n}$$

Basically, a predicate π must be implied either by the bind set Σ , or be an element of the assumed predicates P .

The rule (TS-Lam) types lambda abstractions while (TS-Id) lifts a value into the identity monad. Rule (TS-Let) generalizes over values in the usual way.

Rule (TS-App) ensures that there exists a bind expression that can apply the function (with monadic type m_3) to the argument (m_2) taking it to some monadic type m_4 , and that there exist a bind that can combine the evaluation of the function (of type m_1) with that result monad (m_4) into a final monadic type m_5 . Rule (TS-Do) is similar, but requires only the existence of a bind between the bound expression and the body of the let.

4.2 Principal types

The type rules admit principal types, and there exists an efficient type inference algorithm that finds such types. The way we show this is by a translation of polymonadic terms (and types) to terms (and types) in OML [14] and prove this translation is sound and complete: a polymonadic term is well-typed if and only if its translated OML term has an equivalent type.

We encode terms in our language into OML as shown in Figure 5. We rely on three primitive OML terms that force the typing

$$\begin{aligned} \text{id} & : \forall \alpha. \alpha \rightarrow \text{Id } \alpha \\ \text{do} & : \forall \alpha \beta \mu_1 \mu_2 \mu. ((\mu_1, \mu_2) \triangleright \mu) \\ & \quad \Rightarrow \mu_1 \alpha \rightarrow (\alpha \rightarrow \mu_2 \beta) \rightarrow \mu \beta \\ \text{app} & : \forall \alpha \beta \mu_1 \mu_2 \mu_3 \mu_4 \mu. ((\mu_1, \mu_4) \triangleright \mu, (\mu_2, \mu_3) \triangleright \mu_4) \\ & \quad \Rightarrow \mu_1 (\alpha \rightarrow \mu_3 \beta) \rightarrow \mu_2 \alpha \rightarrow \mu \beta \\ \\ \llbracket x \rrbracket^* & = x \\ \llbracket c \rrbracket^* & = c \\ \llbracket \lambda x. e \rrbracket^* & = \lambda x. \llbracket e \rrbracket \\ \\ \llbracket v \rrbracket & = \text{id } \llbracket v \rrbracket^* \\ \llbracket e_1 e_2 \rrbracket & = \text{app } \llbracket e_1 \rrbracket \llbracket e_2 \rrbracket \\ \llbracket \text{let } x = v \text{ in } e \rrbracket & = \text{let } x = \llbracket v \rrbracket^* \text{ in } \llbracket e \rrbracket \\ \llbracket \text{let } x = e_1 \text{ in } e_2 \rrbracket & = \text{do } \llbracket e_1 \rrbracket \llbracket \lambda x. e_2 \rrbracket^* \quad (\text{with } e_1 \neq v) \end{aligned}$$

Figure 5. Translation of λPM to OML

of the terms to generate the same constraints as our type system does: id for lifting a pure term, do for typing a do-binding, and app for typing an application. Using these primitives, we encode values and expressions of our system into OML.

We write $P \mid \Gamma \vdash_{\text{OML}} e : \tau$ for a derivation in the syntax directed inference system of OML (cf. Jones [14], Fig. 4).

Theorem 2 (Elaboration to OML is sound and complete).

Soundness: *Whenever $P \mid \Gamma \vdash v : \sigma$ we can also derive $P \mid \Gamma \vdash_{\text{OML}} \llbracket v \rrbracket^* : \sigma$ in OML. Similarly, when $P \mid \Gamma \vdash e : m \tau$ we have $P \mid \Gamma \vdash_{\text{OML}} \llbracket e \rrbracket : m \tau$.*

Completeness: *If we can derive $P \mid \Gamma \vdash_{\text{OML}} \llbracket v \rrbracket^* : \sigma$, there also exists a derivation $P \mid \Gamma \vdash v : \sigma$, and similarly, whenever $P \mid \Gamma \vdash_{\text{OML}} \llbracket e \rrbracket : m \tau$, we also have $P \mid \Gamma \vdash e : m \tau$.*

The proof is by straightforward induction on the typing derivation of the term. It is important to note that our system uses the same instantiation and generalization relations as OML which is required for the induction argument. Moreover, the constraint entailment over bind constraints, also satisfies the monotonicity, transitivity and closure under substitution properties required by OML.

As a corollary of the above properties, we have that our system admits principal types via the general-purpose OML type inference algorithm.

5. Coherence

By the result in the previous section, we could perform type inference and rewriting for polymonads by using OML’s algorithm, which is essentially the Haskell type class inference algorithm. Unfortunately, this algorithm is not satisfactory, as it would reject many useful programs. Translated to our setting, the limitation of this algorithm is that it rejects any term whose type $\forall \bar{v}. P \Rightarrow \tau$ binds a variable $\mu \in \bar{v}$ such that μ appears in a constraint $\pi \in P$ but does not appear in the final type τ . We call such variables *open variables* and such constraints *open constraints*. Many polymonadic terms have open constraints; one example was given in Section 3.2.

The reason such terms are rejected in OML is that instantiations of open variables have operational effect—the instantiations determine which *evidence* will be used when evaluating a term (in our setting, the evidence is the binds), and different evidence could have different operational effect.

However, it turns out for polymonads we do not have the same problem: all possible solutions to open constraints will produce terms having the same semantics; i.e., the solutions are *coherent*. To show coherence, we have to show that the programs resulting from evidence translation (i.e. how we insert binds) are operationally equivalent no matter which solution we choose for open

$$\boxed{P \mid \Gamma \vdash v : \tau \quad P \mid \Gamma \vdash e : m \tau}$$

$$\frac{\Gamma(x) = \sigma \quad P \models \sigma \geq \tau}{P \mid \Gamma \vdash x : \tau} \text{ (TS-Var)} \quad \frac{\Gamma(c) = \sigma \quad P \models \sigma \geq \tau}{P \mid \Gamma \vdash c : \tau} \text{ (TS-Const)} \quad \frac{\theta = [\bar{m}/\bar{\mu}][\bar{\tau}/\bar{\alpha}] \quad P \models \theta P_1}{P \models \forall \bar{\mu} \bar{\alpha}. P_1 \Rightarrow \tau \geq \theta \tau} \text{ (TS-Inst)}$$

$$\frac{P \mid \Gamma, x : \tau_1 \vdash e : m \tau_2}{P \mid \Gamma \vdash \lambda x. e : \tau_1 \rightarrow m \tau_2} \text{ (TS-Lam)} \quad \frac{P \mid \Gamma \vdash e_1 : m_1 (\tau_2 \rightarrow m_3 \tau) \quad P \mid \Gamma \vdash e_2 : m_2 \tau_2}{P \models (m_1, m_4) \triangleright m_5 \quad P \models (m_2, m_3) \triangleright m_4} \text{ (TS-App)}$$

$$\frac{P \mid \Gamma \vdash v : \tau}{P \mid \Gamma \vdash v : Id \tau} \text{ (TS-Id)} \quad \frac{P \mid \Gamma \vdash e_1 : m_1 \tau_1 \quad P \mid \Gamma, x : \tau_1 \vdash e_2 : m_2 \tau_2 \quad P \models (m_1, m_2) \triangleright m_3}{P \mid \Gamma \vdash \text{let } x = e_1 \text{ in } e_2 : m_3 \tau_2} \text{ (TS-Do)} \quad \frac{P' \mid \Gamma \vdash v : \tau' \quad \sigma = Gen(\Gamma, P' \Rightarrow \tau') \quad P \mid \Gamma, x : \sigma \vdash e : m \tau}{P \mid \Gamma \vdash \text{let } x = v \text{ in } e : m \tau} \text{ (TS-Let)}$$

Figure 4. Syntax-directed type rules for λ PM. The generalization function is defined as: $Gen(\Gamma, \sigma) = \forall(\text{ftv}(\sigma) \setminus \text{ftv}(\Gamma)).\sigma$

constraints. Unfortunately, doing such a proof directly on the evidence translation for the syntax directed system is difficult: since the rules allow arbitrary instantiations for monadic constraints it is hard to relate the resulting evidence terms semantically. Instead we define a more algorithmic version of the type rules that let us reason about coherence in a more structural and syntactic way.

The remainder of this section presents our proof. First we present a more algorithmic treatment of type inference that shows how source terms are elaborated to target terms containing binds. Next, we define two additional properties needed for coherence, and argue why they are easy to satisfy. We conclude with the presentation of the actual proof.

5.1 Type inference with elaboration

Figure 6 defines a more algorithmic version of the syntax directed type rules of Figure 4. The inference rules are $P \mid \Gamma \vdash e : \tau \rightsquigarrow e$ and $P \mid \Gamma \vdash e : m \tau \rightsquigarrow e$ where the constraints P , type τ (or $m \tau$) and target term e are synthesized. This definition combines a syntax-directed presentation of Hindley-Milner type inference, where the details of unification are elided, with an algorithmic presentation of polymonadic constraint generation and solving.

In particular, in rule (TA-Var) and (TA-Const), the substitutions for the type variables are still chosen arbitrarily as in any standard presentation of Hindley-Milner style type rules. However, the monadic variables are explicitly substituted with fresh monadic variables where the constraints of the type scheme are returned in the predicates P . Similarly, in rule (TA-App) and (TA-Do), the intermediate monadic types are now represented by fresh monadic type variables. Effectively, this ensures that the evidence for each bind is always inserted and solved at specific syntactical locations which is essential to doing the coherence proof. Finally, (TA-Let) performs generalization. It uses a relation $P \xrightarrow{\text{simplify}(\bar{\mu})} P'$; θ to optionally simplify constraints P by using substitution θ to solve (some) monad variables $\bar{\mu}$ with any remaining, unsolved constraints in P' . Solved variables are no longer generalized. We present our particular simplification algorithm in Section 6.

We have not done a full proof yet, but we conjecture that the algorithmic rules of Figure 6 are sound and complete with respect to the syntax directed rules of Figure 4.

The following lemma establishes that our algorithm produces well-typed System F terms. By $\{\{\Gamma\}\}$ we mean the (straightforward) translation of a source typing context to a System F context.

Lemma 3 (Well-typed elaborations). *Given e, τ, m, e, Γ such that if either $P \mid \Gamma \vdash e : \tau \rightsquigarrow e$ or $P \mid \Gamma \vdash e : m \tau \rightsquigarrow e$ then there exists t such that $\{\{\Gamma\}\} \vdash_F \text{abs}(P, e) : t$.*

Proof. See Appendix B. \square

Once we prove $P \mid \Gamma \vdash e : m \tau \rightsquigarrow e$ and instantiate the final type $m \tau$, we must solve the constraints P to produce bind terms \bar{b} that we use to execute e . We can do this using the given polymonadic entailment relation $\Sigma \models \theta P \rightsquigarrow \bar{b}$ based on Σ 's theory Φ . Then we simply execute $e \bar{b}$. We discuss our implementation of solving in Section 7.

5.2 Additional properties

With the evidence translation in place, we can now precisely state what we mean by coherence: suppose we have $P \mid \Gamma \vdash e : m \tau \rightsquigarrow e$. Then if θ and θ' are both solutions to the constraints P , where a solution is a mapping of constraints π to particular binds $\bar{b}_{m_1, m_2, m_3} \in \Sigma$, then both solutions, applied to e , will yield terms with an identical operational effect. Our proof of coherence works by starting with θ and showing how we can iteratively transform it into θ' , all the while proving that each intermediate step does not affect the semantics of the rewritten term. For this proof to work we rely on polymonads satisfying two additional properties.

First, the semantics of a particular bind in Σ must be independent of any type indexes in m ; i.e., type indices are *operationally irrelevant*. That is, as mentioned in Section 3.1, for computations $M \bar{\tau}$ the type indices $\bar{\tau}$ are meant to be phantom, and thus not influence a bind's semantics. To express this idea formally, let (e) represent the type *erasure* of a System F term, so that all type annotations, type abstractions, and type applications are dropped. Then we expect polymonads to satisfy the following property

Param: For all $M_1, M_2, M_3, \bar{\tau}_1, \bar{\tau}_2, \bar{\tau}_3, \bar{\tau}'_1, \bar{\tau}'_2, \bar{\tau}'_3$,

$$(\text{bind}_{M_1 \bar{\tau}_1, M_2 \bar{\tau}_2, M_3 \bar{\tau}_3} e k) \cong (\text{bind}_{M_1 \bar{\tau}'_1, M_2 \bar{\tau}'_2, M_3 \bar{\tau}'_3} e k)$$

Here, we write \cong to mean β/η equality for the untyped lambda calculus (which should match the semantics of the original System F terms). The (Param) property is useful for establishing that intermediate solutions will not change a term's semantics when only the type indices are changing. All of our examples from Section 3 satisfy this property.

Second, we require that the polymonad be *smooth* in that for each pair of constructors $M_1, M_2 \in \mathcal{M}$, we can compute $M = \text{smooth}(M_1, M_2)$ where there exist binds involving M for those binds involving M_1 and/or M_2 :

Definition 4 (Smooth polymonad). *A smooth polymonad $\Sigma_{\mathcal{M}}$ is one for which we can define a function smooth on polymonad constructors $M \in \mathcal{M}$, written $\text{smooth}(M, M') = M^*$. This function has the following properties.*

$P \mid \Gamma \vdash e : m \tau \rightsquigarrow e \quad P \mid \Gamma \vdash v : \tau \rightsquigarrow v$	
$\frac{\Gamma(x) = \forall \bar{\mu}, \bar{\alpha}. P \Rightarrow \tau \quad \theta = [\bar{\mu}'/\bar{\mu}][\bar{\tau}/\bar{\alpha}] \quad \bar{\mu}' \text{ fresh}}{\theta P \mid \Gamma \vdash x : \theta \tau \rightsquigarrow \mathbf{app}(x, \theta P)} \quad (\text{TA-Var})$	$\frac{\Gamma(c) = \forall \bar{\mu}, \bar{\alpha}. P \Rightarrow \tau \quad \theta = [\bar{\mu}'/\bar{\mu}][\bar{\tau}/\bar{\alpha}] \quad \bar{\mu}' \text{ fresh}}{\theta P \mid \Gamma \vdash c : \theta \tau \rightsquigarrow \mathbf{app}(c, \theta P)} \quad (\text{TA-Const})$
$\frac{P \mid \Gamma, x:\tau_1 \vdash e : m \tau_2 \rightsquigarrow e}{P \mid \Gamma \vdash \lambda x. e : \tau_1 \rightarrow m \tau_2 \rightsquigarrow \lambda x:\tau_1. e} \quad (\text{TA-Lam})$	$\frac{P_1 \mid \Gamma \vdash e_1 : m_1 (\tau_2 \rightarrow m_3 \tau) \rightsquigarrow e_1 \quad P_2 \mid \Gamma \vdash e_2 : m_2 \tau_2 \rightsquigarrow e_2 \quad \mu_4, \mu_5 \text{ fresh} \quad P = P_1, P_2, (m_1, \mu_4) \triangleright \mu_5, (m_2, m_3) \triangleright \mu_4}{P \mid \Gamma \vdash e_1 e_2 : \mu_5 \tau \rightsquigarrow \mathbf{b}_{m_1, \mu_4, \mu_5} e_1 (\lambda x:\dots \mathbf{b}_{m_2, m_3, \mu_5} e_2 x)} \quad (\text{TA-App})$
$\frac{P \mid \Gamma \vdash v : \tau \rightsquigarrow v}{P \mid \Gamma \vdash v : \text{Id } \tau \rightsquigarrow (\mathbf{unit}_{\text{Id}} v)} \quad (\text{TA-Id})$	$\frac{e_1 \neq v \quad P_1 \mid \Gamma \vdash e_1 : m_1 \tau_1 \rightsquigarrow e_1 \quad P_2 \mid \Gamma, x:\tau_1 \vdash e_2 : m_2 \tau_2 \rightsquigarrow e_2 \quad \mu_3 \text{ fresh} \quad P = P_1, P_2, (m_1, m_2) \triangleright \mu_3}{P \mid \Gamma \vdash \text{let } x = e_1 \text{ in } e_2 : \mu_3 \tau_2 \rightsquigarrow \mathbf{b}_{m_1, m_2, \mu_3} e_1 (\lambda x:\dots e_2)} \quad (\text{TA-Do})$
$\frac{P' \mid \Gamma \vdash v : \tau' \rightsquigarrow v \quad \bar{\mu}, \bar{\alpha} = \bar{v} = \text{ftv}(P' \Rightarrow \tau') \setminus \text{ftv}(\Gamma) \quad P' \xrightarrow{\text{simplify}(\bar{\mu} \setminus \text{ftv}(\tau'))} P''; \theta \quad \bar{v}' = \bar{v} \setminus \text{dom}(\theta) \quad \sigma = \forall \bar{v}'. P'' \Rightarrow \tau' \quad P \mid \Gamma, x:\sigma \vdash e : m \tau \rightsquigarrow e}{P \mid \Gamma \vdash \text{let } x = v \text{ in } e : m \tau \rightsquigarrow (\lambda x:\dots e) \text{abs}(P'', \theta v)} \quad (\text{TA-Let})$	
<div style="display: flex; justify-content: space-between;"> <div style="width: 45%;"> <p>where</p> $\begin{aligned} \mathbf{app}(e, (P, (m_1, m_2) \triangleright m_3)) &= \mathbf{app}(e, P) \mathbf{b}_{m_1, m_2, m_3} \\ \mathbf{app}(e, \cdot) &= e \\ \mathbf{abs}(((m_1, m_2) \triangleright m, P), e) &= \lambda \mathbf{b}_{m_1, m_2, m:\dots} \mathbf{abs}(P, e) \\ \mathbf{abs}(\cdot, e) &= e \end{aligned}$ </div> </div>	

Figure 6. Algorithmic syntax-directed type rules with evidence translation

1. Given $M_1^* = \text{smooth}(M_1, M_1')$ and $M_2^* = \text{smooth}(M_2, M_2')$.
if $\Sigma \models (M_1 _ , M_2 _) \triangleright M_3 _ \text{ and } \Sigma \models (M_1' _ , M_2' _) \triangleright M_3 _ \text{ then there exist type indices } \bar{\tau}_1, \bar{\tau}_2, \bar{\tau}_3, \bar{\tau}'_1, \bar{\tau}'_2, \bar{\tau}'_3 \text{ such that } \Sigma \models (M_1^* \bar{\tau}_1, M_2^* \bar{\tau}_2) \triangleright M_3 \bar{\tau}_3 \text{ and } \Sigma \models (M_1^* \bar{\tau}'_1, M_2^* \bar{\tau}'_2) \triangleright M_3 \bar{\tau}'_3$
2. Given $M^* = \text{smooth}(M, M')$.
If $\Sigma \models (M_1 _ , M_2 _) \triangleright M _ \text{ or } \Sigma \models (M_1 _ , M_2 _) \triangleright M' _ \text{ then there exist type indices } \bar{\tau}_1, \bar{\tau}_2, \bar{\tau}' \text{ such that } \Sigma \models (M_1 \bar{\tau}_1, M_2 \bar{\tau}_2) \triangleright M^* \bar{\tau}'$.
3. $\text{smooth}(M, M) = M$.

Intuitively, $\text{smooth}(M_1, M_2)$ is like the least upper bound of M_1 and M_2 . This property is useful for generating small modifications to a solution θ to bring it closer, in a coherent fashion, to solution θ' .

We believe that requiring smooth places little practical limitation on polymonads. In particular, when writing programs with just one polymonadic constructor M , the function is trivial: $\text{smooth}(M, M) = M$, $\text{smooth}(\text{Id}, M) = M$, $\text{smooth}(M, \text{Id}) = M$, and $\text{smooth}(\text{Id}, \text{Id}) = \text{Id}$. When using multiple constructors M_1 and M_2 , we must already define how computations interact. For example, suppose we were programming with sessions and state, with constructors $\text{Sess}/2$ and $\text{ST}/0$. We would probably define a third constructor $\text{SessST}/2$ that represents a communicating computation that also modifies the heap, and we would define a bind for lifting $\text{Sess } \alpha \beta \tau$ computations into $\text{SessST } \alpha \beta \tau$ computations, and one for lifting $\text{ST } \tau$ computations into $\text{SessST } \alpha \beta \tau$ computations. By the (Param) property, these liftings will not depend on type indices, and so they are morally monad-style morphisms $\text{Sess} \succ \text{SessST}$ and $\text{ST} \succ \text{SessST}$. Such morphisms, when arranged to form a lattice, easily satisfy the requirements of smooth .

5.3 Proof of coherence

Now we present the final formal details of the proof of coherence. First, some notation: We view constraints P as a directed graph. Nodes are polymonads m , i.e., either monad type constants $M \bar{\tau}$

$\text{up-bnd}((\cdot, \cdot) \triangleright m)$	$= m$
$\text{up-bnds}(P)$	$= \bigcup_{\pi \in P} \{\text{up-bnd}(\pi)\}$
$\text{lo-bnd}((m_1, m_2) \triangleright \cdot)$	$= (m_1, m_2)$
$\text{lo-bnds}(P)$	$= \bigcup_{\pi \in P} \{\text{lo-bnd}(\pi)\}$
$\text{flowsTo}_P \mu$	$= \{\pi \mid \pi \in P \wedge \mu = \text{up-bnd}(\pi)\}$
$\text{flowsFrom}_P \mu$	$= \{\pi \mid \pi \in P \wedge \exists m'. \text{lo-bnd}(\pi) = (\mu, m') \vee \text{lo-bnd}(\pi) = (m', m\mu)\}$

Figure 7. Constraints P as graphs

(where the $\bar{\tau}$ could contain type variables), or monad variables μ . Each bind constraint $(m_1, m_2) \triangleright m$ induces two edges, a *left edge* $m_1 \rightarrow m$ and a *right edge* $m_2 \rightarrow m$. We say the *upper bound* of such a constraint is m , and the lower bound is the pair of monads (m_1, m_2) ; we can think of the former as the constraint's *output* and the latter as its *input*. These notions are defined formally in Figure 7, along with their obvious liftings to constraint sets P . The figure also defines $\text{flowsTo}_P \mu$ —the set of constraints in P that have μ as an upper bound—and $\text{flowsFrom}_P \mu$ —the set of constraints that have μ as a lower bound.

For a given constructor $M/k \in \mathcal{M}$ not all instantiations of M 's type indices may be legal in a given bind. For example, for session types we could not legally define a bind with type $(A \phi \gamma, A \phi \psi) \triangleright A \phi \psi$ if ϕ and γ were different. However, because of (Param), we view the semantics of any bind as due to the default “operational” bind for every triple of polymonad constructors, irrespective of their parameters.

For the sole purpose of coherence proof, which is concerned only with run-time semantics, it is safe to assume the existence of more binds, which only differ in polymonad parameters from the existing binds. We thus *saturate* the initial signature with more binds, provided that they never appear in actual solutions but only used to establish congruence of the actual solutions by transitivity. Saturation is defined formally as follows:

Definition 5 (Saturated signature). Given a bindset Σ , its saturated signature $\Sigma_{sat} = \langle \Sigma \rangle$, is defined as follows:

$$\begin{aligned} \langle \mathbf{b}_{Id, Id, Id} \rangle &= \mathbf{b}_{Id, Id, Id} \\ \langle \mathbf{b} : \forall \bar{\alpha}. \Phi \Rightarrow (M_1 \bar{\tau}_1, M_2 \bar{\tau}_2) \triangleright M_3 \bar{\tau}_3, \Sigma \rangle &= \\ \mathbf{b} : \forall \bar{\beta}_1 \bar{\beta}_2 \bar{\beta}_3. (M_1 \bar{\beta}_1, M_2 \bar{\beta}_2) \triangleright M_3 \bar{\beta}_3, \langle \Sigma \rangle \end{aligned}$$

Finally, following two more definitions, here is our general coherence result for polymonads.

Definition 6 (Well-formedness of a signature). A well-formed signature Σ is one that satisfies the polymonad laws, the (Param) property, and is smooth.

Definition 7 (Ground solution). A solution θ to constraints P is ground for Σ if and only if $co\text{-domain}(\theta)$ contains only ground polymonads M or ground types τ , and $\Sigma \models \theta\pi$ for all $\pi \in P$.

Theorem 8 (Coherence).

Given $\Sigma, \Gamma, P, e, m, \tau, \theta_1, \theta_2, e, \bar{\mu}, \bar{\alpha}$, such that

1. Σ is well-formed and P is cycle-free.
2. $P \mid \Gamma \vdash e : \tau \rightsquigarrow e$ or $P \mid \Gamma \vdash e : m\tau \rightsquigarrow e$.
3. There exist open variables $\bar{\mu}, \bar{\alpha} = ftv(P) \setminus ftv(\tau)$ (or $\bar{\mu}, \bar{\alpha} = ftv(P) \setminus ftv(m\tau)$) such that for all $\mu \in \bar{\mu}$, the sets $flowsTo_P \mu$ and $flowsFrom_P \mu$ are non-empty, i.e., each $\mu \in \bar{\mu}$ has a lower and an upper bound.
4. θ_1 and θ_2 are ground solutions for P such that $\theta_1(\nu) = \theta_2(\nu)$ for all $\nu \notin \bar{\alpha}, \bar{\mu}$.

Then, $\theta_1 e = \theta_2 e$.

Proof. The full proof is shown in Appendix C. Because θ_1 and θ_2 are ground solutions for P under Σ , they are also ground under the saturated signature Σ_{sat} . The proof first constructs a solution θ_3 by combining monad variable substitutions from θ_2 and type variable substitutions from θ_1 ; this is a ground solution under Σ_{sat} . By repeated appeal to the (Param) property for every subterm that differs between $\theta_2 e$ and $\theta_3 e$ we get that $\theta_2 e = \theta_3 e$.

The rest of the proof proceeds by iterating over the constraints solved differently by θ_1 and θ_3 . Because P is cycle-free we can consider each $\pi \in P$ in reverse topological order. We maintain an invariant that each π considered has a (ground) upper bound. At each step we construct two new ground solutions θ'_1 and θ'_3 that only differ from θ_1 and θ_3 in the substitution of the lower bounds for the current constraint π . We assign by $smooth(lo\text{-bnds}(\theta_1\pi), lo\text{-bnds}(\theta_3\pi))$ to these bounds in both θ'_1 and θ'_3 . We prove these are ground solutions due to the properties of $smooth$, and they must have the same semantics as the terms with θ_1 and θ_3 applied, which follows from a corollary of the polymonadic associativity property. We set $\theta_1 = \theta'_1$ and $\theta_3 = \theta'_3$ and continue. At the last step the solutions are exactly the same, so $\theta_1 e = \theta_3 e$ follows. \square

The statement of Theorem 8 requires that there are no cycles in P . It is easy to show the type inference system presented in Figure 6 satisfies this requirement, provided that no type scheme constraints in the typing environment have cycles in them. However cycles would arise if we extended our language to support recursion. Our preliminary investigations suggest we can deal with cycles by annotating recursive definitions with ground types (thus ensuring the upper bound requirement of the proof). Reasoning about coherence of programs with arbitrary recursion is future work.

6. Simplification

In this section we present our simplification algorithm which allows to simplify types prior to generalizing them in (TA-Let) (Figure 6)

$$\begin{array}{c} \frac{}{P \xrightarrow{\text{simplify}(\bar{\mu})} P'; \cdot} \\ \frac{\pi, P \xrightarrow{\text{simplify}(\bar{\mu})} P'; \theta \quad P_2, P_1 \xrightarrow{\text{simplify}(\bar{\mu})} P'; \theta}{\pi, \pi, P \xrightarrow{\text{simplify}(\bar{\mu})} P'; \theta} \quad \frac{P_1, P_2 \xrightarrow{\text{simplify}(\bar{\mu})} P'; \theta}{P_1, P_2 \xrightarrow{\text{simplify}(\bar{\mu})} P'; \theta} \\ \text{S-}\uparrow \frac{\begin{array}{l} up\text{-bnd}(\pi) = \mu \quad \bar{\mu} = \mu, \bar{\mu}' \\ lo\text{-bnd}(\pi) = (Id, m) \vee lo\text{-bnd}(\pi) = (m, Id) \\ flowsFrom_P \mu \neq \{\} \quad flowsTo_P \mu = \{\} \\ \theta = (\mu \mapsto m) \quad \theta P \xrightarrow{\text{simplify}(\bar{\mu}')} P'; \theta' \end{array}}{\pi, P \xrightarrow{\text{simplify}(\bar{\mu})} \theta\pi, P'; \theta, \theta'} \\ \text{S-}\downarrow \frac{\begin{array}{l} \bar{\mu} = \mu, \bar{\mu}' \\ lo\text{-bnd}(\pi) = (Id, \mu) \vee lo\text{-bnd}(\pi) = (\mu, Id) \\ flowsFrom_P \mu = \{\} \quad flowsTo_P \mu \neq \{\} \\ \theta = (\mu \mapsto up\text{-bnd}(\pi)) \quad \theta P \xrightarrow{\text{simplify}(\bar{\mu}')} P'; \theta' \end{array}}{\pi, P \xrightarrow{\text{simplify}(\bar{\mu})} \theta\pi, P'; \theta, \theta'} \end{array}$$

Figure 8. Constraint simplification

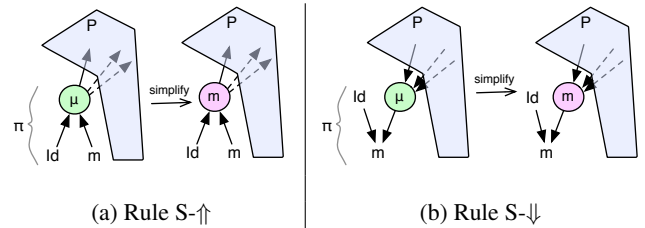


Figure 9. Visual depiction of simplification (Fig. 8)

while eliminating open variables. Simplification makes types easier to read while not reducing their generality.

Figure 8 presents inference rules for the judgment $P \xrightarrow{\text{simplify}(\bar{\mu})} P'; \theta$, which states that constraints P can be simplified to constraints P' (of the same cardinality) according to substitution θ which has a domain that is subset of monad variables $\bar{\mu}$. When referenced in (TA-Let), $\bar{\mu}$ is the list of open variables which appear in constraints P but not the let-bound term's final type. The goal here is to eliminate as many of these variables as possible.

The first rule is the identity rule, performing no solving. The second rule drops duplicate constraints. The third rule permits constraints to be reordered. The last two perform the real work.

Rule S- \uparrow solves monad variable μ with monad m for the situation depicted in Figure 9(a). Here, we have a single constraint π whose upper bound is an open variable μ , and whose lower bounds are some monad m and Id ; it has no other lower bounds in P . If μ also has an upper bound in P then we may substitute μ with its lower bound. This rule is justified by appealing to the polymonad laws. First, the substitution will convert $(Id, m) \triangleright \mu$ into $(Id, m) \triangleright m$ but this creates no additional burden on typing since this identity morphism is always guaranteed to exist by the Functorial Law and the Paired Morphisms law. Second, the substitution will convert all constraints $(\mu, m_1) \triangleright m_2$ in P to $(m, m_1) \triangleright m_2$. Any context that could have satisfied the original constraints can also satisfy these new constraints by the composition closure law: since $m \succ \mu$ and $(\mu, m_1) \triangleright m_2$ then Σ must also contain $(m, m_1) \triangleright m_2$. Note that this reasoning, and the rule, ap-

$$\begin{array}{c}
\text{S-Const} \frac{\Sigma \Vdash \cdot P \quad P' \xrightarrow{\text{simplify}(\bar{\mu})} P''; \theta'}{P, P' \xrightarrow{\text{simplify}(\bar{\mu})} P''; \theta, \theta'} \\
\\
\text{S-Meta} \frac{\pi = (m, Id) \triangleright m \vee \pi = (Id, m) \triangleright m \quad P \xrightarrow{\text{simplify}(\bar{\mu})} P'; \theta}{\pi, P \xrightarrow{\text{simplify}(\bar{\mu})} P'; \theta}
\end{array}$$

Figure 10. Constraint hiding

plies whether we start with $\pi = (Id, m) \triangleright \mu$ or $\pi = (m, Id) \triangleright \mu$ (we have depicted only the first of these).

Rule S- \Downarrow follows similar reasoning, but in the reverse direction. The solution $\mu \mapsto m$ again produces an identity morphism, which always exists, and alters existing constraints $(m_1, m_2) \triangleright \mu$ to be $(m_1, m_2) \triangleright m$. The latter constraints are again justified by composition closure; i.e., since $\mu \succ m$ and $(m_1, m_2) \triangleright \mu$ then Σ must also contain $(m_1, m_2) \triangleright m$.

Example. Recall the session types example we gave in Section 3.2. We apply (S- \Uparrow) several times to simplify the principal type to the more readable one. For example, we apply it to the sixth constraint $(Id, Sess (recv \ \gamma \ int) \ \gamma) \triangleright \mu_6$ and induce the substitution $\mu_6 \mapsto Sess (recv \ \gamma \ int) \ \gamma$ which is applied to the seventh constraint, making it eligible for simplification with (S- \Uparrow) too (substituting for μ_5). The process continues to the second constraint and the fourth, and then stops after substituting $Sess (recv \ \gamma \ int) \ \gamma$ for μ_8 in the first constraint. We can also apply (S- \Uparrow) to the third constraint, which has a different session type, which then propagates to the fifth constraint and the first. At this point we have the following set of constraints:

$$\begin{aligned}
&(Sess (send \ \alpha \ \beta) \ \alpha, Sess (recv \ \gamma \ int) \ \gamma) \triangleright \mu_{10}, \\
&(Sess (recv \ \gamma \ int) \ \gamma, Id) \triangleright Sess (recv \ \gamma \ int) \ \gamma, \\
&(Id, Sess (send \ \alpha \ \beta) \ \alpha) \triangleright Sess (send \ \alpha \ \beta) \ \alpha, \\
&(Id, Sess (recv \ \gamma \ int) \ \gamma) \triangleright Sess (recv \ \gamma \ int) \ \gamma, \\
&(Id, Sess (send \ \alpha \ \beta) \ \alpha) \triangleright Sess (send \ \alpha \ \beta) \ \alpha, \\
&(Id, Sess (recv \ \gamma \ int) \ \gamma) \triangleright Sess (recv \ \gamma \ int) \ \gamma, \\
&(Id, Sess (recv \ \gamma \ int) \ \gamma) \triangleright Sess (recv \ \gamma \ int) \ \gamma,
\end{aligned}$$

Three of these constraints are duplicates so we can drop them, leaving us with the simpler type given in Section 3.2. We can further improve this type, as we show shortly.

Pleasingly, this process yields a simpler type that can be used in the same contexts as the original principal type, so we are not compromising the generality of the code by simplifying its type.

Lemma 9 (Simplification improves types).

Given σ and σ' where σ is $\forall \bar{v}. P \Rightarrow \tau$ and σ' is an improvement of σ' , having form $\forall \bar{v}'. P' \Rightarrow \tau$ where $P \xrightarrow{\text{simplify}(\bar{\mu})} P'; \theta$ and $\bar{v}' = \bar{v} - \text{dom}(\theta)$. Then for all $P'', \Gamma, x, e, m, \tau$, if $P'' \mid \Gamma, x : \sigma \vdash e : m \tau$ such that $\Sigma \Vdash P''$ then there exists some P''' such that $P''' \mid \Gamma, x : \sigma' \vdash e : m \tau$ and $\Sigma \Vdash P'''$.

Proof. The proof is by induction on the derivation $P'' \mid \Gamma, x : \sigma \vdash e : m \tau$. Most cases are by assumption or induction, with the interesting one being (TA-Var) where the variable in question is x , and we know that all of the constraints are solvable according to the reasoning we used to justify the simplifications, above. \square

6.1 Hiding variables and constraints

While constraint simplification is effective, it can leave behind uninteresting constraints. Along with the first and third rules from Figure 8, we can apply the rules in Figure 10 to *hide* uninteresting constraints (along with their type variables, if they only appear

in these constraints). Rule (S-Const) allows us to drop constraints that refer only to constant binds, i.e., those proved with an empty substitution. The intuition is that constant binds could just be in-lined in the function that requires them, so they do not communicate any useful information. Rule (S-Meta) drops constraints of the form $(m, Id) \triangleright m$ and $(Id, m) \triangleright m$ because, by the Functorial and Paired Morphism laws, they must exist for all polymonads m . The constraints are thus not communicating useful information.

Applying these rules to our session types example, we can drop all but the first constraint resulting in a far simpler final type.

$$\begin{aligned}
&\forall \alpha, \beta, \gamma, \mu_{10}. \\
&(Sess (send \ \alpha \ \beta) \ \alpha), (Sess (recv \ \gamma \ int) \ \gamma) \triangleright \mu_{10} \Rightarrow \\
&(\beta \rightarrow \mu_{10} \ int)
\end{aligned}$$

Lemma 9 still holds even when adding in these rules. The structure of the proof is unchanged, and the (TA-Var) case will be augmented to justify the additional simplifications. Note that if we were to include these rules in Simplification directly, the *elaborated* terms would no longer be type correct (we would break Lemma 3) since placeholders for binds inserted in the elaborated term v in the (TA-Let) rule would no longer be abstracted. Therefore we simply imagine these extra constraints being hidden by the IDE.

7. Implementation

We have implemented our type inference algorithm for the simple language given in Figure 3. Our prototype implementation uses the empty theory Φ , and thus constraints must be encoded by enumeration. For example, for the *IST* polymonad from Section 3.3 we define several binds instead of the single *bIST*. In addition to *IST*, we have implemented all of the monad examples from our prior paper [24], the session types example from Section 3.2 and an example of Danielsson's [6] computational complexity monad. We would require non-empty theories for the remaining examples.

Because we have no separate theory, we implement entailment $\Sigma \Vdash \theta \pi \rightsquigarrow \bar{b}$ by unifying π with some $b:s \in \Sigma$ but with its quantified type variables replaced with fresh (unification) variables. Since there may be many $b:s$ against which we can unify π (producing substitution θ), we may need to consider each possibility. Thus, to solve a set of constraints $P = \pi, P'$, we consider each possible substitution θ for π , then attempt to solve the constraints $\theta P'$, composing the result with θ until we find one that works (or fail).

The worst-case complexity of the algorithm is

$$O\left(\sum_{k=1}^{|P|} (|\Sigma|^k)\right)$$

The number of substitutions considered at each step is $|\Sigma|$ when all bind types $b:s$ can be unified with π . (This will happen, for example, when $\pi = (\mu_1, \mu_2) \triangleright \mu_3$.) Thus the performance of the algorithm crucially depends on the number of the monad variables in π , so the order in which the constraints is handled is important. Our prototype employs topological sort for cycle-free subsets of constraints. In this case both lower bounds of a constraint are ground and we only have to solve its upper bound.

When polymonads form a join lattice (as was true of monads in our prior work [24]) we can use a simpler algorithm that performs linearly in the number of constraints, in the absence of cycles in the constraint graph. Because lubs exist for any set of polymonad instances, instead of computing all sound solutions θ we can solve every monad variable to the lub of the lower bounds of constraints that flow to this variable.

8. Related work

A variety of past work has aimed to refine the conventional notion of monads. Several examples, including Atkey's parameterized

monads [2], Wadler and Thiemann’s indexed monads [27], and applications thereof, were cited in the introduction and given in Section 3. Each of these constructions can be viewed as an instance of a polymonad. Filliâtre [10] proposed *generalized monads* as a means to more carefully reason about effects in a monadic style, and his work bears a close resemblance to Wadler and Thiemann’s. Generalized monads can also be seen as instances of polymonads—it is easy to show that the polymonad laws imply Filliâtre’s six required identities. Conversely, it is clear that some useful examples cannot be expressed using any of these prior refinements to monads; for example, our *IST* polymonad cannot be expressed due to its exclusion of certain (information-flow-violating) compositions. Thus polymonads provide greater expressive power.

Kmett’s `Control.Monad.Parameterized` Haskell package [16] provides a typeclass for non-uniform binds, with the goal of generalizing monadic programming. One key limitation is that Kmett’s bind $(m_1, m_2) \triangleright m_3$ must be *functionally dependent*; i.e., m_3 must be a function of m_1 and m_2 . As such, it is not possible to program morphisms between different monadic constructors, i.e., the pair of binds $(m_1, Id) \triangleright m_2$ and $(m_1, Id) \triangleright m_3$ would be forbidden, so there would be no way to convert from m_1 to m_2 and from m_1 to m_3 in the same program. Polymonads do not have this limitation. Kmett does not discuss laws that should govern the proper use of non-uniform binds.

Another line of past work has focused on making monadic programming easier. Haskell’s `do` notation exposes the structure of a monadic computation, and typeclass inference can determine which binds and units should be used, but the placement of morphisms is left to the programmer. The problem is that the use of morphisms (e.g., if defined as a typeclass) would frequently lead to open type variables, which Haskell’s typeclass inference deems ambiguous. Inference with Kmett’s class has the same problems.

For ML, our own prior work [24] showed that no additional notation is needed: left-to-right, call-by-value evaluation order makes the order of operations well-defined, so the syntactic structure of the program indicates where binds, units, and even morphisms should be placed. Moreover, we proved that the monad laws ensured that open variables could be solved arbitrarily without affecting semantics, and so there was no need to reject programs with open types. Our present paper generalizes the approach of this prior paper in two ways, first by applying it to the more general polymonadic construction, and second by showing more rigorously how to reason about two arbitrary, different solutions. Interestingly, the present work arose when we discovered we could not write *IST* as a monad, since a monad would require the existence of binds that could violate an information flow property.

As mentioned in the introduction and Section 2, concurrently with our work Tate developed a general semantic framework called *productors* for describing *producer effect* systems [25]. This framework turns out to match our definition of polymonads in many important cases. Our work is complementary to Tate’s in that we consider practical programming concerns (i.e., type inference and rewriting) in a higher-order functional programming language, whereas he focuses on semantic foundations.

9. Conclusions

We have presented *polymonads*, a generalization of prior monad-like programming idioms. We have shown polymonads to be useful, using them to encode a variety of prior programming constructions. We have also shown how to facilitate programming with them in a direct style—the programmer can use a polymonadic computation $m \tau$ as if it were of type τ and our novel type inference and rewriting algorithm will insert the necessary coercions. Rewritten programs are *coherent*: all solutions to variables not present in the

final type will induce the same semantics. Pleasingly, our algorithm produces general types that are nevertheless simple.

References

- [1] M. Abadi, A. Banerjee, N. Heintze, and J. Riecke. A core calculus of dependency. In *POPL*, volume 26, pages 147–160, 1999.
- [2] R. Atkey. Parameterised notions of computation. *J. Funct. Program.*, 19(3-4):335–376, 2009.
- [3] R. Atkey. Parameterised notions of computation. *Journal of Functional Programming*, 19(3 & 4):335–376, 2009.
- [4] G. Cooper and S. Krishnamurthi. Embedding dynamic dataflow in a call-by-value language. In *ESOP*, 2006.
- [5] K. Crary, A. Kliger, and F. Pfenning. A monadic analysis of information flow security with mutable state. *Journal of functional programming*, 15(02):249–291, 2005.
- [6] N. Danielsson. Lightweight semiformal time complexity analysis for purely functional data structures. In *POPL*, pages 133–144, 2008.
- [7] D. Devriese and F. Piessens. Information flow enforcement in monadic libraries. In *TLDI*, pages 59–72, 2011.
- [8] C. Elliott and P. Hudak. Functional reactive animation. In *ICFP*, pages 263–273, 1997.
- [9] A. Filinski. Representing layered monads. In *POPL*, pages 175–188, 1999.
- [10] J.-C. Filliâtre. A theory of monads parameterized by effects, 1999.
- [11] J.-Y. Girard. *Interprétation fonctionnelle et élimination des coupures de l’arithmétique d’ordre supérieur*. PhD thesis, Université Paris VI I, 1972.
- [12] J. Goguen and J. Meseguer. Security policy and security models. In *Symposium on Security and Privacy*, pages 11–20, 1982.
- [13] G. Hutton and E. Meijer. Monadic Parsing in Haskell. *JFP*, 8(4), 1998.
- [14] M. P. Jones. A theory of qualified types. In *ESOP*, 1992.
- [15] O. Kiselyov and C. Shan. Lightweight monadic regions. In *ACM SIGPLAN Notices*, volume 44, pages 1–12. ACM, 2008.
- [16] E. Kmett. `Control.Monad.Parameterized` package. <http://hackage.haskell.org/packages/archive/monad-param/0.0.4/doc/html/Control-Monad-Parameterized.html>, 2012.
- [17] P. Li and S. Zdancewic. Encoding information flow in Haskell. In *CSFW*, pages 16–27, 2006.
- [18] E. Moggi. Computational lambda-calculus and monads. In *LICS*, 1989.
- [19] Monad laws. http://www.haskell.org/haskellwiki/Monad_Laws, 2012.
- [20] I. Neamtiu, M. Hicks, J. S. Foster, and P. Pratikakis. Contextual effects for version-consistent dynamic software updating and safe concurrent programming. In *POPL*, 2008.
- [21] R. Pucella and J. Tov. Haskell session types with (almost) no class. In *ACM SIGPLAN Notices*, volume 44, pages 25–36. ACM, 2008.
- [22] N. Ramsey and A. Pfeffer. Stochastic lambda calculus and monads of probability distributions. In *POPL*, pages 154–165, 2002.
- [23] A. Russo, K. Claessen, and J. Hughes. A library for light-weight information-flow security in Haskell. In *Haskell*, 2008.
- [24] N. Swamy, N. Guts, D. Leijen, and M. Hicks. Lightweight monadic programming in ML. In *ICFP*, 2011.
- [25] R. Tate. The sequential semantics of producer effect systems, 2012.
- [26] P. Wadler. The essence of functional programming. In *POPL*, 1992.
- [27] P. Wadler and P. Thiemann. The marriage of effects and monads. *ACM Trans. Comput. Logic*, 4:1–32, January 2003.

A. Proof of Lemma 1

We separately prove the two directions of the implication. First we show that for all S such that $\langle S \rangle = \mathcal{M}, \Sigma_{\mathcal{M}}$, if $\models S$ then $\models \langle S \rangle$. (Note that $\text{bind}_{Id, Id, Id} \in \Sigma_{\mathcal{M}}$ as required since $S \models Id \triangleright Id$.)

Proof. Let S be a monadic signature. Suppose that $\models S$. To show $\models \Sigma_{\mathcal{M}}$ we prove the six polymonad laws hold for all the binds defined by $\Sigma_{\mathcal{M}} = \text{Clos}(S)$.

Functorial law: For all $M \in \mathcal{M}$, we have that $S \models Id \triangleright M$ and $S \models M \triangleright M$ so by $\text{Clos}(S)$ we must have that $(M, Id) \triangleright M \in \Sigma_{\mathcal{M}}$.

Left identity: Suppose that for some $M_i, M_j \in \mathcal{M}$, we have

$$\{\text{bind}_{M_i, M_j, M_j}, \text{unit}_{M_i}\} \in \Sigma_{\mathcal{M}}$$

$$\begin{aligned} & \text{bind}_{M_i, M_j, M_j} (\text{unit}_{M_i} e) k \\ = & \text{bind}_{M_j} (f_{M_i, M_j} (\text{unit}_{M_i} e)) (\lambda x. f_{M_j, M_j} (k x)) \quad \text{def} \\ = & \text{bind}_{M_j} (f_{M_i, M_j} (\text{unit}_{M_i} e)) x \quad \text{(vii)} \\ = & \text{bind}_{M_j} (\text{unit}_{M_j} e) k \quad \text{(iv)} \\ = & k e \quad \text{(i)} \end{aligned}$$

Right identity: Suppose that we have

$$\{\text{bind}_{M_i, M_j, M_i}, \text{unit}_{M_j}\} \in \Sigma_{\mathcal{M}}$$

$$\begin{aligned} & \text{bind}_{M_i, M_j, M_i} e \text{unit}_{M_j} \\ = & \text{bind}_{M_i} (f_{M_i, M_i} e) (\lambda x. f_{M_j, M_i} (\text{unit}_{M_j} x)) \quad \text{def} \\ = & \text{bind}_{M_i} e (\lambda x. f_{M_j, M_i} (\text{unit}_{M_j} x)) \quad \text{(vii)} \\ = & \text{bind}_{M_i} e \text{unit}_{M_i} \quad \text{(iv)} \\ = & e \quad \text{(ii)} \end{aligned}$$

Associativity: (1) (\Rightarrow) direction. Suppose that we have

$$\text{bind}_{M_1, M_2, M_{12}} \in \Sigma_S \text{ and } \text{bind}_{M_{12}, M_3, M_{123}} \in \Sigma_S$$

From the definition of closure, we have $M_1 \triangleright M_{12} \in S$, $M_2 \triangleright M_{12} \in S$, $M_{12} \triangleright M_{123} \in S$, and $M_3 \triangleright M_{123} \in S$. Then we can choose $M_{23} = M_{123}$, and easily show $\text{bind}_{M_2, M_3, M_{23}} \in \Sigma_S$ and $\text{bind}_{M_1, M_{23}, M_{123}} \in \Sigma_S$. Indeed, $M_1 \triangleright M_{123} \in S$ and $M_2 \triangleright M_{123} \in S$ follow by transitivity, $M_3 \triangleright M_{123} \in S$ is an assumption, and $M_{123} \triangleright M_{123} \in S$ is an axiom. The direction (\Leftarrow) is analogous.

(2) We want to show that

$$\begin{aligned} & \text{bind}_{M_{12}, M_3, M_{123}} (\text{bind}_{M_1, 2, M_{12}} e_1 k_2) k_{M_3} \\ = & \text{bind}_{M_1, M_j, M_{123}} e_1 (\lambda z. \text{bind}_{2, M_3, M_j} (k_2 e_1) k_{M_3}) \end{aligned}$$

provided that all the mentioned binds exist. We are going to show that both the left and the right sides of this equality are equivalent to the same expression, as shown in Figure 11.

Paired morphisms: Suppose $(M_1, Id) \triangleright M_2 \in \Sigma_{\mathcal{M}}$, then by def. of $\text{Clos}(S)$ we must have that $S \models M_1 \triangleright M_2$ and $S \models Id \triangleright M_2$ which, by def. of $\text{Clos}(S)$ implies $(Id, M_1) \triangleright M_2 \in \Sigma_{\mathcal{M}}$. Same argument goes in the reverse direction.

Composition closure: Suppose $(M_1, M_2) \triangleright M_3 \in \Sigma_{\mathcal{M}}$. By def of closure, we know $S \models M_1 \triangleright M_3$ and $S \models M_2 \triangleright M_3$. Suppose $M'_2 \succ M_2$, $M'_1 \succ M_1$, and $M_3 \succ M'_3$; this implies that $(M'_2, Id) \triangleright M_2 \in \Sigma_{\mathcal{M}}$ and thus that $S \models M'_2 \triangleright M_2$; we similarly know that $S \models M'_1 \triangleright M_1$ and $S \models M_3 \triangleright M'_3$. But since $S \models$ is transitive, we have $S \models M'_2 \triangleright M'_3$ and $S \models M'_1 \triangleright M'_3$ which implies $(M'_1, M'_2) \triangleright M'_3 \in \text{Clos}(S)$. \square

Now we show that for all S such that $\langle S \rangle = \mathcal{M}, \Sigma_{\mathcal{M}}$ then $\models \langle S \rangle$ implies $\models S$.

Proof. To show $\models S$ we prove that monad laws (i-vi) hold for all the binds in S such that $\text{Clos}(S) = \Sigma_{\mathcal{M}}$. For each law, we start by assuming the given binds/units/morphisms are defined in S , which implies the corresponding existence of binds in $\Sigma_{\mathcal{M}}$. Existence of other binds used in proofs of laws (iv)–(vi) is justified at the end.

$$\begin{aligned} \text{(i)} \quad & \text{bind}_M (\text{unit}_M e) k \\ = & \text{bind}_{M, M, M} (\text{unit}_M e) k \quad \text{def} \\ = & k e \quad \text{left id} \end{aligned}$$

$$\begin{aligned} \text{(ii)} \quad & \text{bind}_M e \text{unit}_M \\ = & \text{bind}_{M, M, M} e \text{unit}_M \quad \text{def} \\ = & e \quad \text{right id} \end{aligned}$$

$$\begin{aligned} \text{(iii)} \quad & \text{bind}_M (\text{bind}_M e k_1) k_2 \\ = & \text{bind}_{M, M, M} (\text{bind}_{M, M, M} e k_1) k_2 \quad \text{def} \\ = & \text{bind}_{M, M, M} e (\lambda x. \text{bind}_{M, M, M} (k_1 x) k_2) \quad \text{Assoc} \\ = & \text{bind}_M e (\lambda x. \text{bind}_M (k_1 x) k_2) \quad \text{def} \end{aligned}$$

$$\begin{aligned} \text{(iv)} \quad & f_{M_1, M_2} (\text{unit}_{M_1} e) \\ = & \text{bind}_{M_1, Id, M_2} (\text{bind}_{Id, Id, M_1} e \text{unit}_{Id}) \text{unit}_{Id} \quad \text{def} \\ = & \text{bind}_{Id, Id, M_2} e (\lambda x. \text{bind}_{Id, Id, Id} (\text{unit}_{Id} x) \text{unit}_{Id}) \quad \text{Assoc} \\ = & \text{bind}_{Id, Id, M_2} e \text{unit}_{Id} \quad \text{left id} \\ = & \text{unit}_{M_2} e \quad \text{def} \end{aligned}$$

$$\begin{aligned} \text{(v)} \quad & \text{bind}_{M_2} (f_{M_1, M_2} e) (f_{M_1, M_2} \circ k) \\ = & \text{bind}_{M_2, M_2, M_2} (\text{bind}_{M_1, Id, M_2} e \text{unit}_{Id}) \quad \text{def} \\ & (f_{M_1, M_2} \circ k) \\ = & \text{bind}_{M_1, M_2, M_2} e (\lambda x. \text{bind}_{Id, M_2, M_2} (\text{unit}_{Id} x) (f_{M_1, M_2} \circ k)) \quad \text{Assoc} \\ = & \text{bind}_{M_1, M_2, M_2} e (f_{M_1, M_2} \circ k) \quad \text{left id} \\ = & \text{bind}_{M_1, M_2, M_2} e (\lambda x. \text{bind}_{M_1, Id, M_2} (k x) \text{unit}_{Id}) \quad \text{def} \\ = & \text{bind}_{M_1, Id, M_2} (\text{bind}_{M_1, M_1, M_1} e k) \text{unit}_{Id} \quad \text{Assoc} \\ = & f_{M_1, M_2} (\text{bind}_{M_1} e k) \quad \text{def} \end{aligned}$$

$$\begin{aligned} \text{(vi)} \quad & f_{M_2, M_3} \circ f_{M_1, M_2} \\ = & \lambda x. \text{bind}_{M_2, Id, M_3} (\text{bind}_{M_1, Id, M_2} x \text{unit}_{Id}) \text{unit}_{Id} \quad \text{def} \\ = & \lambda x. \text{bind}_{M_1, Id, M_3} x (\lambda y. \text{bind}_{Id, Id, Id} (\text{unit}_{Id} y) \text{unit}_{Id}) \quad \text{Assoc} \\ = & \lambda x. \text{bind}_{M_1, Id, M_3} x \text{unit}_{Id} \quad \text{left id} \\ = & f_{M_1, M_3} \quad \text{def} \end{aligned}$$

$$\begin{aligned} \text{(vii)} \quad & f_{M, Id, M} e \text{unit}_{Id} \\ = & e \quad \text{def} \\ & \text{right id} \end{aligned}$$

For (iv) and (v), we know that $S \models Id \triangleright M_2$ and $S \models M_2 \triangleright M_2$ and $S \models M_1 \triangleright M_2$, and thus

$$\{\text{bind}_{Id, Id, M_2}, \text{bind}_{Id, M_2, M_2}, \text{bind}_{M_1, M_2, M_2}\} \subseteq \text{Clos}(S)$$

For (vi), we have $S \models Id \triangleright M_3$ and $S \models M_1 \triangleright M_3$ which implies $\text{bind}_{1, Id, 3} \in \text{Clos}(S)$. \square

B. Well-typed elaborations

This section includes the full definitions and proof of Lemma 3.

$$\begin{aligned}
& bind_{M_{12}, M_3, M_{123}} (bind_{M_1, M_2, M_{12}} e_1 k_2) k_3 \\
= & bind_{M_{123}} (f_{M_{12}, M_{123}} (bind_{M_{12}} (f_{M_1, M_{12}} e_1) (\lambda x. f_{M_2, M_{12}} (k_2 x))) (\lambda y. f_{M_3, M_{123}} (k_3 y))) && \text{def of Clos} \\
= & bind_{M_{123}} (bind_{M_{123}} (f_{M_{12}, M_{123}} (f_{M_1, M_{12}} e_1)) (\lambda x. f_{M_{12}, M_{123}} (f_{M_2, M_{12}} (k_2 x)))) (\lambda y. f_{M_3, M_{123}} (k_3 y)) && \text{(v)} \\
= & bind_{M_{123}} (bind_{M_{123}} (f_{M_1, M_{123}} e_1) (\lambda x. (f_{M_2, M_{123}} (k_2 x)))) (\lambda y. f_{M_3, M_{123}} (k_3 y)) && \text{(vi)} \\
= & bind_{M_{123}} (f_{M_1, M_{123}} e_1) (\lambda z. bind_{M_{123}} (f_{M_2, M_{123}} (k_2 z)) (\lambda y. f_{M_3, M_{123}} (k_3 y))) && \text{(vi)} \\
\\
& bind_{M_1, M_j, M_{123}} e_1 (\lambda z. bind_{M_2, M_3, M_j} (k_2 e_1) k_3) \\
= & bind_{M_{123}} (f_{M_1, M_{123}} e_1) (\lambda z. f_{M_j, M_{123}} (bind_{M_j} (f_{M_2, M_j} (k_2 z)) (\lambda y. f_{M_3, M_j} (k_3 y)))) && \text{def of Clos} \\
= & bind_{M_{123}} (f_{M_1, M_{123}} e_1) (\lambda z. (bind_{M_{123}} (f_{M_j, M_{123}} (f_{M_2, M_j} (k_2 z))) (\lambda y. f_{M_j, M_{123}} (f_{M_3, M_j} (k_3 y))))) && \text{by (v)} \\
= & bind_{M_{123}} (f_{M_1, M_{123}} e_1) (\lambda z. (bind_{M_{123}} ((f_{M_2, M_{123}} (k_2 z)) (\lambda y. (f_{M_3, M_{123}} (k_3 y))))) && \text{by (v)}
\end{aligned}$$

Figure 11. Associativity case for proof in Lemma 1

Notations and auxiliary results

$$\begin{aligned}
\{\cdot\} &= \cdot \\
\{\Gamma, x : \forall \bar{v}. P \Rightarrow \tau\} &= \{\{\Gamma\}, x : \forall \bar{v}. \text{abstyp}(P, \tau)\} \\
\text{abstyp}(\cdot, \tau) &= \tau \\
\text{abstyp}(\pi, P, \tau) &= \pi \rightarrow \text{abstyp}(P, \tau) \\
\text{Abs}((m_1, m_2) \triangleright m, P) &= \mathbf{b}_{m_1, m_2, m} : (m_1, m_2) \triangleright m, \text{Abs}(P) \\
\text{Abs}(\cdot) &= \cdot
\end{aligned}$$

We assume the following property about simplification:

Definition 10 (Type-preserving simplification). *Constraint simplification is type-preserving when for all $P, P', \bar{\mu}, \theta$,*

if $P \xrightarrow{\text{simplify}(\bar{\mu})} P'; \theta$ then $\text{dom}(\theta) \subseteq \bar{\mu}$ and $\forall \pi. \pi \in P \Leftrightarrow \theta \pi \in P'$.

The following lemma says that if two multisets of constraints correspond to the same set, then their corresponding typing environment are equivalent:

Lemma 11. *For all Γ, P, P', e, t , if $\{\{\Gamma\}, \text{Abs}(P) \vdash e : t$ and $\forall \pi. \pi \in P \Leftrightarrow \pi \in P'$, then $\{\{\Gamma\}, \text{Abs}(P') \vdash e : t$*

The proof is straightforward by induction on the derivation.

We also use this auxiliary lemma to relate different functions of P :

Lemma 12 (Abstraction and application of evidence). *For all Γ, P, e, t , we have*

1. $\{\{\Gamma\} \vdash_F \text{abs}(P, e) : \text{abstyp}(P, t) \text{ iff } \{\{\Gamma\}, \text{Abs}(P) \vdash_F e : t$
2. $\{\{\Gamma\}, \text{Abs}(P) \vdash_F \text{app}(P, e) : t \text{ iff } \{\{\Gamma\} \vdash_F e : \text{abstyp}(P, t)$

The proof is straightforward by induction on the derivation.

Proof of Lemma 3

Proof. Suppose that there are e, τ, m, e, Γ such that either $P \mid \Gamma \vdash e : \tau \rightsquigarrow e$ or $P \mid \Gamma \vdash e : m\tau \rightsquigarrow e$. We want to show that $\{\{\Gamma\}, \text{Abs}(P) \vdash_F e : \tau$, or $\{\{\Gamma\}, \text{Abs}(P) \vdash_F e : m\tau$.

The proof is by induction on the derivation $P \mid \Gamma \vdash e : \tau \rightsquigarrow e$ (or $P \mid \Gamma \vdash e : m\tau \rightsquigarrow e$).

Sub-case TA-Var: We want to show that

$$\{\{\Gamma\}, \text{Abs}(\theta P) \vdash_F \text{app}(x, \theta P) : \theta \tau \bar{\mu}' \bar{\tau}$$

From the hypothesis of (TA-Var) we know that $\Gamma(x) = \forall \bar{\mu}. \bar{\alpha}. P \Rightarrow \tau$, and there is $\theta = [\bar{\mu}' / \bar{\mu}][\bar{\tau} / \bar{\alpha}]$, and $\bar{\mu}'$ fresh. So in the translated environment we get $\{\{\Gamma\}(x) = \forall \bar{\mu}. \bar{\alpha}. \text{abstyp}(P, \tau)$. We can instantiate this type with θ and get $\{\{\Gamma\} \vdash_F x : \text{abstyp}(\theta P, \theta \tau \bar{\mu}' \bar{\tau})$, and then conclude using Lemma 12.2.

Sub-case TA-Lam and TA-Val: By induction.

Sub-case TA-Do: We need to show that

$$\{\{\Gamma\}, \text{Abs}(P) \vdash_F \mathbf{b}_{m_1, m_2, \mu_3} \tau_1 \tau_2 \tau_1 (\lambda x : \mathbf{t}_1. e_2) : m \tau_2$$

when by induction hypothesis we have

$$P = P_1, P_2, (m_1, m_2) \triangleright \mu_3$$

$$\{\{\Gamma\}, \text{Abs}(P_1) \vdash_F e_1 : m_1 \tau_1$$

$$\{\{\Gamma\}, x : \tau_1, \text{Abs}(P_2) \vdash e_2 : m_2 \tau_2$$

Since $(m_1, m_2) \triangleright \mu_3 \in P$, by definition $\text{Abs}(P)$ contains the binding $\mathbf{b}_{m_1, m_2, \mu_3} : (m_1, m_2) \triangleright \mu_3$. We can type the goal using instantiation and application rules of \vdash_F .

Sub-case TA-Let: We want to show that

$$\{\{\Gamma\}, \text{Abs}(P) \vdash_F (\lambda x : \forall \bar{v}'. \text{abstyp}(P'', \tau'). e) \text{abs}(P'', \theta v) : m \tau$$

By induction hypothesis we know

$$(IH1) \quad \{\{\Gamma\}, \text{Abs}(P') \vdash_F v : \tau'$$

$$(IH2) \quad \{\{\Gamma, x : \forall \bar{v}'. \text{abstyp}(P'', \tau')\}, \text{Abs}(P) \vdash_F e : m \tau$$

when

$$\bar{\mu}, \bar{\alpha} = \bar{v} = \text{ftv}(P' \Rightarrow \tau') \setminus \text{ftv}(\Gamma)$$

$$P' \xrightarrow{\text{simplify}(\bar{\mu} \setminus \text{ftv}(\tau'))} P''; \theta$$

$$\bar{v}' = \bar{v} \setminus \text{dom}(\theta)$$

Using (IH2) we type the function:

$$\{\{\Gamma\}, \text{Abs}(P) \vdash_F (\lambda x : \forall \bar{v}'. \text{abstyp}(P'', \tau'). e) : (\forall \bar{v}'. \text{abstyp}(P'', \tau')) \rightarrow m \tau$$

We apply θ to (IH1) and note that $\text{dom}(\theta) \cap (\text{ftv}(\tau') \cup \text{ftv}(\Gamma)) = \emptyset$ by Definition 10. So we get

$$\{\{\Gamma\}, \theta \text{Abs}(P') \vdash_F \theta v : \tau'$$

Since $\forall \pi. \pi \in P'' \Leftrightarrow \theta P'$, we apply Lemma 11, and get $\{\{\Gamma\}, \text{Abs}(P'') \vdash_F \theta v : \tau'$. We can weaken the typing environment with $\text{Abs}(P)$ and get

$$\{\{\Gamma\}, \text{Abs}(P), \text{Abs}(P'') \vdash_F \theta v : \tau'$$

So we can type the argument of the application in the goal:

$$\{\{\Gamma\}, \text{Abs}(P) \vdash_F \text{abs}(P'', \theta v) : \text{abstyp}(P'', \tau')$$

So the function application in the goal can be given type $m \tau$. \square

C. Definitions and proofs for Section 5

The statement of Theorem 8 assumes that the constraints generated by our type inference algorithm are cycle-free. We show that this property always holds under some restrictions on the shape of type schemes in the typing environment. Intuitively, this holds because our language does not include recursion.

Definition 13 (Cycle-free type environment). *We say that there is a path from m to m' in P when there exist m_1, m_2 such that*

$$(m, m_1) \triangleright m_2 \in P \text{ or } (m_1, m) \triangleright m_2 \in P$$

and either $m_2 = m'$ or there is a path from m_2 to m' . We say that there is a cycle in P when there exists m , and a path from m to m . Γ is cycle-free when for all $x : \forall \bar{v}. P \Rightarrow \tau$, P is cycle free.

Lemma 14 (Type inference produces cycle-free constraints). *For all Γ, P, e, m, τ, e , if Γ is cycle-free and $\Gamma|P \vdash e : \tau \rightsquigarrow e$ (or $\Gamma|P \vdash e : m\tau \rightsquigarrow e$), then P is cycle-free.*

Proof. By induction on the typing derivation $\Gamma|P \vdash e : \tau \rightsquigarrow e$ (or $\Gamma|P \vdash e : m\tau \rightsquigarrow e$).

Case TA-Var and TA-Const: Γ must contain the scheme $\forall \bar{\mu}, \bar{\alpha}. P' \Rightarrow \tau$ and the resulting constraints are $P = [\bar{\mu}'/\bar{\mu}][\bar{\tau}/\bar{\alpha}]P'$ for some fresh $\bar{\mu}'$. By assumption, P' has no cycles and the substitution only introduces fresh variables, so P has no cycles.

Sub-case TA-Lam and TA-ID: Induction hypothesis applies.

Sub-case TA-App and TA-Do: The resulting substitution combines constraints from subterms, which are cycle-free by the induction hypothesis, with new constraints that only add edges to fresh variables, so no cycles are created.

Sub-case TA-Let: The final constraints are the same as in the second inductive hypothesis. We can apply the induction hypothesis, because the type environment extended with the type scheme for the let bound term is still cycle free. The derivation of the bound term produces cycle free constraints P' which then get simplified, and simplification does not introduce new constraints, so no cycles. \square

Coherence for local modifications We show coherence by iteratively constructing ground solutions that have fewer differences and preserve semantics of the elaboration. We rely on a coherence of local modifications, to show that at every step of our construction we can *smoothly* modify the lower bounds of a constraint, and preserve groundness of the solution.

The proof of local coherence relies on the following lemma:

Lemma 15. *For all polymonad instances $m_1 \dots m_7$, and m'_3, m'_6 ,*

$$\{\text{bind}_{m_3 m_6 m_7}, \text{bind}_{m_1 m_2 m_3}, \text{bind}_{m_4 m_5 m_6}, \text{bind}_{m'_3 m'_6 m_7}, \text{bind}_{m_1 m_2 m'_3}, \text{bind}_{m_4 m_5 m'_6}\} \in \Sigma$$

implies

$$\text{bind}_{m_3 m_6 m_7} (\text{bind}_{m_1 m_2 m_3} e_1 k_2) (\lambda x. \text{bind}_{m_4 m_5 m_6} e_4 k_5) = \text{bind}_{m'_3 m'_6 m_7} (\text{bind}_{m_1 m_2 m'_3} e_1 k_2) (\lambda x. \text{bind}_{m_4 m_5 m'_6} e_4 k_5)$$

Proof. From the Associativity(1) polymonad law we know that there exist monads m_8, m_9 such that

$$\{\text{bind}_{m_1 m_8 m_7}, \text{bind}_{m_9 m_5 m_8}, \text{bind}_{m_2 m_4 m_9}\} \in \Sigma$$

Applying Associativity(2) to the lhs of the goal we get

$$\begin{aligned} & \text{bind}_{m_3 m_6 m_7} (\text{bind}_{m_1 m_2 m_3} e_1 k_2) (\lambda x. \text{bind}_{m_4 m_5 m_6} e_4 k_5) = \\ & \text{bind}_{m_1 m_8 m_7} e_1 \\ & (\lambda y. \text{bind}_{m_9 m_5 m_8} (\text{bind}_{m_2 m_4 m_9} (k_2 y) (\lambda z. e_4))) k_5 \end{aligned}$$

Applying Associativity(2) to the rhs of the goal we also get

$$\begin{aligned} & \text{bind}_{m'_3 m'_6 m_7} (\text{bind}_{m_1 m_2 m'_3} e_1 k_2) (\lambda x. \text{bind}_{m_4 m_5 m'_6} e_4 k_5) = \\ & \text{bind}_{m_1 m_8 m_7} e_1 \\ & (\lambda y. \text{bind}_{m_9 m_5 m_8} (\text{bind}_{m_2 m_4 m_9} (k_2 y) (\lambda z. e_4))) k_5 \end{aligned}$$

The goal follows by transitivity of =. \square

Next, we formalize the notion of a local modification θ_2 of a ground solution θ_1 to constraint set.

Definition 16 (Local modification of a solution). *Given a ground solution θ_1 to a constraint set (π, P) a local modification to θ_1 is a ground solution θ_2 for which the following conditions are true:*

1. $\theta_1 lo\text{-}bnd(\pi) \neq \theta_2 lo\text{-}bnd(\pi)$;
2. $\theta_1 P = \theta_2 P$ and $\theta_1 up\text{-}bnd(\pi) = \theta_2 up\text{-}bnd(\pi)$;
3. for all $\mu \in lo\text{-}bnd(\pi)$, if $\theta_1 \mu \neq \theta_2 \mu$ then $flowsTo_P \mu \neq \{\}$.

Conditions (1) and (2) establish that the substitutions differ for the variables that occur in $lo\text{-}bnd(\pi)$, and coincide on all other variables. Condition (3) states that the variables whose values differ in θ_1 and θ_2 must have incoming edges in P .

We suppose that the constraints generated at type inference are labeled with the name of the typing rule where they have been introduced, for example $(m_1, m_2) \stackrel{Do}{\triangleright} m_3$ denotes a bind introduced by the rule (TA-Do).

Theorem 17 (Local coherence).

Given $\Sigma_0, \Gamma, P, e, m, \tau, \theta_1, \theta_2, e$, such that

1. Signature Σ is well-formed.
2. $P \mid \Gamma \vdash e : \tau \rightsquigarrow e$ or $P \mid \Gamma \vdash e : m\tau \rightsquigarrow e$.
3. θ_1 is a ground solution for P , and θ_2 is a local modification of θ_1 .
4. $dom(\theta_1) \subseteq (ftv(P) \setminus (ftv(\tau)))$ (or $dom(\theta_1) \subseteq (ftv(P) \setminus (ftv(m\tau)))$)

Then, $\theta_1 e = \theta_2 e$.

Proof. (Sketch) Since θ_2 is a local modification of θ_1 , we have (from condition (1) of Def. 16) that $P = \pi, P'$ and $\theta_1 lo\text{-}bnd(\pi) \neq \theta_2 lo\text{-}bnd(\pi)$. We proceed by cases on the shape of π .

Case π is an Do bundle: We have $\theta_1 \pi = (m_L, m_R) \stackrel{Do}{\triangleright} m_U$ and $\theta_2 \pi = (m'_L, m'_R) \stackrel{Do}{\triangleright} m_U$. From Condition (3) of Definition 16, if $m_L \neq m'_L$, then it's a variable μ_L , and $flowsTo_P \mu_L \neq \{\}$. Similarly, if $m_R \neq m'_R$, then it's a variable μ_R , and $flowsTo_P \mu_R \neq \{\}$.

We proceed by cases on the shape of each of the constraints (π', π'') in sets $flowsTo_P \mu_L$ and $flowsTo_P \mu_R$, respectively.

Sub-case π' and π'' are Do bundles: From examining the typing rule TA-Do, we have $\theta_1 \pi' = (m_1, m_2) \stackrel{Do}{\triangleright} \mu_L$ and $\theta_2 \pi' = (m_1, m_2) \stackrel{Do}{\triangleright} \mu'_L$ for some m_1, m_2 ; we similarly have $\theta_1 \pi'' = (m_1, m_2) \stackrel{Do}{\triangleright} \mu_R$ and $\theta_2 \pi'' = (m'_1, m'_2) \stackrel{Do}{\triangleright} \mu'_R$ for some m'_1, m'_2 . From the shape of the constraints, we reason that we have a source term e of this form:

$$(\text{let } y = (\text{let } z = e_3 \text{ in } e_4) \text{ in let } x = e_1 \text{ in } e_2)$$

that is elaborated to the terms shown below, where e_1, e_2, e_3 , and e_4 are the elaborated of the sub-terms. The solution θ_1 yields the following term \hat{e} :

$$\text{bind}_{m_L, m_R, m_U} (\text{bind}_{m_1, m_2, m_L} e_3 (\lambda z. \dots e_4)) (\lambda y. \dots \text{bind}_{m'_1, m'_2, m_R} e_1 (\lambda x. \dots e_2))$$

whereas the solution θ_2 yields the following term:

$$\text{bind}_{m'_L, m'_R, m_U} (\text{bind}_{m_1, m_2, m'_L} e_3 (\lambda z. \dots e_4)) (\lambda y. \dots \text{bind}_{m'_1, m'_2, m'_R} e_1 (\lambda x. \dots e_2))$$

The equality of the terms $\theta_1 \hat{e}$ and $\theta_2 \hat{e}$ follows from Lemma 15. \square

Case Other cases: To do. \square

Proof of coherence (Theorem 8)

Proof. Suppose that for some $\Sigma, \Gamma, P, e, m, \tau, \theta_1, \theta_2, e, \bar{\mu}, \bar{\alpha}$, we have

1. Σ is well-formed and P is cycle-free.
2. $P \mid \Gamma \vdash e : \tau \rightsquigarrow e$ or $P \mid \Gamma \vdash e : m\tau \rightsquigarrow e$.
3. There exist open variables $\bar{\mu}, \bar{\alpha} = ftv(P) \setminus ftv(\tau)$ (or $\bar{\mu}, \bar{\alpha} = ftv(P) \setminus ftv(m\tau)$) such that for all $\mu \in \bar{\mu}$, the sets $flowsTo_P \mu$ and $flowsFrom_P \mu$ are non-empty, i.e., each $\mu \in \bar{\mu}$ has a lower and an upper bound.
4. θ_1 and θ_2 are ground solutions for P such that $\theta_1(\nu) = \theta_2(\nu)$ for all $\nu \notin \bar{\alpha}, \bar{\mu}$.

Since θ_1 and θ_2 are ground solutions for P under the original signature Σ , they are also ground solutions under the saturated signature Σ_{sat} . We can construct a solution θ_3 with $\text{dom}(\theta_3) = \text{dom}(\theta_2)$, such that it assigns the same monads as θ_2 to monad variables: $\forall \mu \in \bar{\mu}. \theta_3 \mu = \theta_2 \mu$, and the same types as θ_1 to type variables: $\forall \alpha \in \bar{\alpha}. \theta_3 \alpha = \theta_1 \alpha$. By construction, solutions θ_3 and θ_2 only differ for $\bar{\alpha}$, and thus θ_3 is also ground under the saturated signature Σ_{sat} . By successively applying the (Param) property (required to hold by Condition 1 of Definition 6) to every constraint that contains a variable from $\bar{\alpha}$, we show that $\theta_3 e = \theta_2 e$.

Now we prove that $\theta_1 e = \theta_3 e$ by induction on the number of monad variables whose values differ in the solutions. Notationally, we write $\theta \setminus A$ to be the substitution with domain $\text{dom}(\theta) \setminus A$ such that $(\theta \setminus A)(\nu) = \theta(\nu)$ for all $\nu \in \text{dom}(\theta \setminus A)$.

The rest of the proof is by iteration on the number n of constraints in P with variables μ in their lower bounds and which are solved differently by θ_1 and θ_3 :

$$\bar{\pi} = \{\pi \mid \pi \in P, \text{ftv}(\text{lo-bnd}(\theta_1 \pi)) \cap \bar{\mu} \neq \{\}, \theta_1 \pi \neq \theta_3 \pi\}$$

Because P is cycle-free, we can sort it topologically. For the proof we number the constraints $\bar{\pi}$ in reverse topological order:

$$\forall i, j. \text{up-bnd}(\pi_i) \in \text{lo-bnd}(\pi_j) \implies i > j \quad (\text{Topo-}\pi)$$

At each iteration i we construct a θ_1^i and θ_3^i for which the following **invariant** holds:

1. θ_1^i and θ_3^i are ground solutions to P under the saturated signature Σ_{sat}
2. $\theta_1^i e \cong \theta_1 e$ and $\theta_3^i e \cong \theta_3 e$
3. $\forall k. 1 \leq k \leq i \implies \theta_1^i \pi_k = \theta_3^i \pi_k$
4. $\forall k. i + 1 \leq k \leq n \implies \theta_1^i \pi_k = \theta_1 \pi_k \wedge \theta_3^i \pi_k = \theta_3 \pi_k$
5. $\forall \pi \in P. \exists k. \pi_k = \pi \implies \theta_1^i \pi = \theta_3^i \pi$

Base case: if $n = 0$, the two solutions do not differ and the lemma holds trivially. Now suppose that the invariant holds for i iterations. We show how to construct θ_1^{i+1} and θ_3^{i+1} .

Let us consider the constraint $\pi_{i+1} = (m_L, m_R) \triangleright m_U$. From its definition we know that $\text{ftv}(\pi_{i+1}) \cap \bar{\mu} \neq \{\}$, so we consider all positions where open variables may appear, and the values assigned to them by the substitutions θ_1^{i+1} and θ_3^{i+1} .

From (Topo- π) we have that if there is a π_k s.t. $m_U \in \text{lo-bnd}(\pi_k)$ then $i + 1 > k$. From condition (3) of the invariant, we have that $\theta_1^i \pi_k = \theta_3^i \pi_k$. So m_U is either a constant, or a variable that has the same value in θ_1^i and θ_3^i : $\theta_1^{i+1} m_U = \theta_3^{i+1} m_U$. If there is no such π_k , then we can draw the same conclusion since we are given that all open μ variables have an upper bound.

We consider the case when m_L and m_R are equal to some open variables μ_L and μ_R , resp., and have different values in θ_1 and θ_3 . From Condition (2) of Definition 6 we know that there is a pair m_L^* and m_R^* such that $(m_L^*, m_R^*) = \text{smooth}((\theta_1 \mu_L, \theta_1 \mu_R), (\theta_3 \mu_L, \theta_3 \mu_R))$. We note that if either m_L or m_R is a constant, m_L^* or m_R^* will be equal to the same constant, respectively, as *smooth* preserves constants by definition. So from here on we assume m_L is some variable μ_L and m_R is some variable μ_R .

Now we define θ_1^{i+1} and θ_3^{i+1} as follows:

$$\begin{aligned} \theta^{i+1} &= (\mu_L \mapsto m_L^*) (\mu_R \mapsto m_R^*) \\ \theta_1^{i+1} &= \theta^{i+1} \theta_1^i \setminus \{\mu_L, \mu_R\} \\ \theta_3^{i+1} &= \theta^{i+1} \theta_3^i \setminus \{\mu_L, \mu_R\} \end{aligned}$$

The last three conditions of the invariant follow directly from this construction. Next we show why θ_1^{i+1} and θ_3^{i+1} are ground solution, as in Definition 7. Since the co-domain of θ_1^{i+1} and θ_3^{i+1} is constructed from the co-domain of θ_1^i and θ_3^i , which are themselves ground solutions, and from the results of *smooth* function, we know that co-domain of θ_1^{i+1} and θ_3^{i+1} only contains ground

polymonads and types. Now we have to show that $\Sigma_{sat} \models \theta_1^{i+1} \pi$ and $\Sigma_{sat} \models \theta_3^{i+1} \pi$ for all $\pi \in P$.

From the induction hypothesis we know that θ_1^i and θ_3^i are ground; so $\Sigma_{sat} \models \theta_1^i \pi$ and $\Sigma_{sat} \models \theta_3^i \pi$ for all $\pi \in P$.

We consider different cases for all the constraints $\pi \in P$.

Sub-case $\pi = \pi_{i+1}$: The first property of lubs applies: since we have $\theta_1^i((m_L, m_R) \triangleright m_U)$ and

$\theta_3^i((m_L, m_R) \triangleright m_U)$ and

the invariant establishes that $\theta_1^i m_U = \theta_1^{i+1} m_U = \theta_3^i m_U = \theta_3^{i+1} m_U$, so we have

$\theta_1^{i+1}((m_L, m_R) \triangleright m_U)$ and

$\theta_3^{i+1}((m_L, m_R) \triangleright m_U)$.

Sub-case π such that *up-bnd*(π) \in *lo-bnds*(π_{i+1}): Second property of lubs: for all m, m' if $(m, m') \triangleright \theta_1^i m_L$ or $(m, m') \triangleright \theta_3^i m_L$ then $(m, m') \triangleright \theta^{i+1} m_L$. Similarly, for all m, m' if $(m, m') \triangleright \theta_1^i m_R$ or $(m, m') \triangleright \theta_3^i m_R$ then $(m, m') \triangleright \theta^{i+1} m_R$.

Sub-case Other $\pi \in P$: are unchanged, so they are still entailed.

So the solutions θ_1^{i+1} and θ_3^{i+1} are both ground. Now we prove that the second condition of the invariant holds for θ_1^{i+1} and θ_3^{i+1} . By construction, θ_1^{i+1} (and θ_3^{i+1} respectively) are local modifications of θ_1^i (and θ_3^i). By Theorem 17, we know that $\theta_1^i e = \theta_1^{i+1} e$ and $\theta_3^i e = \theta_3^{i+1} e$. By transitivity of $=$, $\theta_1^{i+1} e = \theta_1 e$ and $\theta_3^{i+1} e = \theta_3 e$.

At the end of the iterations, we have $\theta_1^n e = \theta_3^n e$. By transitivity, $\theta_1 e = \theta_1^n e = \theta_3^n e = \theta_3 e$. We combine this result with $\theta_3 e = \theta_2 e$, we get $\theta_1 e = \theta_2 e$. \square

D. Categorical foundations

In this section we give some preliminary details of our categorical study of polymonads. Whilst the main focus of this paper is the generalization of monads for programming, there has been considerable semantic work on generalizing monads; in particular, we mention Filinski's layered monads [9] and Atkey's parameterised monads [2].

Definition 18 (Polymonads). *Given a (cartesian) category, \mathbb{C} , a polymonad is given by two collections:*

1. A collection of endofunctors over \mathbb{C} :

$$\mathbf{T} = \{T_i : \mathbb{C} \rightarrow \mathbb{C} \mid i \in \{0..n\}\}$$

where we require each endofunctor T_i to come equipped with a strength $\tau_{A,B}^i : A \times T_i(B) \rightarrow T_i(A \times B)$ that satisfies the following two diagrams.

$$\begin{array}{ccc} T_i(A) & & \\ \text{snd} \uparrow & \swarrow T_i(\text{snd}) & \\ 1 \times T_i(A) & \xrightarrow{\tau_A^i} & T_i(1 \times A) \end{array}$$

$$\begin{array}{ccc} (A \times B) \times T_i(C) & \xrightarrow{\tau_{A \times B, C}^i} & T_i((A \times B) \times C) \\ \alpha \downarrow & & \downarrow T_i(\alpha) \\ A \times (B \times T_i(C)) & \xrightarrow{Id \times \tau_{B, C}^i} A \times T_i(B \times C) \xrightarrow{\tau_{A, B \times C}^i} & T_i(A \times (B \times C)) \end{array}$$

By convention in any given collection, T_0 is always the identity functor.

2. A collection of natural transformations:

$$\otimes_{\mathbf{T}} = \{\star_A^j : T_{j_2}(T_{j_1}(A)) \rightarrow T_{j_3}(A) \mid j \in \{0..m\}, \\ T_{j_1} \in \mathbf{T}, T_{j_2} \in \mathbf{T}, T_{j_3} \in \mathbf{T}\}$$

where the collection must satisfy the (strong) polynomad laws. Again, it is by convention that \star^0 is the identity natural transformation between identity functors, i.e. $\star_A^0 : T_0(T_0(A)) \rightarrow T_0(A)$.

A collection $\otimes_{\mathbf{T}}$ is said to be polynomadic if it satisfies the following.

- $\forall T \in \mathbf{T}. \exists \star^1 \in \otimes_{\mathbf{T}}$ such that $\star^1 : T(T_0(A)) \rightarrow T(A)$ and $\star^1 = Id$.
- $\forall \star^1 \in \otimes_{\mathbf{T}}$ such that $\star_A^1 : T_1(T_0(A)) \rightarrow T_2(A)$ then $\exists \star^2 \in \otimes_{\mathbf{T}}$ such that $\star_A^2 : T_0(T_1(A)) \rightarrow T_2(A)$ and $\star^1 = \star^2$, and vice versa.
- $\forall \star^1, \star^2 \in \otimes_{\mathbf{T}}$ such that $\star_A^1 : T_1(T_2(A)) \rightarrow T_3(A)$ and $\star_A^2 : T_1(T_2(A)) \rightarrow T_3(A)$ then $\star^1 = \star^2$.
- $\forall \star^1, \star^2, \star^3, \star^4 \in \otimes_{\mathbf{T}}$ if $\star^1 : T_2(T_1(A)) \rightarrow T_3(A)$, $\star^2 : T_0(T_1'(A)) \rightarrow T_1(A)$, $\star^3 : T_0(T_2'(A)) \rightarrow T_2(A)$ and $\star^4 : T_0(T_3(A)) \rightarrow T_3'(A)$ then $\exists \star^5 \in \otimes_{\mathbf{T}}$ such that $\star^5 : T_2'(T_1'(A)) \rightarrow T_3'(A)$.
- $\forall \star^1, \star^2, \star^3, \star^4 \in \otimes_{\mathbf{T}}$ such that $\star_A^1 : T_1(T_2(A)) \rightarrow T_4(A)$, $\star_A^2 : T_4(T_3(A)) \rightarrow T_5(A)$, $\star_A^3 : T_2(T_3(A)) \rightarrow T_6(A)$ and $\star_A^4 : T_1(T_6(A)) \rightarrow T_5(A)$, the following diagram commutes.

$$\begin{array}{ccc} T_1(T_2(T_3(A))) & \xrightarrow{\star_{T_3(A)}^1} & T_4(T_3(A)) \\ \downarrow T_1(\star_A^3) & & \downarrow \star_A^2 \\ T_1(T_6(A)) & \xrightarrow{\star_A^4} & T_5(A) \end{array}$$

A polynomadic collection $\otimes_{\mathbf{T}}$ is said to be strong if it in addition satisfies the following laws.

- $\forall \star_A^i : T_2(T_1(A)) \rightarrow T_3(A) \in \otimes_{\mathbf{T}}$ the following diagram commutes.

$$\begin{array}{ccccc} A \times T_3(B) & \xrightarrow{\tau_{A,B}^3} & T_3(A \times B) & & \\ \uparrow Id \times \star_B^i & & \uparrow \star_{A \times B}^i & & \\ A \times T_2(T_1(B)) & \xrightarrow{\tau_{A,T_1(B)}^2} & T_2(A \times T_1(B)) & \xrightarrow{\tau_{T_1(B)}^2} & T_2(T_1(A \times B)) \end{array}$$

- $\forall \star_A^i : T_0(T_0(A)) \rightarrow T_1(A) \in \otimes_{\mathbf{T}}$ the following diagram commutes.

$$\begin{array}{ccc} A \times B & & \\ \downarrow Id \times \star_A^i & \searrow \star_{A \times B}^i & \\ A \times T_1(B) & \xrightarrow{\tau_{A,B}^1} & T_1(A \times B) \end{array}$$

The definition is relatively straightforward. The polynomad law (a) requires that there is an identity natural transformation between every functor in the endofunctor set. Polynomad law (b) essentially ensures that the identity functor acts like the unit in the monoid where functor composition acts as the multiplication. Polynomad law (c) reflects the coherence restriction of our setting that there can only be only bind operation between any three functors. Polynomad law (d) is the analog of the composition closure law defined in Section 2.2.

Polynomad law (e) is the generalization of the associativity law for monads, and laws (f) and (g) are the generalizations of the strong monad laws.

It is interesting to note that we do not require the generalization of monad triangle laws. In fact, these are derivable.

Lemma 19. For every polynomad, the following diagrams commute.

- For $\star^1 : T_0(T_0(A)) \rightarrow T_2(A)$, $\star^2 : T_1(T_2(A)) \rightarrow T_3(A)$ and $\star^3 : T_0(T_1(A)) \rightarrow T_3(A)$

$$\begin{array}{ccc} T_1(A) & \xrightarrow{T_1(\star_A^1)} & T_1(T_2(A)) \\ & \searrow \star_A^3 & \downarrow \star_A^2 \\ & & T_3(A) \end{array}$$

- For $\star^1 : T_0(T_0(A)) \rightarrow T_1(A)$, $\star^2 : T_1(T_2(A)) \rightarrow T_3(A)$ and $\star^3 : T_0(T_2(A)) \rightarrow T_3(A)$

$$\begin{array}{ccc} T_2(A) & \xrightarrow{\star_{T_2(A)}^1} & T_1(T_2(A)) \\ & \searrow \star_A^3 & \downarrow \star_A^2 \\ & & T_3(A) \end{array}$$

Proof. 1. In the following diagram, the square commutes as it is an instance of the polynomad law (d). Laws (a) and (b) combine to enforce that $\star_A^4 : T_1(T_0(A)) \rightarrow T_1(A)$ both exists in the polynomad and is the identity natural transformation. Hence the left triangle commutes.

$$\begin{array}{ccc} T_1(A) & \xlongequal{\quad} & T_1(T_0(T_0(A))) \xrightarrow{T_1(\star_A^1)} T_1(T_2(A)) \\ & \searrow & \downarrow \star_{T_0(A)}^4 \quad \downarrow \star^2 \\ & & T_1(T_0(A)) \xrightarrow{\star_{T_0(A)}^4} T_3(A) \end{array}$$

- In the following diagram, the polynomad laws (a) and (b) combine to enforce that $\star_A^4 : T_0(T_3(A)) \rightarrow T_3(A)$ both exists in the polynomad and is the identity natural transformation. The square commutes as it is an instance of the polynomad law (d).

$$\begin{array}{ccc} T_2(A) & \xlongequal{\quad} & T_0(T_0(T_2)) \xrightarrow{\star_{T_2(A)}^1} T_1(T_2(A)) \\ & & \downarrow T_0(\star^3) \quad \downarrow \star_A^2 \\ & & T_0(T_3(A)) \xrightarrow{\star^4} T_3(A) \end{array}$$

□

Atkey [3] proposed an generalization of a strong monad where the functor is no longer an endofunctor but a functor $T : \mathbb{S}^{\text{op}} \times \mathbb{S} \times \mathbb{C} \rightarrow \mathbb{C}$, where \mathbb{S} is an additional category whose objects are denote states and morphisms denote state transitions.

Theorem 20. Every \mathbb{S} -parameterized monad (T, η, μ) on \mathbb{C} [3] is a polynomad.

By currying the functor $T : \mathbb{S}^{\text{op}} \times \mathbb{S} \rightarrow (\mathbb{C} \rightarrow \mathbb{C})$, the polynomad functor and natural transformation collections can be defined by enumeration over the objects in \mathbb{S} . Checking that the polynomad laws are satisfied is routine.