

# Polymonadic programming

Michael Hicks<sup>1</sup>   Gavin Bierman<sup>2</sup>   Nataliya Guts<sup>1</sup>   Daan Leijen<sup>2</sup>   Nikhil Swamy<sup>2</sup>

<sup>1</sup>University of Maryland, College Park   <sup>2</sup>Microsoft Research

mwh@cs.umd.edu, {gmb, daan, nswamy}@microsoft.com, nyuguts@gmail.com

Monads are a popular tool for the working functional programmer to structure effectful computations. This paper presents *polymonads*, a generalization of monads. Polymonads give the familiar monadic bind the more general type  $\forall a, b. \mathsf{L} \ a \rightarrow (a \rightarrow \mathsf{M} \ b) \rightarrow \mathsf{N} \ b$ , to compose computations with three different kinds of effects, rather than just one. Polymonads subsume monads and parameterized monads, and can express other constructions, including precise type-and-effect systems and information flow tracking; more generally, polymonads correspond to Tate’s *productoid* semantic model. We show how to equip a core language (called  $\lambda\mathsf{PM}$ ) with syntactic support for programming with polymonads. Type inference and elaboration in  $\lambda\mathsf{PM}$  allows programmers to write polymonadic code directly in an ML-like syntax—our algorithms compute principal types and produce elaborated programs wherein the binds appear explicitly. Furthermore, we prove that the elaboration is *coherent*: no matter which (type-correct) binds are chosen, the elaborated program’s semantics will be the same. Pleasingly, the inferred types are easy to read: the polymonad laws justify (sometimes dramatic) simplifications, but with no effect on a type’s generality.

## 1 Introduction

Since the time that Moggi first connected them to effectful computation [20], *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 [24], and functional reactivity [7, 4].

Monads (and morphisms between them) are not a panacea, however, and so researchers have proposed various extensions. Examples include Wadler and Thiemann’s [29] indexed monad for typing effectful computations; Filliâtre’s generalized monads [10]; Atkey’s parameterized monad [3], which has been used to encode disciplines like regions [17] and session types [23]; Devriese and Piessens’ [6] monad-like encodings for information flow controls; and many others. Oftentimes these extensions are needed to prove stronger properties about computations, for instance to prove the absence of information leaks or memory errors.

Unfortunately, these extensions do not enjoy the same status as monads in terms of language support. For example, the conveniences that Haskell provides for monadic programs (e.g., the `do` notation combined with type-class inference) do not apply to these extensions. One might imagine adding specialized support for each of these extensions on a case-by-case basis, but a unifying construction into which all of them, including normal monads, fit is clearly preferable.

This paper proposes just such a unifying construction, making several contributions. Our first contribution is the definition of a *polymonad*, a new way to structure effectful computations. Polymonads give the familiar monadic bind (having type  $\forall a, b. \mathsf{M} \ a \rightarrow (a \rightarrow \mathsf{M} \ b) \rightarrow \mathsf{M} \ b$ ) the more general type  $\forall a, b. \mathsf{L} \ a \rightarrow (a \rightarrow \mathsf{M} \ b) \rightarrow \mathsf{N} \ b$ . That is, a polymonadic bind can compose computations with three different types to a monadic bind’s one. Section 2 defines polymonads formally, along with the *polymonad laws*, which we prove are a generalization of the monad and morphism laws. To precisely characterize their expressiveness, we prove that polymonads correspond to Tate’s *productoids* [27] (Theorem 2), a recent

semantic model general enough to capture most known effect systems, including all the constructions listed above.<sup>1</sup>

Whereas Tate’s interest is in semantically modeling sequential compositions of effectful computations, our interest is in supporting practical programming in a higher-order language. Our second contribution is the definition of  $\lambda_{\text{PM}}$  (Section 3), an ML-like programming language well-suited to programming with polymonads. We work out several examples in  $\lambda_{\text{PM}}$ , including novel polymonadic constructions for stateful information flow tracking, contextual type and effect systems [21], and session types.

Our examples are made practical by  $\lambda_{\text{PM}}$ ’s support for type inference and elaboration, which allows programs to be written in a familiar ML-like notation while making no mention of the bind operators. Enabling this feature, our third contribution (Section 4) is an instantiation of Jones’ theory of qualified types [14] to  $\lambda_{\text{PM}}$ . In a manner similar to Haskell’s type class inference, we show that type inference for  $\lambda_{\text{PM}}$  computes *principal types* (Theorem 3). Our inference algorithm is equipped with an elaboration phase, which translates source terms by inserting binds where needed. We prove that elaboration is *coherent* (Theorem 10), meaning that when inference produces constraints that could have several solutions, when these solutions are applied to the elaborated terms the results will have equivalent semantics, thanks to the polymonad laws. This property allows us to do better than Haskell, which does not take such laws into account, and so needlessly rejects programs it thinks might be ambiguous. Moreover, as we show in Section 5, the polymonad laws allow us to dramatically simplify types, making them far easier to read without compromising their generality. A prototype implementation of  $\lambda_{\text{PM}}$  is available from the first author’s web page and has been used to check all the examples in the paper.

Put together, our work lays the foundation for providing practical support for advanced monadic programming idioms in typed, functional languages.

## 2 Polymonads

We begin by defining polymonads formally. We prove that a polymonad generalizes a collection of monads and morphisms among those monads. We also establish a correspondence between polymonads and productoids, placing our work on a semantic foundation that is known to be extremely general.

**Definition 1.** A *polymonad*  $(\mathcal{M}, \Sigma)$  consists of (1) a collection  $\mathcal{M}$  of unary type constructors, with a distinguished element  $\text{Id} \in \mathcal{M}$ , such that  $\text{Id } \tau = \tau$ , and (2) a collection,  $\Sigma$ , of bind operators such that the laws below hold, where  $(M, N) \triangleright P \triangleq \forall a. b. M a \rightarrow (a \rightarrow N b) \rightarrow P b$ .

For all  $M, N, P, Q, R, S, T, U \in \mathcal{M}$ .

**(Functor)**  $\exists b. b : (M, \text{Id}) \triangleright M \in \Sigma$  and  $b m (\lambda y. y) = m$

**(Paired morphisms)**  $\exists b_1 : (M, \text{Id}) \triangleright N \in \Sigma \iff \exists b_2 : (\text{Id}, M) \triangleright N \in \Sigma$  and  
 $\forall b_1 : (M, \text{Id}) \triangleright N, b_2 : (\text{Id}, M) \triangleright N. b_1 (f v) (\lambda y. y) = b_2 v f$

**(Diamond)**  $\exists P, b_1, b_2. \{b_1 : (M, N) \triangleright P, b_2 : (P, R) \triangleright T\} \subseteq \Sigma \iff$   
 $\exists S, b_3, b_4. \{b_3 : (N, R) \triangleright S, b_4 : (M, S) \triangleright T\} \subseteq \Sigma$

**(Associativity)**  $\forall b_1, b_2, b_3, b_4. \text{If}$   
 $\{b_1 : (M, N) \triangleright P, b_2 : (P, R) \triangleright T, b_3 : (N, R) \triangleright S, b_4 : (M, S) \triangleright T\} \subseteq \Sigma$   
 $\text{then } b_2 (b_1 m f) g = b_4 m (\lambda x. b_3 (f x) g)$

**(Closure)**  $\text{If } \exists b_1, b_2, b_3, b_4.$   
 $\{b_1 : (M, N) \triangleright P, b_2 : (S, \text{Id}) \triangleright M, b_3 : (T, \text{Id}) \triangleright N, b_4 : (P, \text{Id}) \triangleright U\} \subseteq \Sigma$   
 $\text{then } \exists b. b : (S, T) \triangleright U \in \Sigma$

<sup>1</sup>We discovered the same model concurrently with Tate and independently of him, though we have additionally developed supporting algorithms for (principal) type inference, (provably coherent) elaboration, and (generality-preserving) simplification. Nevertheless, our presentation here has benefited from conversations with him.

Definition 1 may look a little austere, but there is a simple refactoring that recovers the structure of functors and monad morphisms from a polymonad.<sup>2</sup> Given  $(\mathcal{M}, \Sigma)$ , we can easily construct the following sets:

$$\begin{aligned} (\text{Maps}) \quad M &= \{(\lambda f m. \text{bind } m \ f) : (a \rightarrow b) \rightarrow M \ a \rightarrow M \ b \mid \text{bind} : (M, \text{Id}) \triangleright M \in \Sigma\} \\ (\text{Units}) \quad U &= \{(\lambda x. \text{bind } x \ (\lambda y. y)) : a \rightarrow M \ a \mid \text{bind} : (\text{Id}, \text{Id}) \triangleright M \in \Sigma\} \\ (\text{Lifts}) \quad L &= \{(\lambda x. \text{bind } x \ (\lambda y. y)) : M \ a \rightarrow N \ a \mid \text{bind} : (M, \text{Id}) \triangleright N \in \Sigma\} \end{aligned}$$

It is fairly easy to show that the above structure satisfies generalizations of the familiar laws for monads and monad morphisms. For example, one can prove  $\text{bind} (\text{unit } e) \ f = f \ e$ , and  $\text{lift} (\text{unit}_1 \ e) = \text{unit}_2 \ e$  for all suitably typed  $\text{unit}_1, \text{unit}_2 \in U$ ,  $\text{lift} \in L$  and  $\text{bind} \in \Sigma$ .

With these intuitions in mind, one can see that the **Functor** law ensures that each  $M \in \Sigma$  has a map in  $M$ , as expected for monads. From the construction of  $L$ , one can see that a  $\text{bind} (M, \text{Id}) \triangleright N$  is just a morphism from  $M$  to  $N$ . Since this comes up quite often, we write  $M \hookrightarrow N$  as a shorthand for  $(M, \text{Id}) \triangleright N$ . The **Paired morphisms** law amounts to a coherence condition that all morphisms can be re-expressed as binds. The **Associativity** law is the familiar associativity law for monads generalized for both our more liberal typing for bind operators and for the fact that we have a *collection* of binds rather than a single bind. The **Diamond** law essentially guarantees a coherence property for associativity, namely that it is always possible to complete an application of **Associativity**. The **Closure** law ensures closure under composition of monad morphisms with binds, also for coherence.

It is easy to prove that every collection of monads and monad morphisms is also a polymonad. In fact, in Appendix A, we prove a stronger result that relates polymonads to Tate’s *productoids* [27].

**Theorem 2.** *Every polymonad gives rise to a productoid, and every productoid that contains an Id element and whose joins are closed with respect to the lifts, is a polymonad.*

Tate developed productoids as a categorical foundation for effectful computation. He demonstrates the expressive power of productoids by showing how they subsume other proposed extensions to monads [29, 8, 3]. This theorem shows polymonads can be soundly interpreted using productoids. Strictly speaking, productoids are more expressive than polymonads, since they do not, in general, need to have an *Id* element, and only satisfy a slightly weaker form of our **Closure** condition. However, these restrictions are mild, and certainly in categories that are Cartesian closed, these conditions are trivially met for all productoids. Thus, for programming purposes, polymonads and productoids have exactly the same expressive power. The development of the rest of this paper shows, for the first time, how to harness this expressive power in a higher-order programming language, tackling the problem of type inference, elaborating a program while inserting binds, and proving elaboration coherent.

### 3 Programming with polymonads

This section presents  $\lambda_{\text{PM}}$ , an ML-like language for programming with polymonads. We also present several examples that provide a flavor of programming in  $\lambda_{\text{PM}}$ . As such, we aim to keep our examples as simple as possible while still showcasing the broad applicability of polymonads. For a formal characterization of the expressiveness of polymonads, we appeal to Theorem 2.

<sup>2</sup>An online version of this paper provides an equivalent formulation of Definition 1 in terms of join operators instead of binds. It can be found here: <http://research.microsoft.com/en-us/um/people/nswamy/papers/polymonads.pdf>. The join-based definition is perhaps more natural for a reader with some familiarity with category theory; the bind-based version shown here is perhaps more familiar for a functional programmer.

<i>Signatures</i> ( $\mathcal{M}, \Sigma$ ):		$k$ -ary constructors	$\mathcal{M} ::= \cdot \mid M/k, \mathcal{M}$
		ground constructor	$M ::= M \bar{\tau}$
		bind set	$\Sigma ::= \cdot \mid b:s, \Sigma$
		bind specifications	$s ::= \forall \bar{a}. \Phi \Rightarrow (M_1, M_2) \triangleright M_3$
		theory constraints	$\Phi$
<i>Terms</i> :		values	$v ::= x \mid c \mid \lambda x. e$
		expressions	$e ::= v \mid e_1 e_2 \mid \text{let } x = e_1 \text{ in } e_2$ $\quad \mid \text{if } e \text{ then } e_1 \text{ else } e_2 \mid \text{letrec } f = v \text{ in } e$
<i>Types</i> :		monadic types	$m ::= M \mid \rho$
		value types	$\tau ::= a \mid T \bar{\tau} \mid \tau_1 \rightarrow m \tau_2$
		type schemes	$\sigma ::= \forall \bar{a} \bar{\rho}. P \Rightarrow \tau$
		bag of binds	$P ::= \cdot \mid \pi, P$
		bind type	$\pi ::= (m_1, m_2) \triangleright m_3$

Figure 1:  $\lambda_{\text{PM}}$ : Syntax for signatures, types, and terms

**Polymonadic signatures.** A  $\lambda_{\text{PM}}$  *polymonadic signature*  $(\mathcal{M}, \Sigma)$  (Figure 1) amends Definition 1 in two ways. Firstly, each element  $M$  of  $\mathcal{M}$  may be *type-indexed*—we write  $M/k$  to indicate that  $M$  is a  $(k+1)$ -ary type constructor (we sometimes omit  $k$  for brevity). For example, constructor  $W/1$  could represent an effectful computation so that  $W \varepsilon \tau$  characterizes computations of type  $\tau$  that have effect  $\varepsilon$ . Type indexed constructors (rather than large enumerations of non-indexed constructors) are critical for writing reusable code, e.g., so we can write functions like  $\text{app} : \forall a, b, \varepsilon. (a \rightarrow W \varepsilon b) \rightarrow a \rightarrow W \varepsilon b$ . We write  $M$  to denote *ground constructors*, which are monadic constructors applied to all their type indexes; e.g.,  $W \varepsilon$  is ground. Secondly, a bind set  $\Sigma$  is not specified intensionally as a set, but rather extensionally using a language of *theory constraints*  $\Phi$ . In particular,  $\Sigma$  is a list of mappings  $b:s$  where  $s$  contains a triple  $(M_1, M_2) \triangleright M_3$  along with constraints  $\Phi$ , which determine how the triple’s constructors may be instantiated. For example, a mapping  $\text{sube} : \forall \varepsilon_1, \varepsilon_2. \varepsilon_1 \subseteq \varepsilon_2 \Rightarrow (W \varepsilon_1, \text{ld}) \triangleright W \varepsilon_2$  specifies the set of binds involving type indexes  $\varepsilon_1, \varepsilon_2$  such that the theory constraint  $\varepsilon_1 \subseteq \varepsilon_2$  is satisfied.

$\lambda_{\text{PM}}$ ’s type system is parametric in the choice of theory constraints  $\Phi$ , which allows us to encode a variety of prior monad-like systems with  $\lambda_{\text{PM}}$ . To interpret a particular set of constraints,  $\lambda_{\text{PM}}$  requires a theory entailment relation  $\models$ . Elements of this relation, written  $\Sigma \models \pi \rightsquigarrow b; \theta$ , state that there exists  $b: \forall \bar{a}. \Phi \Rightarrow (M_1, M_2) \triangleright M_3$  in  $\Sigma$  and a substitution  $\theta'$  such that  $\theta\pi = \theta'(M_1, M_2) \triangleright M_3$ , and the constraints  $\theta'\Phi$  are satisfiable. Here,  $\theta$  is a substitution for the free (non-constant) variables in  $\pi$ , while  $\theta'$  is an instantiation of the abstracted variables in the bind specification. Thus, the interpretation of  $\Sigma$  is the following set of binds:  $\{b:\pi \mid \Sigma \models \pi \rightsquigarrow b; \cdot\}$ . Signature  $(\mathcal{M}, \Sigma)$  is a polymonad if this set satisfies the polymonad laws (where each ground constructor is treated distinctly).

Our intention is that type indices are *phantom*, meaning that they are used as a type-level representation of some property of the polymonad’s current state, but a polymonadic bind’s implementation does not depend on them. For example, we would expect that binds treat objects of type  $W \varepsilon \tau$  uniformly, for all  $\varepsilon$ ; different values of  $\varepsilon$  could statically prevent unsafe operations like double-frees or dangling pointer dereferences. Of course, a polymonad may include other constructors distinct from  $W$  whose bind operators could have a completely different semantics. For example, if an object has different states that would affect the semantics of binds, or if other effectful features like exceptions were to be modeled, the programmer can use a different constructor  $M$  for each such feature. As such, our requirement that the type indices are phantom does not curtail expressiveness.

*Signature:*

$$\begin{aligned} \mathcal{M} &= IST/2 \\ \Phi &::= l_1 \sqsubseteq l_2 \mid \Phi_1, \Phi_2 \\ \Sigma &= \text{bld} : \quad \text{ld} \hookrightarrow \text{ld}, \\ &\quad \text{unitIST} : \quad \forall p, l. \text{ld} \hookrightarrow IST\ p\ l, \\ &\quad \text{mapIST} : \quad \forall p_1, l_1, p_2, l_2. p_2 \sqsubseteq p_1, l_1 \sqsubseteq l_2 \Rightarrow \\ &\quad \quad \quad IST\ p_1\ l_1 \hookrightarrow IST\ p_2\ l_2, \\ &\quad \text{applIST} : \quad \forall p_1, l_1, p_2, l_2. p_2 \sqsubseteq p_1, l_1 \sqsubseteq l_2 \Rightarrow \\ &\quad \quad \quad (\text{ld}, IST\ p_1\ l_1) \triangleright IST\ p_2\ l_2, \\ &\quad \text{blIST} : \quad \forall p_1, l_1, p_2, l_2, p_3, l_3. \\ &\quad \quad \quad l_1 \sqsubseteq p_2, l_1 \sqsubseteq l_3, l_2 \sqsubseteq l_3, \\ &\quad \quad \quad p_3 \sqsubseteq p_1, p_3 \sqsubseteq p_2 \Rightarrow \\ &\quad \quad \quad (IST\ p_1\ l_1, IST\ p_2\ l_2) \triangleright IST\ p_3\ l_3 \end{aligned}$$

*Types and auxiliary functions:*

$$\begin{aligned} \tau &: \quad \dots \mid \text{intref } \tau \mid L \mid H \\ \text{read} &: \quad \forall l. \text{intref } l \rightarrow IST\ H\ l\ \text{int} \\ \text{write} &: \quad \forall l. \text{intref } l \rightarrow \text{int} \rightarrow IST\ l\ L\ () \end{aligned}$$

*Example program:*

```

let add_interest =  $\lambda$ savings.  $\lambda$ interest.
  let currinterest = read interest in
    if currinterest > 0 then
      let currbalance = read savings in
        let newbalance =
          currbalance + currinterest in
          write savings newbalance
    else ()

```

Figure 2: Polymonad *IST*, implementing stateful information flow control

**Terms and types.**  $\lambda_{\text{PM}}$ 's term language is standard.  $\lambda_{\text{PM}}$  programs do not explicitly reference binds, but are written in *direct style*, with implicit conversions between computations of type  $m\ \tau$  and their  $\tau$ -typed results. Type inference determines the bind operations to insert (or abstract) to type check a program.

To make inference feasible, we rely crucially on  $\lambda_{\text{PM}}$ 's call-by-value structure. Following our prior work on monadic programming for ML [26], we restrict the shape of types assignable to a  $\lambda_{\text{PM}}$  program by separating value types  $\tau$  from the types of polymonadic computations  $m\ \tau$ . Here, metavariable  $m$  may be either a ground constructor  $M$  or a polymonadic type variable  $\rho$ . 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  $\tau \rightarrow \text{ld } \tau'$ . We also include types  $T\ \bar{\tau}$  for fully applied type constructors, e.g., list *int*.

Programs can also be given type schemes  $\sigma$  that are polymorphic in their polymonads, e.g.,  $\forall a, b, \rho. (a \rightarrow \rho b) \rightarrow a \rightarrow \rho b$ . Here, the variable  $a$  ranges over value types  $\tau$ , while  $\rho$  ranges over ground constructors  $M$ . Type schemes may also be qualified by a set  $P$  of bind constraints  $\pi$ . For example,  $\forall \rho. (\rho, \text{ld}) \triangleright M \Rightarrow (\text{int} \rightarrow \rho\ \text{int}) \rightarrow M\ \text{int}$  is the type of a function that abstracts over a bind having shape  $(\rho, \text{ld}) \triangleright M$ . Notice that  $\pi$  triples may contain polymonadic type variables  $\rho$  while specification triples  $s \in \Sigma$  may not. Moreover,  $\Phi$  constraints never appear in  $\sigma$ , which is thus entirely independent of the choice of the theory.

### 3.1 Polymonadic information flow controls

Polymonads are appealing because they can express many interesting constructions as we now show.

Figure 2 presents a polymonad *IST*, which implements *stateful* information flow tracking [6, 25, 19, 5, 1]. The idea is that some program values are secret and some are public, and no information about the former should be learned by observing the latter—a property called noninterference [11]. In the setting of *IST*, we are worried about leaks via the heap.

Heap-resident storage cells are given type *intref*  $l$  where  $l$  is the secrecy label of the referenced cell. Labels  $l \in \{L, H\}$  form a lattice with order  $L \sqsubseteq H$ . A program is acceptable if data labeled  $H$  cannot flow, directly or indirectly, to computations or storage cells labeled  $L$ . In our polymonad implementation,  $L$  and  $H$  are just types  $T$  (but only ever serve as indexes), and the lattice ordering is implemented by theory constraints  $l_1 \sqsubseteq l_2$  for  $l_1, l_2 \in \{L, H\}$ .

The polymonadic constructor *IST*/2 uses secrecy labels for its type indexes. A computation with type *IST*  $p\ l\ \tau$  potentially writes to references labeled  $p$  and returns a  $\tau$ -result that has security label  $l$ ;

we call  $p$  the *write label* and  $l$  the *output label*. Function `read` reads a storage cell, producing a  $IST\ H\ l\ int$  computation—the second type index  $l$  matches that of  $l$ -labeled storage cell. Function `write` writes a storage cell, producing a  $IST\ l\ L\ ()$  computation—the first type index  $l$  matches the label of the written-to storage cell.  $H$  is the most permissive write label and so is used for the first index of `read`, while  $L$  is the most permissive output label and so is used for the second index of `write`.

Aside from the identity bind `bld`, implemented as `reverse apply`, there are four kinds of binds. Unit `unit!ST`  $p\ l$  lifts a normal term into an  $IST$  computation. Bind `map!ST`  $p\ l$  lifts a computation into a more permissive context (i.e.,  $p_2$  and  $l_2$  are at least as permissive as  $l_1$  and  $l_2$ ), and `applST`  $p\ l$  does likewise, and are implemented using `map!ST` as follows: `applST`  $p\ l = \lambda x. \lambda f. \text{map!ST } p\ l\ (f\ x)\ (\lambda x. x)$ . Finally, bind `b!ST` composes a computation  $IST\ p_1\ l_1\ \alpha$  with a function  $\alpha \rightarrow IST\ p_2\ l_2\ \beta$ . The constraints ensure safe information flow:  $l_1 \sqsubseteq p_2$  prevents the second computation from leaking information about its  $l_1$ -secure  $\alpha$ -typed argument into a reference cell that is less than  $l_1$ -secure. Dually, the constraints  $l_1 \sqsubseteq l_3$  and  $l_2 \sqsubseteq l_3$  ensure that the  $\beta$ -typed result of the composed computation is at least as secure as the results of each component. The final constraints  $p_3 \sqsubseteq p_1$  and  $p_3 \sqsubseteq p_2$  ensure that the write label of the composed computation is a lower bound of the labels of each component.

Proving  $(\mathcal{M}, \Sigma)$  satisfies the polymonad laws is straightforward. The functor and paired morphism laws are easy to prove. The diamond law is more tedious: we must consider all possible pairs of binds that compose. This reasoning involves consideration of the theory constraints as implementing a lattice, and so would work for any lattice of labels, not just  $H$  and  $L$ . In all, there were ten cases to consider. We prove the associativity law for the same ten cases. This proof is straightforward as the implementation of  $IST$  ignores the indexes: `read`, `write` and various binds are just as in a normal state monad, while the indexes serve only to prevent illegal flows. Finally, proving closure is relatively straightforward—we start with each possible bind shape and then consider correctly-shaped flows into its components; in all there were eleven cases.

**Example.** The lower right of Figure 2 shows an example use of  $IST$ . The `add_interest` function takes two reference cells, `savings` and `interest`, and modifies the former by adding to it the latter if it is non-negative.<sup>3</sup> Notice that expressions of type  $IST\ p\ l\ \tau$  are used as if they merely had type  $\tau$ —see the branch on `currinterest`, for example. The program is rewritten during type inference to insert, or abstract, the necessary binds so that the program type checks. This process results in the following type for `add_interest`:<sup>4</sup>

$$\forall \rho_6, \rho_{27}, a_1, a_2. P \Rightarrow \text{intref } a_1 \rightarrow \text{intref } a_2 \rightarrow \rho_{27}\ () \\ \text{where } P = (\text{Id}, \text{Id}) \triangleright \rho_6, (IST\ H\ a_1, IST\ a_1\ L) \triangleright \rho_6, (IST\ H\ a_2, \rho_6) \triangleright \rho_{27}$$

The rewritten version of `add_interest` starts with a sequence of  $\lambda$  abstractions, one for each of the bind constraints in  $P$ . If we imagine these are numbered  $b_1 \dots b_3$ , e.g., where  $b_1$  is a bind with type  $(\text{Id}, \text{Id}) \triangleright \rho_6$ , then the term looks as follows (notation  $\dots$  denotes code elided for simplicity):

```

λsavings. λinterest. b3 (read interest)
  (λ currinterest. if currinterest > 0 then (b2 ...) else (b1 ()) (λ z. z)))

```

In a program that calls `add_interest`, the bind constraints will be solved, and actual implementations of these binds will be passed in for each of  $b_i$  (using a kind of dictionary-passing style as with Haskell type classes).

Looking at the type of `add_interest` we can see how the constraints prevent improper information flows. In particular, if we tried to call `add_interest` with  $a_1 = L$  and  $a_2 = H$ , then the last two constraints become

<sup>3</sup>For ease of presentation, the program in Figure 2 uses **let** to sequence computations. This is not essential, e.g., we need not have **let**-bound `currbalance`.

<sup>4</sup>This and other example types were generated by our prototype implementation.

$\mathcal{M}$	$= CE/3$	Types and theory constraints:
$\Sigma$	$= \text{bld} : (\text{Id}, \text{Id}) \triangleright \text{Id},$ $\text{unitce} : (\text{Id}, \text{Id}) \triangleright CE \top \emptyset \top$ $\text{appce} : \forall \alpha_1, \alpha_2, \varepsilon_1, \varepsilon_2, \omega_1, \omega_2.$ $(\alpha_2 \subseteq \alpha_1, \varepsilon_1 \subseteq \varepsilon_2, \omega_2 \subseteq \omega_1) \Rightarrow$ $(\text{Id}, CE \alpha_1 \varepsilon_1 \omega_1) \triangleright CE \alpha_2 \varepsilon_2 \omega_2$ $\text{mapce} : \forall \alpha_1, \alpha_2, \varepsilon_1, \varepsilon_2, \omega_1, \omega_2.$ $(\alpha_2 \subseteq \alpha_1, \varepsilon_1 \subseteq \varepsilon_2, \omega_2 \subseteq \omega_1) \Rightarrow$ $(CE \alpha_1 \varepsilon_1 \omega_1, \text{Id}) \triangleright CE \alpha_2 \varepsilon_2 \omega_2$ $\text{bindce} : \forall \alpha_1, \varepsilon_1, \omega_1, \alpha_2, \varepsilon_2, \omega_2, \varepsilon_3.$ $\varepsilon_2 \cup \omega_2 = \omega_1, \varepsilon_1 \cup \alpha_1 = \alpha_2, \varepsilon_1 \cup \varepsilon_2 = \varepsilon_3 \Rightarrow$ $(CE \alpha_1 \varepsilon_1 \omega_1, CE \alpha_2 \varepsilon_2 \omega_2) \triangleright CE \alpha_1 \varepsilon_3 \omega_2$	$\tau ::= \dots \mid \{A_1\} \dots \{A_n\} \mid \emptyset \mid \top \mid \tau_1 \cup \tau_2$ $\Phi ::= \tau \subseteq \tau' \mid \tau = \tau' \mid \Phi, \Phi$
		Auxiliary functions:
		$\text{read} : \forall \alpha, \omega, r. \text{intref } r \rightarrow CE \alpha r \omega \text{ int}$ $\text{write} : \forall \alpha, \omega, r. \text{intref } r \rightarrow \text{int} \rightarrow CE \alpha r \omega ()$

Figure 3: Polymonad expressing contextual type and effect systems

$(IST\ H\ L, IST\ L\ L) \triangleright \rho_6, (IST\ H\ H, \rho_6) \triangleright \rho_{27}$ , and so we must instantiate  $\rho_6$  and  $\rho_{27}$  in a way allowed by the signature in Figure 2. While we can legally instantiate  $\rho_6 = IST\ L\ l_3$  for any  $l_3$  to solve the second constraint, there is then no possible instantiation of  $\rho_{27}$  that can solve the third constraint. After substituting for  $\rho_6$ , this constraint has the form  $(IST\ H\ H, IST\ L\ l_3) \triangleright \rho_{27}$ , but this form is unacceptable because the  $H$  output of the first computation could be leaked by the  $L$  side effect of the second computation. On the other hand, all other instantiations of  $a_1$  and  $a_2$  (e.g.,  $a_1 = H$  and  $a_2 = L$  to correspond to a secret savings account but a public interest rate) do have solutions and do not leak information. Having just discussed the latter two constraints, consider the first,  $(\text{Id}, \text{Id}) \triangleright \rho_6$ . This constraint is important because it says that  $\rho_6$  must have a unit, which is needed to properly type the else branch; units are not required of a polymonad in general.

The type given above for `add_interest` is not its principal type, but an *improved* one. As it turns out, the principal type is basically unreadable, with 19 bind constraints! Fortunately, Section 5 shows how some basic rules can greatly simplify types without reducing their applicability, as has been done above. Moreover, our coherence result (given in the next section) assures that the corresponding changes to the elaborated term do not depend on the particular simplifications: the polymonad laws ensure all such elaborations will have the same semantics.

### 3.2 Contextual type and effect systems

Wadler and Thiemann [29] showed how a monadic-style construct can be used to model type and effect systems. Polymonads can model standard effect systems, but more interestingly can be used to model *contextual effects* [21], which augment traditional effects with the notion of *prior* and *future* effects of an expression within a broader context. As an example, suppose we are using a language that partitions memory into *regions*  $R_1, \dots, R_n$  and reads/writes of references into region  $R$  have effect  $\{R\}$ . Then in the context of the program `read  $r_1$ ; read  $r_2$` , where  $r_1$  points into region  $R_1$  and  $r_2$  points into region  $R_2$ , the contextual effect of the subexpression `read  $r_1$`  would be the triple  $[\emptyset; \{R_1\}; \{R_2\}]$ : the prior effect is empty, the present effect is  $\{R_1\}$ , and the future effect is  $\{R_2\}$ .

Figure 3 models contextual effects as the polymonad  $CE\ \alpha\ \varepsilon\ \omega\ \tau$ , for the type of a computation with prior, present, and future effects  $\alpha$ ,  $\varepsilon$ , and  $\omega$ , respectively. 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. As an example source of effects, we include `read` and `write` functions on references into

$\mathcal{M}$	$= \text{Id}, A/2$	<i>Types:</i>
$\Sigma$	$= \text{bld} : (\text{Id}, \text{Id}) \triangleright \text{Id},$ $\text{mapA} : \forall p, r. (A\ p\ r, \text{Id}) \triangleright A\ p\ r,$ $\text{appA} : \forall p, r. (\text{Id}, A\ p\ r) \triangleright A\ p\ r,$ $\text{unitA} : \forall p. (\text{Id}, \text{Id}) \triangleright A\ p\ p,$ $\text{bindA} : \forall p, q, r. (A\ p\ q, A\ q\ r) \triangleright A\ p\ r$	$\tau ::= \dots \mid \text{send } \tau_1\ \tau_2 \mid \text{recv } \tau_1\ \tau_2 \mid \text{end}$
		<i>Auxiliary functions:</i>
		$\text{send} : \forall a, q. a \rightarrow A(\text{send } a\ q)\ q()$
		$\text{recv} : \forall a, q. () \rightarrow A(\text{recv } a\ q)\ q\ a$

Figure 4: Parameterized monad for session types, expressed as a polymonad

region sets  $r$ . The bind unitce ascribes a pure computation as having an empty effect and any prior and future effects. The binds appce and mapce express that it is safe to consider an additional effect for the current computation (the  $\varepsilon$ s are covariant), and fewer effects for the prior and future computations ( $\alpha$ s and  $\omega$ 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.

### 3.3 Parameterized monads, and session types

Finally, we show  $\lambda_{\text{PM}}$  can express Atkey’s parameterized monad [3], which has been used to encode disciplines like regions [17] and session types [23]. The type constructor  $A\ p\ q\ \tau$  can be thought of (informally) as the type of a computation producing a  $\tau$ -typed result, with a pre-condition  $p$  and a post-condition  $q$ .

As a concrete example, Figure 4 gives a polymonadic expression of Pucella and Tov’s notion of session types [23]. The type  $A\ p\ q\ \tau$  represents a computation involved in a two-party session which starts in protocol state  $p$  and completes in state  $q$ , returning a value of type  $\tau$ . The key element of the signature  $\Sigma$  is the  $\text{bindA}$ , which permits composing two computations where the first’s post-condition matches the second’s precondition. We use the type index  $\text{send } \tau\ q$  to denote a protocol state that requires a message of type  $\tau$  to be sent, and then transitions to  $q$ . Similarly, the type index  $\text{recv } \tau\ r$  denotes the protocol state in which once a message of type  $\tau$  is received, the protocol transitions to  $r$ . We also use the index  $\text{end}$  to denote the protocol end state. The signatures of two primitive operations for sending and receiving messages capture this behavior.

As an example, the following  $\lambda_{\text{PM}}$  program implements one side of a simple protocol that sends a message  $x$ , waits for an integer reply  $y$ , and returns  $y+1$ .

$\text{let go} = \lambda x. \text{let } \_ = \text{send } x \text{ in incr } (\text{recv } ())$   
 Simplified type:  $\forall a, b, q, \rho. (A(\text{send } a\ b)\ b, A(\text{recv } \text{int } q)\ q)) \triangleright \rho \Rightarrow (a \rightarrow \rho\ \text{int})$

There are no specific theory constraints for session types: constraints simply arise by unification and are solved as usual when instantiating the final program (e.g., to call  $\text{go } 0$ ).

## 4 Coherent type inference for $\lambda_{\text{PM}}$

This section defines our declarative type system for  $\lambda_{\text{PM}}$  and proves that type inference produces principal types, and that elaborated programs are coherent.

Figure 5 gives a syntax-directed type system, organized into two main judgments. The value-typing judgment  $P \mid \Gamma \vdash v : \tau \rightsquigarrow e$  types a value  $v$  in an environment  $\Gamma$  (binding variables  $x$  and constants  $c$  to



$$\begin{array}{c}
\boxed{P \models P'} \quad \frac{\forall \pi \in P'. \pi \in P \vee \pi \in \Sigma}{P \models P'} \quad (\text{TS-Entail}) \\
\\
\boxed{P \models \sigma \geq \tau \rightsquigarrow f} \quad \frac{\theta = [\bar{\tau}/\bar{a}][\bar{m}/\bar{p}] \quad P \models \theta P_1}{P \models (\forall \bar{a} \bar{p}. P_1 \Rightarrow \tau) \geq \theta \tau \rightsquigarrow \text{app}(\theta P_1)} \quad (\text{TS-Inst}) \\
\\
\boxed{P | \Gamma \vdash v : \tau \rightsquigarrow e} \quad \frac{v \in \{x, c\} \quad P \models \Gamma(v) \geq \tau \rightsquigarrow f}{P | \Gamma \vdash v : \tau \rightsquigarrow f v} \quad (\text{TS-XC}) \\
\\
\frac{P | \Gamma, x : \tau_1 \vdash e : m \tau_2 \rightsquigarrow e}{P | \Gamma \vdash \lambda x. e : \tau_1 \rightarrow m \tau_2 \rightsquigarrow \lambda x. e} \quad (\text{TS-Lam}) \\
\\
\boxed{P | \Gamma \vdash e : m \tau \rightsquigarrow e} \quad \frac{P | \Gamma \vdash v : \tau \rightsquigarrow e}{P, \text{ld} \hookrightarrow m | \Gamma \vdash v : m \tau \rightsquigarrow \text{b}_{\text{ld}, \text{ld}, m} e (\lambda x. x)} \quad (\text{TS-V}) \\
\\
\frac{P_1 | \Gamma, x : \tau \vdash v : \tau \rightsquigarrow e_1 \quad (\sigma, e_2) = \text{Gen}(\Gamma, P_1 \Rightarrow \tau, e_1) \quad P | \Gamma, x : \sigma \vdash e : m \tau' \rightsquigarrow e_3}{P | \Gamma \vdash \text{letrec } x = v \text{ in } e : m \tau' \rightsquigarrow \text{letrec } x = e_2 \text{ in } e_3} \quad (\text{TS-Rec}) \\
\\
\frac{P_1 | \Gamma \vdash v : \tau \rightsquigarrow e_1 \quad (\sigma, e_2) = \text{Gen}(\Gamma, P_1 \Rightarrow \tau, e_1) \quad P | \Gamma, x : \sigma \vdash e : m \tau' \rightsquigarrow e_3}{P | \Gamma \vdash \text{let } x = v \text{ in } e : m \tau' \rightsquigarrow \text{let } x = e_2 \text{ in } e_3} \quad (\text{TS-Let}) \\
\\
\frac{P | \Gamma \vdash e_1 : m_1 \tau_1 \rightsquigarrow e_1 \quad P | \Gamma, x : \tau_1 \vdash e_2 : m_2 \tau_2 \rightsquigarrow e_2 \quad e_1 \neq v \quad P \models (m_1, m_2) \triangleright m_3}{P | \Gamma \vdash \text{let } x = e_1 \text{ in } e_2 : m_3 \tau_2 \rightsquigarrow \text{b}_{m_1, m_2, m_3} e_1 (\lambda x. e_2)} \quad (\text{TS-Do}) \\
\\
\frac{P | \Gamma \vdash e_1 : m_1 (\tau_2 \rightarrow m_3 \tau) \rightsquigarrow e_1 \quad P | \Gamma \vdash e_2 : m_2 \tau_2 \rightsquigarrow e_2 \quad P \models (m_1, m_4) \triangleright m_5 \quad P \models (m_2, m_3) \triangleright m_4}{P | \Gamma \vdash e_1 e_2 : m_5 \tau \rightsquigarrow \text{b}_{m_1, m_4, m_5} e_1 (\text{b}_{m_2, m_3, m_4} e_2)} \quad (\text{TS-App}) \\
\\
\frac{P | \Gamma \vdash e_1 : m_1 \text{bool} \rightsquigarrow e_1 \quad P | \Gamma \vdash e_2 : m_2 \tau \rightsquigarrow e_2 \quad P | \Gamma \vdash e_3 : m_3 \tau \rightsquigarrow e_3 \quad P \models m_2 \hookrightarrow m, m_3 \hookrightarrow m, (m_1, m) \triangleright m'}{P | \Gamma \vdash \text{if } e_1 \text{ then } e_2 \text{ else } e_3 : m' \tau \rightsquigarrow \text{b}_{m_1, m, m'} e_1 (\lambda b. \text{if } b \text{ then } \text{b}_{m_2, \text{ld}, m} e_2 (\lambda x. x) \text{ else } \text{b}_{m_3, \text{ld}, m} e_3 (\lambda x. x))} \quad (\text{TS-If}) \\
\\
\begin{array}{ll}
\text{Gen}(\Gamma, P \Rightarrow \tau, e) & = (\forall (\text{ftv}(P \Rightarrow \tau) \setminus \text{ftv}(\Gamma)). P \Rightarrow \tau, \text{abs}(P, e)) \\
\text{abs}((m_1, m_2) \triangleright m_3, P), e & = \lambda \text{b}_{m_1, m_2, m_3}. \text{abs}(P, e) \\
\text{abs}(\cdot, e) & = e \\
\text{app}(P, (m_1, m_2) \triangleright m_3)) & = \lambda f. \text{app}(P) (f \text{b}_{m_1, m_2, m_3}) \\
\text{app}(\cdot) & = \lambda x. x
\end{array}
\end{array}$$

Figure 5: Syntax-directed type rules for  $\lambda_{\text{PM}}$ , where  $\Sigma$  is an implicit parameter.

type schemes) at the type  $\tau$ , provided the constraints  $P$  are satisfiable. Moreover, it *elaborates* the value  $v$  into a lambda term  $e$  that explicitly contains binds, lifts, and evidence passing (as shown in Section 3.1). However, note that the elaboration is independent and we can read just the typing rules by ignoring the elaborated terms. The expression-typing judgment  $P | \Gamma \vdash e : m \tau \rightsquigarrow e$  is similar, except that it yields a computation type. Constraint satisfiability  $P \models P'$ , defined in the figure, states that  $P'$  is satisfiable under the hypothesis  $P$  if  $P' \subseteq P \cup \Sigma$  where we consider  $\pi \in \Sigma$  if and only if  $\Sigma \models \pi \rightsquigarrow b; \cdot$  (for some  $b$ ).

The rule (TS-XC) types a variable or constant at an instance of its type scheme in the environment.

$$\begin{array}{ll}
\llbracket x \rrbracket^* & = x \\
\llbracket c \rrbracket^* & = c \\
\llbracket \lambda x. e \rrbracket^* & = \lambda x. \llbracket e \rrbracket \\
\\ 
\llbracket v \rrbracket & = \text{ret } \llbracket v \rrbracket^* \\
\llbracket e_1 e_2 \rrbracket & = \text{app } \llbracket e_1 \rrbracket \llbracket e_2 \rrbracket \\
\llbracket \text{let } x = v \text{ in } e \rrbracket & = \text{let } x = \llbracket v \rrbracket^* \text{ in } \llbracket e \rrbracket \\
\llbracket \text{let } x = e_1 \text{ in } e_2 \rrbracket & = \text{do } \llbracket e_1 \rrbracket \llbracket \lambda x. e_2 \rrbracket^* \quad (\text{when } e_1 \neq v) \\
\llbracket \text{if } e_1 \text{ then } e_2 \text{ else } e_3 \rrbracket & = \text{cond } \llbracket e_1 \rrbracket \lambda(). \llbracket e_2 \rrbracket \lambda(). \llbracket e_3 \rrbracket \\
\llbracket \text{letrec } f = v \text{ in } e \rrbracket & = \text{letrec } f = \llbracket v \rrbracket^* \text{ in } \llbracket e \rrbracket \\
\\ 
\text{ret} & : \forall \alpha \rho. (\text{Id} \hookrightarrow \rho) \Rightarrow \alpha \rightarrow \rho \alpha \\
\text{do} & : \forall \alpha \beta \rho_1 \rho_2 \rho. ((\rho_1, \rho_2) \triangleright \rho) \Rightarrow \rho_1 \alpha \rightarrow (\alpha \rightarrow \rho_2 \beta) \rightarrow \rho \beta \\
\text{app} & : \forall \alpha \beta \rho_1 \rho_2 \rho_3 \rho_4 \rho. ((\rho_1, \rho_4) \triangleright \rho, (\rho_2, \rho_3) \triangleright \rho_4) \Rightarrow \rho_1 (\alpha \rightarrow \rho_3 \beta) \rightarrow \rho_2 \alpha \rightarrow \rho \beta \\
\text{cond} & : \forall \alpha \rho_1 \rho_2 \rho_3 \rho \rho'. (\rho_2 \hookrightarrow \rho, \rho_3 \hookrightarrow \rho, (\rho_1, \rho) \triangleright \rho') \\
& \quad \Rightarrow \rho_1 \text{ bool} \rightarrow ((\rightarrow \rho_2 \alpha) \rightarrow ((\rightarrow \rho_3 \alpha) \rightarrow \rho' \alpha)
\end{array}$$

Figure 6: Type inference for  $\lambda_{\text{PM}}$  via elaboration to OML

The instance relation for type schemes  $P \models \sigma \geq \tau \rightsquigarrow f$  is standard: it instantiates the bound variables, and checks that the abstracted constraints are entailed by the hypothesis  $P$ . The elaborated  $f$  term supplies the instantiated evidence using the `app` form. The rule (TS-Lam) is straightforward where the bound variable is given a value type and the body a computation type.

The rule (TS-V) allows a value  $v : \tau$  to be used as an expression by lifting it to a computation type  $m \tau$ , so long as there exists a morphism (or unit) from the `Id` functor to  $m$ . The elaborated term uses  $b_{\text{Id}, \text{Id}, m}$  to lift explicitly to monad  $m$ . Note that for evidence we make up names ( $b_{\text{Id}, \text{Id}, m}$ ) based on the constraint  $(\text{Id} \hookrightarrow m)$ . This simplifies our presentation but an implementation would name each constraint explicitly [16]. We use the name  $b_{m_1, \text{Id}, m_2}$  for morphism constraints  $m_1 \hookrightarrow m_2$ , and use  $b_{m_1, m_2, m_3}$  for general bind constraints  $(m_1, m_2) \triangleright m_3$ .

(TS-Rec) types a recursive let-binding by typing the definition  $v$  at the same (mono-)type as the `letrec`-bound variable  $f$ . When typing the body  $e$ , we generalize the type of  $f$  using a standard generalization function  $\text{Gen}(\Gamma, P \Rightarrow \tau, e)$ , which closes the type relative to  $\Gamma$  by generalizing over its free type variables. However, in contrast to regular generalization, we return both a generalized type, as well as an elaboration of  $e$  that takes all generalized constraints as explicit evidence parameters (as defined by rule `abs`). (TS-Let) is similar, although somewhat simpler since there is no recursion involved.

(TS-Do) is best understood by looking at its elaboration: since we are in a call-by-value setting, we interpret a **let**-binding as forcing and sequencing two computations using a single bind where  $e_1$  is typed monomorphically.

(TS-App) is similar to (TS-Do), where, again, since we use call-by-value, in the elaboration we sequence the function and its argument using two bind operators, and then apply the function. (TS-If) is also similar, since we sequence the expression  $e$  in the guard with the branches. As usual, we require the branches to have the same type. This is achieved by generating morphism constraints,  $m_2 \hookrightarrow m$  and  $m_3 \hookrightarrow m$  to coerce the type of each branch to a functor  $m$  before sequencing it with the guard expression.

## 4.1 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 Overloaded ML (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. OML's type inference algorithm is known to enjoy principal types, so a corollary of our translation is that principal types exist for our system too.

We encode terms in our language into OML as shown in Figure 6. We rely on four primitive OML terms that force the typing of the terms to generate the same constraints as our type system does: `ret` for lifting a pure term, `do` for typing a `do`-binding, `app` for typing an application, and `cond` for conditionals. Using these primitives, we encode values and expressions of our system into OML.

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

**Theorem 3** (Encoding to OML is sound and complete).

**Soundness:** *Whenever  $P|\Gamma \vdash v : \tau$  we have  $P|\Gamma \vdash_{\text{OML}} \llbracket v \rrbracket^* : \tau$ . Similarly, whenever  $P|\Gamma \vdash e : m \tau$  then we have  $P|\Gamma \vdash_{\text{OML}} \llbracket e \rrbracket : m \tau$ .*

**Completeness:** *Whenever  $P|\Gamma \vdash_{\text{OML}} \llbracket v \rrbracket^* : \tau$ , then we have  $P|\Gamma \vdash v : \tau$ . Similarly, whenever  $P|\Gamma \vdash_{\text{OML}} \llbracket e \rrbracket : m \tau$ , then we have  $P|\Gamma \vdash e : m \tau$ .*

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, our system admits principal types via the general-purpose OML type inference algorithm.

## 4.2 Ambiguity

Seeing the previous OML translation, one might think we could directly translate our programs into Haskell since Haskell uses OML style type inference. Unfortunately, in practice, Haskell would reject many useful programs. In particular, Haskell rejects as ambiguous any term whose type  $\forall \tilde{\alpha}. P \Rightarrow \tau$  includes a variable  $\alpha$  that occurs free in  $P$  but not in  $\tau$ ;<sup>5</sup> we call such type variables *open*. Haskell, in its generality, must reject such terms since the instantiation of an open variable can have operational effect, while at the same time, since the variable does not appear in  $\tau$ , the instantiation for it can never be uniquely determined by the context in which the term is used. A common example is the term `show . read` with the type  $(\text{Show } a, \text{Read } a) \Rightarrow \text{String} \rightarrow \text{String}$ , where  $a$  is open. Depending on the instantiation of  $a$ , the term may parse and show integers, or doubles, etc.

Rejecting all types that contain open variables works well for type classes, but it would be unacceptable for  $\lambda_{\text{PM}}$ . Many simple terms have principal types with open variables. For example, the term  $\lambda f. \lambda x. f x$  has type  $\forall ab\rho_1\rho_2\rho_3. ((\text{Id}, \rho_1) \triangleright \rho_2, (\text{Id}, \rho_2) \triangleright \rho_3) \Rightarrow (a \rightarrow \rho_1 b) \rightarrow \alpha \rightarrow \rho_3 b$  where type variable  $\rho_2$  is open.

In the special case where there is only one polymonadic constructor available when typing the program, the coherence problem is moot, e.g., say, if the whole program were to only be typed using only the *IST* polymonad of Section 3.1. However, recall that polymonads generalize monads and morphisms, for which there can be coherence issues (as is well known), so polymonads must address them. As an example, imagine combining our *IST* polymonad (which generalizes the state monad) with an exception monad `Exn`, resulting in an *ISTExn* polymonad. Then, an improperly coded `bind` that composed *IST* with `Exn` could sometimes reset the heap, and sometimes not (a similar example is provided by Filinski [9]).

<sup>5</sup>The actual ambiguity rule in Haskell is more involved due to functional dependencies and type families but that does not affect our results.

A major contribution of this paper is that for binds that satisfy the polymonad laws, we need not reject all types with open variables. In particular, by appealing to the polymonadic laws, we can prove that programs with open type variables in bind constraints are indeed unambiguous. Even if there are many possible instantiations, the semantics of each instantiation is equivalent, enabling us to solve polymonadic constraints much more aggressively. This coherence result is at the essence of making programming with polymonads practical.

### 4.3 Coherence

The main result of this section (Theorem 10) establishes that for a certain class of polymonads, the ambiguity check of OML can be weakened to accept more programs while still ensuring that programs are coherent. Thus, for this class of polymonads, programmers can reliably view our syntax-directed system as a specification without being concerned with the details of how the type inference algorithm is implemented or how programs are elaborated.

The proof of Theorem 10 is a little technical—the following roadmap summarizes the structure of the development.

- We define the class of *principal* polymonads for which unambiguous typing derivations are coherent. All polymonads that we know of are principal.
- Given  $P \mid \Gamma \vdash e : t \rightsquigarrow e$  (with  $t \in \{\tau, m \tau\}$ ), the predicate  $\text{unambiguous}(P, \Gamma, t)$  characterizes when the derivation is unambiguous. This notion requires interpreting  $P$  as a graph  $G_P$ , and ensuring (roughly) that all open variables in  $P$  have non-zero in/out-degree in  $G_P$ .
- A *solution*  $S$  to a constraint graph with respect to a polymonad  $(\mathcal{M}, \Sigma)$  is an assignment of ground polymonad constructors  $M \in \mathcal{M}$  to the variables in the graph such that each instantiated constraint is present in  $\Sigma$ . We give an equivalence relation on solutions such that  $S_1 \cong S_2$  if they differ only on the assignment to open variables in a manner where the composition of binds still computes the same function according to the polymonad laws.
- Finally, given  $P \mid \Gamma \vdash e : t \rightsquigarrow e$  and  $\text{unambiguous}(P, \Gamma, t)$ , we prove that all solutions to  $P$  that agree on the free variables of  $\Gamma$  and  $t$  are in the same equivalence class.

While Theorem 10 enables our type system to be used in practice, this result is not the most powerful theorem one can imagine. Ideally, one might like a theorem of the form  $P \mid \Gamma \vdash e : t \rightsquigarrow e$  and  $P' \mid \Gamma \vdash e : t \rightsquigarrow e'$  implies  $e$  is extensionally equal to  $e'$ , given that both  $P$  and  $P'$  are satisfiable. While we conjecture that this result is true, a proof of this property is out of our reach, at present. There are at least two difficulties. First, a coherence result of this form is unknown for qualified type systems in a call-by-value setting. In an unpublished paper, Jones [15] proves a coherence result for OML, but his technique only applies to call-by-name programs. Jones also does not consider reasoning about coherence based on an equational theory for the evidence functions (these functions correspond to our binds). So, proving the ideal coherence theorem would require both generalizing Jones' approach to call-by-value and then extending it with support for equational reasoning about evidence. In the meantime, Theorem 10 provides good assurance and lays the foundation for future work in this direction.

**Defining and analyzing principality.** We introduce a notion of principal polymonads that corresponds to Tate's "principalled productoids." Informally, in a principal polymonad, if there is more than one way to combine pairs of computations in the set  $F$  (e.g.,  $(M, M') \triangleright M_1$  and  $(M, M') \triangleright M_2$ ), then there must be a "best" way to combine them. This best way is called the principal join of  $F$ , and all other ways

to combine the functors are related to the principal join by morphisms. All the polymonadic libraries we have encountered so far are principal polymonads. It is worth emphasizing that principality does not correspond to functional dependency—it is perfectly reasonable to combine  $M$  and  $M'$  in multiple ways, and indeed, for applications like sub-effecting, this expressiveness is important. We only require that there be an ordering among the choices. In the definition below, we take  $\downarrow \mathcal{M}$  to be set of ground instances of all constructors in  $\mathcal{M}$ .

**Definition 4** (Principal polymonad). *A polymonad  $(\mathcal{M}, \Sigma)$  is a principal polymonad if and only if for any set  $F \subseteq \downarrow \mathcal{M}^2$ , and any  $\{M_1, M_2\} \subseteq \downarrow \mathcal{M}$  such  $\{(M, M') \triangleright M_1 \mid (M, M') \in F\} \subseteq \Sigma$  and  $\{(M, M') \triangleright M_2 \mid (M, M') \in F\} \subseteq \Sigma$ , then there exists  $\hat{M} \in \downarrow \mathcal{M}$  such that  $\{\hat{M} \hookrightarrow M_1, \hat{M} \hookrightarrow M_2\} \subseteq \Sigma$ , and  $\{(M, M') \triangleright \hat{M} \mid (M, M') \in F\} \subseteq \Sigma$ . We call  $\hat{M}$  the principal join of  $F$  and write it as  $\sqcup F$*

**Definition 5** (Graph-view of a constraint-bag  $P$ ). *A graph-view  $G_P = (V, A, E_{\triangleright}, E_{eq})$  of a constraint-bag  $P$  is a graph consisting of a set of vertices  $V$ , a vertex assignment  $A : V \rightarrow m$ , a set of directed edges  $E_{\triangleright}$ , and a set of undirected edges  $E_{eq}$ , where:*

- $V = \{\pi.0, \pi.1, \pi.2 \mid \pi \in P\}$ , i.e., each constraint contributes three vertices.
- $A(\pi.i) = m_i$  when  $\pi = (m_0, m_1) \triangleright m_2$ , for all  $\pi.i \in V$
- $E_{\triangleright} = \{(\pi.0, \pi.2), (\pi.1, \pi.2) \mid \pi \in P\}$
- $E_{eq} = \{(v, v') \mid v, v' \in V \wedge v \neq v' \wedge \exists \rho. \rho = A(v) = A(v')\}$

**Notation** We use  $v$  in this section to stand for a graph vertex, rather than a value in a program. We also make use of a pictorial notation for graph views, distinguishing the two flavors of edges in a graph. Each constraint  $\pi \in P$  induces two edges in  $E_{\triangleright}$ . These edges are drawn with solid lines, with a triangle for orientation. Unification constraints arise from correlated variable occurrences in multiple constraints—we depict these with double dotted lines. For example, the pair of constraints  $m_1 \searrow \triangleright \rho$  and  $m_2 \searrow \triangleright \rho'$  contributes four unification edges, two for  $\rho$  and two for  $\rho'$ . We show its graph view alongside.

Unification constraints reflect the dataflow in a program. Referring back to Figure 5, in a principal derivation using (TS-App), correlated occurrences of unification variables for  $m_4$  in the constraints indicate how the two binds operators compose. The following definition captures this dataflow and shows how to interpret the composition of bind constraints using unification edges as a lambda term (in the expected way).<sup>6</sup>

**Definition 6** (Functional view of a flow edge). *Given a constraint graph  $G = (V, A, E_{\triangleright}, E_{eq})$ , an edge  $\eta = (\pi.2, \pi'.i) \in E_{eq}$ , where  $i \in \{0, 1\}$  and  $\pi \neq \pi'$  is called a flow edge. The flow edge  $\eta$  has a functional interpretation  $F_G(\eta)$  defined as follows:*

$$\begin{aligned} \text{If } i = 0, \quad F_G(\eta) &= \lambda(x:A(\pi.0) \ a) \ (y:a \rightarrow A(\pi.1) \ b) \ (z:b \rightarrow A(\pi'.1) \ c). \\ &\quad \text{bind}_{A(\pi'.0), A(\pi'.1), A(\pi'.2)}(\text{bind}_{A(\pi.0), A(\pi.1), A(\pi.2)} \ x \ y) \ z \\ \text{If } i = 1, \quad F_G(\eta) &= \lambda(x:A(\pi'.0) \ a) \ (y:a \rightarrow A(\pi.0) \ b) \ (z:b \rightarrow A(\pi.1) \ c). \\ &\quad \text{bind}_{A(\pi'.0), A(\pi'.1), A(\pi'.2)} \ x \ (\lambda a. \text{bind}_{A(\pi.0), A(\pi.1), A(\pi.2)} (y \ a) \ z) \end{aligned}$$

We can now define our ambiguity check—a graph is unambiguous if it contains a sub-graph that has no cyclic dataflows, and where open variables only occur as intermediate variables in a sequence of binds.

<sup>6</sup>Note, for the purposes of our coherence argument, unification constraints between value-type variables  $a$  are irrelevant. Such variables may occur in two kinds of contexts. First, they may constrain some value type in the program, but these do not depend on the solutions to polymonadic constraints. Second, they may constrain some index of a polymonadic constructor; but, as mentioned previously, these indices are phantom and do not influence the semantics of elaborated terms.

**Definition 7** (Unambiguous constraints). *Given  $G_P = (V, A, E_{\triangleright}, E_{eq})$ , the predicate  $\text{unambiguous}(P, \Gamma, t)$  holds if and only if there exists  $E'_{eq} \subseteq E_{eq}$ , such that in the graph  $G' = (V, A, E_{\triangleright}, E'_{eq})$  all of the following are true.*

1. *For all  $\pi \in P$ , there is no path from  $\pi.2$  to  $\pi.0$  or  $\pi.1$ .*
2. *For all  $v \in V$ , if  $A(v) \in \text{ftv}(P) \setminus \text{ftv}(\Gamma, t)$ , then there exists a flow edge that connects to  $v$ .*

We call  $G'$  a core of  $G_P$ .

**Definition 8** (Solution to a constraint graph). *For a polymonadic signature  $(\mathcal{M}, \Sigma)$ , a solution to a constraint graph  $G = (V, A, E_{\triangleright}, E_{eq})$ , is a vertex assignment  $S : V \rightarrow \mathcal{M}$  such that all of the following are true.*

1. *For all  $v \in V$ , if  $A(v) \in \mathcal{M}$  then  $S(v) = A(v)$*
2. *For all  $(v_1, v_2) \in E_{eq}$ ,  $S(v_1) = S(v_2)$ .*
3. *For all  $\{(\pi.0, \pi.2), (\pi.1, \pi.2)\} \subseteq E_{\triangleright}$ ,  $(S(\pi.0), S(\pi.1)) \triangleright S(\pi.2) \in \Sigma$ .*

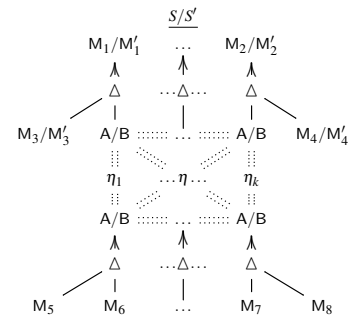
We say that two solutions  $S_1$  and  $S_2$  to  $G$  agree on  $\rho$  if for all vertices  $v \in V$  such that  $A(v) = \rho$ ,  $S_1(v) = S_2(v)$ .

Now we define  $\cong_R$ , a notion of equivalence of two solutions which captures the idea that the differences in the solutions are only to the internal open variables while not impacting the overall function computed by the binds in a constraint. It is easy to check that  $\cong_R$  is an equivalence relation.

**Definition 9** (Equivalence of solutions). *Given a polymonad  $(\mathcal{M}, \Sigma)$  and constraint graph  $G = (V, A, E_{\triangleright}, E_{eq})$ , two solutions  $S_1$  and  $S_2$  to  $G$  are equivalent with respect to a set of variables  $R$  (denoted  $S_1 \cong_R S_2$ ) if and only if  $S_1$  and  $S_2$  agree on all  $\rho \in R$  and for each vertex  $v \in V$  such that  $S_1(v) \neq S_2(v)$  for all flow edges  $\eta$  incident on  $v$ ,  $F_{G_1}(\eta) = F_{G_2}(\eta)$ , where  $G_i = (V, S_i, E_{\triangleright}, E_{eq})$ .*

**Theorem 10** (Coherence). *For all principal polymonads, derivations  $P|\Gamma \vdash e : t \rightsquigarrow e$  such that  $\text{unambiguous}(P, \Gamma, t)$ , and for any two solutions  $S$  and  $S'$  to  $G_P$  that agree on  $R = \text{ftv}(\Gamma, t)$ , we have  $S \cong_R S'$ .*

(Sketch; full version in appendix) The main idea is to show that all solutions in the core of  $G_P$  are in the same equivalence class (the solutions to the core include  $S$  and  $S'$ ). The proof proceeds by induction on the number of vertices at which  $S$  and  $S'$  differ. For the main induction step, we take vertices in topological order, considering the least (in the order) set of vertices  $Q$ , all related by unification constraints, and whose assignment in  $S$  is  $A$  and in  $S'$  is  $B$ , for some  $A \neq B$ . The vertices in  $Q$  are shown in the graph alongside, all connected to each other by double dotted lines (unification constraints), and their neighborhood is shown as well. Since vertices are considered in topological order, all the vertices below  $Q$  in the graph have the same assignment in  $S$  and in  $S'$ . We build solutions  $S_1$  and  $S'_1$  from  $S$  and  $S'$  respectively, that instead assign the principal join  $J = \sqcup\{(M_5, M_6), \dots, (M_7, M_8)\}$  to the vertices in  $Q$ , where  $S_1 \cong_R S'_1$  by the induction hypothesis. Finally, we prove  $S \cong_R S_1$  and  $S' \cong_R S'_1$  by showing that the functional interpretation of each of the flow edges  $\eta_i$  are equal according to the polymonad laws, and conclude  $S \cong_R S'$  by transitivity.



$$\begin{array}{c}
\begin{array}{c}
\pi = (\text{Id}, m) \triangleright \rho \vee \pi = (m, \text{Id}) \triangleright \rho \\
\rho \in \bar{\rho} \quad \text{flowsFrom}_{P, P'} \rho \neq \{\} \\
\text{flowsTo}_{P, P'} \rho = \{\}
\end{array}
\quad
\begin{array}{c}
\pi = (\text{Id}, \rho) \triangleright m \vee \pi = (\rho, \text{Id}) \triangleright m \\
\rho \in \bar{\rho} \quad \text{flowsFrom}_{P, P'} \rho = \{\} \\
\text{flowsTo}_{P, P'} \rho \neq \{\}
\end{array} \\
\text{S-}\Uparrow \frac{}{P, \pi, P' \xrightarrow{\text{simplify}(\bar{\rho})} \rho \mapsto m} \quad
\text{S-}\Downarrow \frac{}{P, \pi, P' \xrightarrow{\text{simplify}(\bar{\rho})} \rho \mapsto m} \\
\text{S-}\sqcup \frac{
\begin{array}{c}
F = \text{flowsTo}_P \rho \\
m \in F \Rightarrow m = M \\
\text{for some } M
\end{array}
}{P \xrightarrow{\text{simplify}(\bar{\rho})} \rho \mapsto \sqcup F} \quad
\frac{
\begin{array}{c}
P \xrightarrow{\text{simplify}(\bar{\rho})} \theta \\
\theta P \xrightarrow{\text{simplify}(\bar{\rho})} \theta'
\end{array}
}{P \xrightarrow{\text{simplify}(\bar{\rho})} \theta' \theta} \quad
\frac{}{P \xrightarrow{\text{simplify}(\bar{\rho})} .} \\
\text{where } \begin{array}{l} \text{flowsTo}_P \rho = \{ (m_1, m_2) \mid (m_1, m_2) \triangleright \rho \in P \} \\ \text{flowsFrom}_P \rho = \{ m \mid \exists m'. \pi \in P \wedge (\pi = (\rho, m') \triangleright m \vee \pi = (m', \rho) \triangleright m) \} \end{array}
\end{array}$$

Figure 7: Eliminating open variables in constraints

## 5 Simplification and solving

Before running a program, we must solve the constraints produced during type inference, and apply the appropriate evidence for these constraints in the elaborated program. We also perform *simplification* on constraints prior to generalization to make types easier to read, but without compromising their utility.

A simple syntactic transformation on constraints can make inferred types easier to read. For example, we can hide duplicate constraints, identity morphisms (which are trivially satisfiable), and constraints that are entailed by the signature. More substantially, we can find instantiations for open variables in a constraint set before generalizing a type (and at the top-level, before running a program). To do this, we introduce below a modified version of (TS-Let) (from Figure 5); a similar modification is possible for (TS-Rec).

$$\frac{
\begin{array}{c}
P_1 \mid \Gamma \vdash v : \tau \rightsquigarrow e_1 \\
P_1 \xrightarrow{\text{simplify}(\bar{\rho} \setminus \text{ftv}(\tau))} \theta \quad (\sigma, e_2) = \text{Gen}(\Gamma, \theta P_1 \Rightarrow \tau, e_1) \quad P \mid \Gamma, x : \sigma \vdash e : m \tau' \rightsquigarrow e_3
\end{array}
}{P \mid \Gamma \vdash \text{let } x = v \text{ in } e : m \tau' \rightsquigarrow \text{let } x = e_2 \text{ in } e_3}$$

This rule employs the judgment  $P \xrightarrow{\text{simplify}(\bar{\rho})} \theta$ , defined in Figure 7, to simplify constraints by eliminating some open variables in  $P$  (via the substitution  $\theta$ ) before type generalization. There are three main rules in the judgment, (S- $\Uparrow$ ), (S- $\Downarrow$ ) and (S- $\sqcup$ ), while the last two simply take the transitive closure.

Rule (S- $\Uparrow$ ) solves monad variable  $\rho$  with monad  $m$  if we have a constraint  $\pi = (\text{Id}, m) \triangleright \rho$ , where the only edges directed inwards to  $\rho$  are from  $\text{Id}$  and  $m$ , although there may be many out-edges from  $\rho$ . (The case where  $\pi = (m, \text{Id}) \triangleright \rho$  is symmetric.) Such a constraint can always be solved without loss of generality using an identity morphism, which, by the polyonad laws is guaranteed to exist. Moreover, by the closure law, any solution that chooses  $\rho = m'$ , for some  $m' \neq m$  could just as well have chosen  $\rho = m$ . Thus, this rule does not impact solvability of the constraints. Rule S- $\Downarrow$  follows similar reasoning in the reverse direction. Finally, we the rule (S- $\sqcup$ ) exploits the properties of a principal polyonad. Here we have a variable  $\rho$  such that all its in-edges are from pairs of ground constructors  $M_i$ , so we can simply apply the join function to compute a solution for  $\rho$ . For a principal polyonad, if such a solution exists, this simplification does not impact solvability of the rest of the constraint graph.

**Example.** Recall the information flow example we gave in Section 3.1, in Figure 2. Its principal type is the following, which is hardly readable:

$$\begin{aligned} \forall \bar{\rho}_i, a_1, a_2. P_0 \Rightarrow \text{intref } a_1 \rightarrow \text{intref } a_2 \rightarrow \rho_{27} () \\ \text{where } P_0 = (\text{Id}, \rho_3) \triangleright \rho_2, (\text{Id}, \text{IST } H \ a_2) \triangleright \rho_3, (\rho_{26}, \text{Id}) \triangleright \rho_4, (\text{Id}, \text{Id}) \triangleright \rho_4, \\ (\rho_8, \rho_4) \triangleright \rho_6, (\text{Id}, \rho_9) \triangleright \rho_8, (\text{Id}, \text{Id}) \triangleright \rho_9, (\rho_{11}, \rho_{25}) \triangleright \rho_{26}, \\ (\text{Id}, \rho_{12}) \triangleright \rho_{11}, (\text{Id}, \text{IST } H \ a_1) \triangleright \rho_{12}, (\rho_{17}, \rho_{23}) \triangleright \rho_{25}, (\rho_{14}, \rho_{18}) \triangleright \rho_{17}, \\ (\text{Id}, \text{Id}) \triangleright \rho_{18}, (\text{Id}, \rho_{15}) \triangleright \rho_{14}, (\text{Id}, \text{Id}) \triangleright \rho_{15}, (\rho_{20}, \rho_{24}) \triangleright \rho_{23}, \\ (\text{Id}, \text{IST } a_1 \ L) \triangleright \rho_{24}, (\text{Id}, \rho_{21}) \triangleright \rho_{20}, (\text{Id}, \text{Id}) \triangleright \rho_{21}. \end{aligned}$$

After applying (S- $\uparrow$ ) and (S- $\downarrow$ ) several times, and then hiding redundant constraints, we simplify  $P_0$  to  $P$  which contains only three constraints. If we had fixed  $a_1$  and  $a_2$  (the labels of the function parameters) to  $H$  and  $L$ , respectively, we could do even better. The three constraints would be  $(\text{IST } HL, \rho_6) \triangleright \rho_{27}, (\text{Id}, \text{Id}) \triangleright \rho_6, (\text{IST } HH, \text{IST } HL) \triangleright \rho_6$ . Then, applying (S- $\sqcup$ ) to  $\rho_6$  we would get  $\rho_6 \mapsto \text{IST } HH$ , which when applied to the other constraints leaves only  $(\text{IST } HL, \text{IST } HH) \triangleright \rho_{27}$ , which cannot be simplified further, since  $\rho_{27}$  appears in the result type.

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 11** (Simplification improves types). *For a principal polymonad, given  $\sigma$  and  $\sigma'$  where  $\sigma$  is  $\forall \bar{a} \bar{\rho}. P \Rightarrow \tau$  and  $\sigma'$  is an improvement of  $\sigma$ , having form  $\forall \bar{a}' \bar{\rho}'. \theta P \Rightarrow \tau$  where  $P \xrightarrow{\text{simplify}(\bar{\rho})} \theta$  and  $\bar{a}' \bar{\rho}' = (\bar{a} \bar{\rho}) - \text{dom}(\theta)$ . Then for all  $P'', \Gamma, x, e, m, \tau$ , if  $P'' \mid \Gamma, x: \sigma \vdash e : m \tau$  such that  $\models P''$  then there exists some  $P'''$  such that  $P''' \mid \Gamma, x: \sigma' \vdash e : m \tau$  and  $\models P'''$ .*

Note that our  $\xrightarrow{\text{simplify}(\bar{\rho})}$  relation is non-deterministic in the way it picks constraints to analyze, and also in the order in which rules are applied. In practice, for an acyclic constraint graph, one could consider nodes in the graph in topological order and, say, apply (S- $\sqcup$ ) first, since, if it succeeds, it eliminates a variable. For principal polymonads and acyclic constraint graphs, this process would always terminate.

However, if unification constraints induce cycles in the constraint graph, simply computing joins as solutions to internal variables may not work. This should not come as a surprise. In general, finding solutions to arbitrary polymonadic constraints is undecidable, since, in the limit, they can be used to encode the correctness of programs with general recursion. Nevertheless, simple heuristics such as unrolling cycles in the constraint graph a few times may provide good mileage, as would the use of domain-specific solvers for particular polymonads, and such approaches are justified by our coherence proof.

## 6 Related work and conclusions

This paper has presented *polymonads*, a generalization of monads and morphisms, which, by virtue of their relationship to Tate’s *productoids*, are extremely powerful, subsuming monads, parameterized monads, and several other interesting constructions. Thanks to supporting algorithms for (principal) type inference, (provably coherent) elaboration, and (generality-preserving) simplification (none of which Tate considers), this power comes with strong supports for the programmer. Like monads before them, we believe polymonads can become a useful and important element in the functional programmer’s toolkit.

Constructions resembling polymonads have already begun to creep into languages like Haskell. Notably, Kmett’s `Control.Monad.Parameterized` Haskell package [18] provides a type class for bind-like operators that have a signature resembling our  $(m_1, m_2) \triangleright m_3$ . One key limitation is that Kmett’s binds must be *functionally dependent*:  $m_3$  must be functionally determined from  $m_1$  and  $m_2$ . As such, it



is not possible to program morphisms between different constructors, i.e., the pair of binds  $(m_1, \text{Id}) \triangleright m_2$  and  $(m_1, \text{Id}) \triangleright m_3$  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. Kmett also requires units into  $\text{Id}$ , which may later be lifted, but such lifting only works for first-order code before running afoul of Haskell’s ambiguity restriction. Polymonads do not have either limitation. Kmett does not discuss laws that should govern the proper use of non-uniform binds. As such, our work provides the formal basis to design and reason about libraries that functional programmers have already begun developing.

While polymonads subsume a wide range of prior monad-like constructions, and indeed can express any system of *producer effects* [27], as might be expected, other researchers have explored generalizing monadic effects along other dimensions that are incomparable to polymonads. For example, Altenkirch et al. [2] consider *relative monads* that are not endofunctors; each polymonad constructor must be an endofunctor. Uustalu and Vene [28] suggest structuring computations comonadically, particularly to work with context-dependent computations. This suggests a loose connection with our encoding of contextual effects as a polymonad, and raises the possibility of a “co-polymonad”, something we leave for the future. Still other generalizations include reasoning about effects equationally using Lawvere theories [22] or with arrows [12]—while each of these generalize monadic constructions, they appear incomparable in expressiveness to polymonads. A common framework to unify all these treatments of effects remains an active area of research—polymonads are a useful addition to the discourse, covering at least one large area of the vast design space.

## References

- [1] M. Abadi, A. Banerjee, N. Heintze & J.G. Riecke (1999): *A core calculus of dependency*. In: *POPL*.
- [2] Thorsten Altenkirch, James Chapman & Tarmo Uustalu (2010): *Monads need not be endofunctors*. In: *FOSSACS*.
- [3] Robert Atkey (2009): *Parameterised notions of computation*. *Journal of Functional Programming* 19(3 & 4), pp. 335–376.
- [4] Greg Cooper & Shriram Krishnamurthi (2006): *Embedding dynamic dataflow in a call-by-value language*. In: *ESOP*.
- [5] K. Crary, A. Kliger & F. Pfenning (2005): *A monadic analysis of information flow security with mutable state*. *J. Funct. Program.* 15(02), pp. 249–291.
- [6] D. Devriese & F. Piessens (2011): *Information flow enforcement in monadic libraries*. In: *TLDI*.
- [7] Conal Elliott & Paul Hudak (1997): *Functional reactive animation*. In: *ICFP*.
- [8] A. Filinski (1999): *Representing layered monads*. In: *POPL*.
- [9] Andrzej Filinski (1994): *Representing monads*. In: *POPL*.
- [10] Jean-Christophe Filliâtre (1999): *A Theory of Monads Parameterized By Effects*.
- [11] J.A. Goguen & J. Meseguer (1982): *Security policy and security models*. In: *Symposium on Security and Privacy*.
- [12] John Hughes (2000): *Generalising monads to arrows*. *Sci. Comput. Program.* 37(1-3).
- [13] Graham Hutton & Erik Meijer (1998): *Monadic Parsing in Haskell*. *J. Funct. Program.* 8(4), pp. 437–444.
- [14] Mark P. Jones (1992): *A theory of qualified types*. In: *ESOP*.
- [15] Mark P. Jones (1993): *Coherence for Qualified Types*. Technical Report YALEU/DCS/RR-989, Yale University.
- [16] Mark P. Jones (1994): *Simplifying and Improving Qualified Types*. Technical Report YALEU/DCS/RR-1040, Yale University.
- [17] O. Kiselyov & C. Shan (2008): *Lightweight monadic regions*. In: *Haskell Symposium*.
- [18] Edward Kmett (2012): *Control.Monad.Parameterized package*. On Hackage repository.
- [19] P. Li & S. Zdancewic (2006): *Encoding information flow in Haskell*. In: *CSFW*.
- [20] Eugenio Moggi (1989): *Computational lambda-calculus and monads*. In: *LICS*.

- [21] Iulian Neamtiu, Michael Hicks, Jeffrey S. Foster & Polyvios Pratikakis (2008): *Contextual Effects for Version-Consistent Dynamic Software Updating and Safe Concurrent Programming*. In: *POPL*.
- [22] Gordon D. Plotkin & John Power (2001): *Semantics for Algebraic Operations*. *Electr. Notes Theor. Comput. Sci.* 45.
- [23] R. Pucella & J.A. Tov (2008): *Haskell session types with (almost) no class*. In: *Haskell Symposium*.
- [24] Norman Ramsey & Avi Pfeffer (2002): *Stochastic lambda calculus and monads of probability distributions*. In: *POPL*.
- [25] Alejandro Russo, Koen Claessen & John Hughes (2008): *A library for lightweight information-flow security in Haskell*. In: *Haskell Symposium*.
- [26] Nikhil Swamy, Nataliya Guts, Daan Leijen & Michael Hicks (2011): *Lightweight Monadic Programming in ML*. In: *ICFP*.
- [27] Ross Tate (2013): *The Sequential Semantics of Producer Effect Systems*. In: *POPL*.
- [28] Tarmo Uustalu & Varmo Vene (2008): *Comonadic Notions of Computation*. *Electr. Notes Theor. Comput. Sci.* 203(5).
- [29] Philip Wadler & Peter Thiemann (2003): *The marriage of effects and monads*. *ACM Trans. Comput. Logic* 4, pp. 1–32.

## Appendix

### A Polymonads are productoids and vice versa

Given a polymonad  $(\mathcal{M}, \Sigma)$ , we can construct a 4-tuple  $(\mathcal{M}, U, L, B)$  as follows:

**(Units)**  $U = \{(\lambda x. \text{bind } x (\lambda y. y)) : a \rightarrow M a \mid \text{bind} : (Id, Id) \triangleright M \in \Sigma\},$

**(Lifts)**  $L = \{(\lambda x. \text{bind } x (\lambda y. y)) : M a \rightarrow N a \mid \text{bind} : M \hookrightarrow N \in \Sigma\},$

**(Binds)** The set  $B = \Sigma - \{\text{bind} \mid \text{bind} : (Id, Id) \triangleright M \text{ or } \text{bind} : (M, Id) \triangleright N \in \Sigma\}.$

It is fairly easy to show that the above structure satisfies generalizations of the familiar laws for monads and monad morphisms.

**Theorem 12.** *Given a polymonad  $(\mathcal{M}, \Sigma)$ , the induced 4-tuple  $(\mathcal{M}, U, L, B)$  satisfies the following properties.*

**(Left unit)**  $\forall \text{unit} \in U, \text{bind} \in B.$  if  $\text{unit} : \forall a. a \rightarrow M a$  and  $\text{bind} : (M, N) \triangleright N$  then  $\text{bind} (\text{unit } e) f = f(e)$  where  $e : \tau$  and  $f : \tau \rightarrow N \tau'.$

**(Right unit)**  $\forall \text{unit} \in U, \text{bind} \in B.$  if  $\text{unit} : \forall a. a \rightarrow N a$  and  $\text{bind} : (M, N) \triangleright M$  then  $\text{bind } m (\text{unit}) = m$  where  $m : M \tau.$

**(Associativity)**  $\forall \text{bind}_1, \text{bind}_2, \text{bind}_3, \text{bind}_4 \in B.$  if  $\text{bind}_1 : (M, N) \triangleright P,$   
 $\text{bind}_2 : (P, R) \triangleright T, \text{bind}_3 : (M, S) \triangleright T,$  and  $\text{bind}_4 : (N, R) \triangleright S$  then  
 $\text{bind}_2 (\text{bind}_1 m f) g = \text{bind}_3 m (\lambda x. \text{bind}_4 (f x) g)$   
 where  $m : M \tau, f : \tau \rightarrow N \tau'$  and  $g : \tau' \rightarrow R \tau''$

**(Morphism 1)**  $\forall \text{unit}_1, \text{unit}_2 \in U, \text{lift} \in L.$  if  $\text{unit}_1 : \forall a. a \rightarrow M a,$   $\text{unit}_2 : \forall a. a \rightarrow N a$  and  $\text{lift} : \forall a. M a \rightarrow N a$  then  $\text{lift} (\text{unit}_1 e) = \text{unit}_2 e$  where  $e : \tau.$

**(Morphism 2)**  $\forall \text{bind}_1, \text{bind}_2 \in B, \text{lift}_1, \text{lift}_2, \text{lift}_3 \in L.$  if  $\text{bind}_1 : (M, P) \triangleright S,$   
 $\text{bind}_2 : (N, Q) \triangleright T, \text{lift}_1 : \forall a. M a \rightarrow N a, \text{lift}_2 : \forall a. P a \rightarrow Q a$  and  $\text{lift}_3 : \forall a. S a \rightarrow T a$  then  $\text{lift}_3 (\text{bind}_1 m f) =$   
 $\text{bind}_2 (\text{lift}_1 m) (\lambda x. \text{lift}_2 (f x))$   
 where  $m : M \tau$  and  $f : \tau \rightarrow P \tau'.$

Now we show how this definition can be used to relate polymonads to Tate's *productoids* [27]. The definition of a productoid is driven by an underlying algebraic structure: the effectoid [27, Theorem 1]. An effectoid  $(E, U, \leq, \mapsto)$  is a set  $E$ , with an identified subset  $U \subseteq E$  and relations  $\leq \subseteq E \times E$  and  $(.; -) \mapsto - \subseteq E \times E \times E$ , that satisfies six monoid-like conditions. It is possible to show that a polymonad directly induces an effectoid structure and hence a productoid.

**Lemma 13.** *Given a polymonad  $(\mathcal{M}, U, L, B)$  we can define an effectoid  $(E, U, \leq, (.; -) \mapsto -)$  as follows.*

$$\begin{aligned} E &= \mathcal{M} & U &= \{M \mid \text{unit} : a \rightarrow M a \in U\} \\ \leq &= \{(M, N) \mid \text{lift} : M a \rightarrow N a \in L\} & (.; -) \mapsto - &= \{(M, N, P) \mid (M, N) \triangleright P \in B\} \end{aligned}$$

**Lemma 14.** *Every polymonad gives rise to a productoid.*

*Proof.* We have shown that a polymonad gives rise to an effectoid. Given an effectoid  $(E, U, \leq, (.; -) \mapsto -)$  a productoid is defined as a collection of functors indexed by the collection  $E$ , and three collections of natural transformations indexed by the three relations. These functors and natural transformations are required to satisfy five addition properties [27, Theorem 2]. The five properties are the five properties of Theorem 12, so the proof is immediate.  $\square$

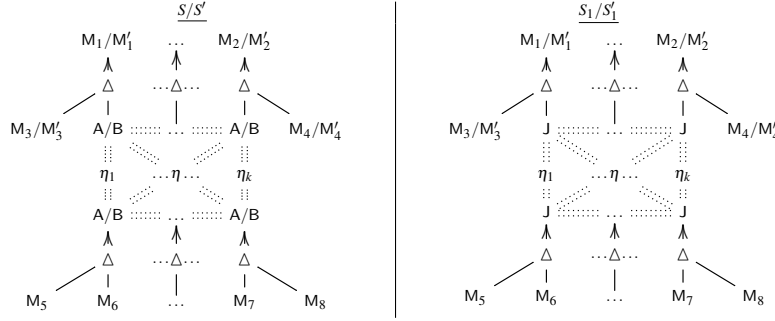


Figure 8: Constraint graphs used to illustrate the proof of coherence (Theorem 10)

Interestingly, we can identify conditions where the opposite direction also holds.

**Lemma 15.** *A productoid  $(\mathbf{C}, \{F_e : \mathbf{C} \rightarrow \mathbf{C}\}_{e \in E}, \{\eta : 1 \Rightarrow F_e\}_{e \in U}, \{\mu : F_{e_1} \circ F_{e_2} \Rightarrow F_{e_3}\}_{(e_1; e_2) \mapsto e_3}, \{\sigma : F_{e_1} \Rightarrow F_{e_2}\}_{e_1 \leq e_2})$  that in addition satisfies the following conditions gives rise to a polymonad.*

1.  $\text{Id} \in E$  and  $F_{\text{Id}} = 1$
2.  $\text{Id} \in U$
3. For all  $e \in E$ ,  $(e; \text{Id}) \mapsto e$
4. For all  $e \in E$ ,  $\mu : F_e \circ 1 \Rightarrow F_e = \text{Id}$
5.  $(e_1; e_2) \mapsto e \wedge e'_1 \leq e_1 \Rightarrow (e'_1; e_2) \mapsto e$
6.  $(e_1; e_2) \mapsto e \wedge e'_2 \leq e_2 \Rightarrow (e_1; e'_2) \mapsto e$

These additional conditions are fairly mild: (1)-(4) simply ensure that the  $\text{Id}$  element is interpreted as the identity functor. Conditions (5)-(6) are also quite straightforward; certainly if the category is cartesian closed then the extra natural transformations are always defined.

## B Coherence of solutions

**Lemma 16** (Solutions to a core). *For a polymonad  $(\mathcal{M}, \Sigma)$ , and a constraint graph  $G$  with a core  $G'$ , the set of all solutions  $\mathcal{S}'$  to  $G'$  includes all the solutions  $\mathcal{S}$  of  $G$ .*

*Proof.* (Sketch) This is easy to see, since  $G'$  differs from  $G$  only in that it includes fewer unification constraints. So, all solutions to  $G$  are also solutions to  $G'$ .  $\square$

**Theorem 17** (Coherence). *For all principal polymonads, derivations  $P|\Gamma \vdash e : t \rightsquigarrow e$  such that  $\text{unambiguous}(P, \Gamma, t)$ , and for any two solutions  $S$  and  $S'$  to  $G_P$  that agree on  $R = \text{ftv}(\Gamma, t)$ , we have  $S \cong_R S'$ .*

*Proof.* We consider the set  $\mathcal{S}$  of all solutions to the core of  $G_P$  that agree on  $\text{ftv}(\Gamma, t)$ , and prove that all these solutions are in the same equivalence class. By Lemma 16,  $\{S, S'\} \subseteq \mathcal{S}$ , establishing our goal.

Let  $G = (V, A, E_{\triangleright}, E_{eq})$  be a core of  $G_P$  and let  $S$  and  $S'$  be arbitrary elements of  $\mathcal{S}$ .  $S$  and  $S'$  may only differ on the open variables of  $P$ . Since  $G$  is unambiguous, the nodes associated with these variables all have non-zero in- and out-degree. Let  $U_{S, S'} = \{v \mid v \in V \wedge S(v) \neq S'(v)\}$ ; the proof proceeds by induction on the size of  $U$ .

**Base case**  $|U_{S, S'}| = 0$ : Trivial, since we have  $S(v) = S'(v)$ , for all  $v$ .

**Induction step**  $|U_{S, S'}| = i$ : From the induction hypothesis: All solutions  $S_1$  and  $S'_1$  such that  $|U_{S_1, S'_1}| < i$ , we have  $S_1 \cong S'_1$ .

Topologically sort  $G$ , such that all vertices in the same connected component following edges in  $E_{eq}$  have the same index, and each vertex  $v$  is assigned an index greater than the index of all vertices  $v'$  such that  $(v, v')$  is an edge in  $E_{\triangleright}$ . That is, “leaf” nodes have the highest indices.

Pick a vertex  $v$  with the maximal index, such that  $S(v) = A$  and  $S'(v) = B$ , for  $A \neq B$ , and let  $I$  be the set of vertices reachable from  $v$  via unification edges. Since both  $S$  and  $S'$  are solutions, there must exist

an open variable  $\rho$  such that  $A(v) = \rho$ , and since  $G$  is a core, there must be some non-empty set of flow edges incident on  $v$ .

Thus, the neighborhood of  $v$  in the graphs  $G$ , under assignment  $S$  and  $S'$  has a shape as shown in graph at left in Figure 8. All the nodes in  $I$  are shown connected by double dotted lines—they each have assignment  $A/B$  in  $S/S'$ . Since all the nodes in  $I$  have an index greater than the index of any variable that differs among  $S$  and  $S'$ , all their immediate predecessors have identical assignments in the two solutions (i.e.,  $M_5, \dots, M_8$ ). However, the other assignments may differ, (e.g., the top-left node could be assigned  $M_1$  in  $S$  and  $M'_1$  in  $S'$ , etc.) Each flow-edge  $\{\eta_1, \dots, \eta_k\}$  incident upon one of the nodes with the same index as  $v$  is also labeled.

Now, since we have a principal polymonad, there exists a principal join of  $\{(M_5, M_6), \dots, (M_7, M_8)\}$ —call it  $J$ . Consider the assignment  $S_1$  (resp.  $S'_1$ ) that differs from  $S$  (resp.  $S'$ ) only by assigning  $J$  to each vertex in  $I$  instead of  $A$  (resp.  $B$ ).

We first show that  $S_1$  (resp.  $S'_1$ ) is a solution and that  $S \cong S_1$  (resp.  $S' \cong S'_1$ ). Then, we note that since  $S_1$  and  $S'_1$  agree on all the vertices in  $I$ ,  $|U_{S_1, S'_1}| < i$ , so we apply the induction hypothesis to show that  $S_1 \cong S'_1$  and conclude with transitivity of  $\cong$ .

To show that  $S_1$  (resp.  $S'_1$ ) is a solution, since  $J$  is a join of  $M_5, M_6, \dots$ , then  $\{(M_5, M_6) \triangleright J, \dots, (M_7, M_8) \triangleright J\}$  all exist, as well as  $J \hookrightarrow A$  (resp.  $J \hookrightarrow B$ ). By the Closure property, for every  $(M, A) \triangleright M'$  (resp.  $B$ ) there also exists  $(M, J) \triangleright M'$ . Thus, the assignment of  $J$  to  $I$  is valid for a solution.

To show that  $S \cong S_1$  (resp.  $S' \cong S'_1$ ), we have to show that  $F_S(\eta_i) = F_{S_1}(\eta_i)$  (resp.  $F_{S'}(\eta_i) = F_{S'_1}(\eta_i)$ ), for all  $i$ . Taking  $\eta_k$  as a representative case (the other cases are similar), we need to show the identity below, which is an immediate corollary of Associativity 1 and 2 (resp. for  $B, M'_4, M'_2$ ).

$$\text{bind}_{A, M_4, M_2}(\text{bind}_{M_7, M_8, A} x y) z = \text{bind}_{J, M_4, M_2}(\text{bind}_{M_7, M_8, J} x y) z$$

□