

# Kinds Are Calling Conventions

PAUL DOWNEN and ZENA M. ARIOLA, University of Oregon, USA

SIMON PEYTON JONES, Microsoft Research, UK

RICHARD A. EISENBERG, Bryn Mawr College, USA and Tweag I/O, UK

A language supporting polymorphism is a boon to programmers: they can express complex ideas once and reuse functions in a variety of situations. However, polymorphism is pain for compilers tasked with producing efficient code that manipulates concrete values.

This paper presents a new intermediate language that allows for efficient static compilation, while still supporting flexible polymorphism. Specifically, it permits polymorphism over not only the types of values, but also the representation of values, the arity of primitive machine functions, and the evaluation order of arguments—all three of which are useful in practice. The key insight is to encode information about a value's calling convention in the *kind* of its type, rather than in the type itself.

CCS Concepts: • **Software and its engineering** → *Semantics; Compilers*.

Additional Key Words and Phrases: arity, levity, representation, polymorphism, type systems

## ACM Reference Format:

Paul Downen, Zena M. Ariola, Simon Peyton Jones, and Richard A. Eisenberg. 2020. Kinds Are Calling Conventions. *Proc. ACM Program. Lang.* 4, ICFP, Article 104 (August 2020), 29 pages. <https://doi.org/10.1145/3408986>

## 1 INTRODUCTION

Polymorphism supports reuse by allowing one piece of code to work with values of many different types. But ubiquitous polymorphism usually comes with a runtime cost: all values must share a common representation, usually a pointer to a “boxed” (heap-allocated) object. This is sometimes *much* less efficient than a monomorphic version of the same code, specialized to a particular representation (such as an unboxed 64-bit word).

One approach is to specialize code to a single type. But we would get more reuse if we could specialize to, say, “any type represented by an unboxed 64-bit word.” Since *kinds* classify types, perhaps we can write code that is monomorphic in the kind, but polymorphic in the type. Hence our slogan: kinds are calling conventions. For example, consider the function *twice*:

$$\textit{twice } f \ x = f \ (f \ x)$$

To control performance, we would like to have a say in matters like: Can  $f$  be a thunk? How many arguments does  $f$  expect (its arity)? Can  $x$  be a thunk? How is  $x$  represented? Moreover, we want to express the answers to these questions in a type system with strong static guarantees.

A major insight of this paper is the discovery that we can refine the vague notion of “ways in which we want to classify types” along three different axes:

---

Authors' addresses: Paul Downen; Zena M. Ariola, University of Oregon, Eugene, Oregon, USA, [pdownen@cs.uoregon.edu](mailto:pdownen@cs.uoregon.edu), [ariola@cs.uoregon.edu](mailto:ariola@cs.uoregon.edu); Simon Peyton Jones, Microsoft Research, Cambridge, UK, [simonpj@microsoft.com](mailto:simonpj@microsoft.com); Richard A. Eisenberg, Bryn Mawr College, Bryn Mawr, PA, USA, Tweag I/O, Cambridge, UK, [rae@richarde.dev](mailto:rae@richarde.dev).

---

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

© 2020 Copyright held by the owner/author(s).

2475-1421/2020/8-ART104

<https://doi.org/10.1145/3408986>

- *Representation*. How is this argument represented at runtime?
- *Levity*. What is the evaluation strategy of this argument (e.g., call-by-value or call-by-need)?
- *Arity*. For functions, how many arguments are needed before its code can be executed?

Many functions can be polymorphic in some of these axes, but not in others.

Our focus is on an *intermediate language*. The programmer may write in a uniform language, but the compiler needs an intermediate language that can express low-level representation choices, and expose those choices to the optimizer. For example, the programmer might work exclusively with boxed integer values, say of type `Int`, but the intermediate language can have an unboxed type `Int#`, together with explicit operations to box and unbox integers. This allows the optimizer to eliminate many box-followed-by-unbox chains [Peyton Jones and Launchbury 1991].

This paper builds directly on several earlier works that statically keep track of different representations [Eisenberg and Peyton Jones 2017] and function arities [Downen et al. 2019] within a type and kind system. Our new contribution is to bring them together into a single framework, more powerful and more precise than any of its predecessors. Specifically:

- We introduce a polymorphic intermediate language that statically captures *calling conventions in kinds*, and has polymorphism over the *representation*, *levity*, and *arity* of types (Section 4).
- Our intermediate language is equally well-suited for *both eager and lazy* functional languages. Concretely we show how to compile two higher-level, polymorphic source languages—call-by-name and call-by-value System F—to our intermediate language (Section 5).
- We show how to compile our polymorphic intermediate language to a lower-level language with multiple representations (e.g., pointers versus integers) and multi-arity functions, but *not* polymorphism (Section 6). Compilation is driven by kinds and keeps type erasure and code reuse; typing restrictions ensures polymorphic code is compiled to monomorphic code.
- We provide evidence of correctness for the full compilation process (Theorems 1 to 4) from both call-by-name and call-by-value source languages to machine code.
- We describe a small extension to our intermediate language to allow for dynamic checks on the arities of closures, so that we can use the best arity available at runtime (Section 7).

## 2 FUNCTION ARITY

We identify three axes of classification, above. Of these three, we peel off arity to explain it first; mixed representations and evaluation strategies have existed for longer and may be more familiar. This section gives a high-level overview of our approach; the details will be nailed down in Section 4.

### 2.1 What Is Arity?

What is a function’s “arity?” In practical terms, the arity of a function determines what code a compiler will generate to call that function. All modern architectures support calling conventions that allow for efficiently calling subroutines by passing multiple arguments at once. If a function has arity  $n$ , then a call to that function will pass  $n$  arguments simultaneously.

In a naively compiled curried language, every function has arity 1. This implementation of functional languages is unacceptably inefficient. Instead, compilers must somehow map curried surface-language functions onto multi-argument machine-language functions.

To do this, we need two things. First, an *arity analysis*; and second, an *intermediate language* in which arity is explicit, so that we can express and memorialise the results of the analysis.

Arity analysis is not always straightforward. For example:

$$\begin{array}{ll}
 f_1, f_2, f_3 : \text{Int} \rightarrow \text{Int} \rightarrow \text{Int} & \\
 f_1 = \lambda x. \lambda y. \text{let } z = \text{expensive } x \text{ in } y + z & f_2 = \lambda x. f_1 \ x \\
 f_3 = \lambda x. \text{let } z = \text{expensive } x \text{ in } \lambda y. y + z & f_4 = \lambda x. f_3 \ x
 \end{array}$$

Here,  $f_1$  clearly has arity 2 because it starts with two  $\lambda$ s: it takes two arguments before computing a result. Function  $f_2$  has only one  $\lambda$ , but still has arity 2, because we could safely  $\eta$ -expand it to have two  $\lambda$ s. In other words,  $f_2$  has arity 2, just like  $f_1$ , because it also requires two arguments before doing serious work. The function  $f_3$  has the same type as  $f_1$ , but it *must* take its arguments one at a time: in the call  $(\text{map } (f_3 \ 100) \ xs)$  we expect  $(\text{expensive } 100)$  to be computed at most once, whereas the same call for  $f_1$  would recompute  $(\text{expensive } 100)$  for each element of  $xs$ . Similarly,  $f_4$  has arity 1: it cannot be  $\eta$ -expanded without the risk of computing  $(\text{expensive } 100)$  repeatedly.

These choices become particularly clear in a call-by-value language with side effects. For example, if  $(\text{expensive } 100)$  printed something on the screen, the fact that it is evaluated only once—rather than once for each element of  $xs$ —is a matter of semantics, not mere efficiency. Even in a pure language like Haskell, an optimizing compiler should still treat computation as a sort of effect; it must, for example, avoid changing the asymptotic efficiency of the program.

In light of these examples, here are two informal definitions of arity:

- An expression  $e$  has arity  $n$  when it can be soundly  $\eta$ -expanded to  $(\lambda x_1 \dots x_n. e \ x_1 \dots x_n)$ . “Soundness” concerns semantics in an effectful language, but “only” efficiency in a pure one.
- An expression  $e$  has arity  $n$  when it does no “useful work” until it is applied to  $n$  arguments; hence those arguments can be passed simultaneously.

Since the type of a function does not describe its arity (compare  $f_1$  and  $f_3$  above), practical compilers like GHC perform arity analysis based on intensional properties: the form of the expression determines its arity. In its simplest form, we can just count  $\lambda$ s. Since that is pessimistic on examples like  $f_2$ , GHC uses a variety of simple static heuristic analyses [e.g., [Breitner 2014](#)]. The focus of this paper is not on arity analysis, however, but rather on the intermediate language in which we can memorialise the results of that analysis.

## 2.2 Arity in the Intermediate Language

However arity analysis is done, we need a way to express its results in the intermediate language. In GHC, this is done through an informal decoration on each binder describing its arity.<sup>1</sup> This turns out to be extremely unsatisfactory in practice: GHC has lots and lots of dark corners as a result of this rather squishy notion of arity.

It would be much better if arity were a solid, statically-checked part of the intermediate language. How is that possible? In the world of  $\lambda$ -calculi, we are familiar with calculi having unrestricted  $\beta$  and  $\eta$  rules (such as the call-by-name  $\lambda$ -calculus), as well as calculi having restricted  $\beta$  and  $\eta$  (such as the call-by-value [[Sabry and Wadler 1997](#)] or call-by-need  $\lambda$ -calculus [[Ariola and Felleisen 1997](#)]). The latter are often used as intermediate languages in compilers, to avoid recomputation of, say, integers. The exact restrictions depend on the language being compiled: its evaluation strategy, whether there are side effects, and so on. Our key idea, introduced by [Downen et al. \[2019\]](#), is this:

Define an intermediate language ( $\mathcal{IL}$ ) that has *unrestricted  $\eta$  expansion* for functions, while allowing for restricted  $\beta$  reduction on other types of expressions.

This is an unusual choice, so we use a different notation for functions,  $\lambda(x:\tau).e : \tau \rightsquigarrow \sigma$ , where the ( $\rightsquigarrow$ ) arrow denotes the new primitive function type. Now arity can be read from types: you can freely  $\eta$ -expand *any* expression  $e : \tau \rightsquigarrow \sigma$ , without reference to the form of  $e$ . The previous paper [[Downen et al. 2019](#)] and this one are simply working out the consequences of this one idea.

## 2.3 Currying

If the primitive function type of our intermediate language allows unrestricted  $\eta$ -expansion, how can we express functions like  $f_3$  that intentionally use currying as a way to avoid work duplication?

<sup>1</sup>For aficionados, it is part of the binder’s `IdInfo`

We can do so like this:

$$\begin{aligned} f_3 &: \text{Int} \rightsquigarrow \{\text{Int} \rightsquigarrow \text{Int}\} \\ f_3 &= \lambda x. \text{let } z = \text{expensive } x \text{ in Clos } (\lambda y. y + z) \end{aligned}$$

The type  $\text{Int} \rightsquigarrow \{\text{Int} \rightsquigarrow \text{Int}\}$  makes explicit that  $f_3$  is an arity-1 function that returns a closure, denoted by the curly braces, inside which is another arity-1 function. In the term language,  $\lambda$  builds a *primitive function* of type  $\tau \rightsquigarrow \sigma$ , while  $\text{Clos}$  builds a *closure* of type  $\{\tau\}$ . Here are the (slightly simplified) introduction (-I) and elimination (-E) rules:

$$\begin{array}{c} \frac{\Gamma, x:\tau \vdash e:\sigma}{\Gamma \vdash \lambda x:\tau. e:\tau \rightsquigarrow \sigma} \text{FUN-I} \qquad \frac{\Gamma \vdash e:\tau}{\Gamma \vdash \text{Clos } e:\{\tau\}} \text{CLO-I} \\ \frac{\Gamma \vdash e:\tau \rightsquigarrow \sigma \quad \Gamma \vdash e':\tau}{\Gamma \vdash e e':\sigma} \text{FUN-E} \qquad \frac{\Gamma \vdash e:\{\tau\}}{\Gamma \vdash \text{App } e:\tau} \text{CLO-E} \end{array}$$

The App form unboxes the primitive function wrapped up by Clos.

Our plan, then, is to desugar the source-language function type  $\tau \rightarrow \sigma$  into the intermediate-language type as  $\{\tau \rightsquigarrow \sigma\}$ , adding the corresponding Clos and App constructs in the terms.<sup>2</sup> Then we can perform arity analysis, and express its results by transforming the intermediate language program into one with fewer intermediate Clos nodes.

## 2.4 Functions Are Called, Not Evaluated

If we are able to freely  $\eta$ -expand, we must only evaluate functions when they are called. Consider:

$$x = \text{let } f : \text{Int} \rightsquigarrow \text{Int} = \text{expensive } 100 \text{ in } \dots f \dots f \dots$$

When might *expensive* 100 be evaluated? In a strict language, it is evaluated right away, so the value (some first-class function) can be bound to  $f$  before continuing with the body of the *let*. In a lazy language like Haskell, *expensive* 100 might be forced long before  $f$  needs to be called, as in *seq*  $f$   $y$  or with a strict pattern. But in both cases, evaluating  $f$  without calling it violates the unrestricted  $\eta$  we desire. After all, the definition of  $x$  is  $\eta$ -equivalent to:

$$x' = \text{let } f : \text{Int} \rightsquigarrow \text{Int} = \lambda y. \text{expensive } 100 \ y \text{ in } \dots f \dots f \dots$$

After  $\eta$ -expansion, the only way to evaluate *expensive* 100 is to *call*  $f$  with an argument. Unrestricted  $\eta$  means that  $x$  and  $x'$  *must* be the same—in both semantics and asymptotic cost.

Unrestricted  $\eta$ -expansion requires a matching evaluation order for function bindings; the evaluator must treat every expression  $e : \tau \rightsquigarrow \sigma$  the same as  $(\lambda x. e \ x)$ . Yet, we do not want to change the evaluation order for expressions of other types, like  $\text{Int}$ . Thus, our language's semantics becomes type-directed: it is only the type of  $f$  (in this case,  $\text{Int} \rightsquigarrow \text{Int}$ ) that tells us not to evaluate the right-hand side. In a precise sense (spelled out formally in Section 4), we cannot *evaluate* functions (e.g., due to strictness, *seq*, etc.); functions can only be *called*. Furthermore, we will see in Section 3.7 that functions are, in fact,  $\eta$ -expanded during compilation to a lower level. So a semantics that prematurely evaluates functions before they are called is not supported by our machine language.

## 3 WHY KINDS ARE CALLING CONVENTIONS

So far we have reviewed the ideas of Downen et al. [2019]. We are now ready to introduce the problem we tackle in this paper, and our solution to it.

Previous research asserted *types are calling conventions* [Bolingbroke and Peyton Jones 2009]. We respectfully disagree, instead claiming that *kinds* are calling conventions. As we shall see, this principle offers a unified framework combining several previous works on representations

<sup>2</sup>In Downen et al. [2019], the intermediate language itself had two arrows,  $(\rightarrow)$  as well as  $(\rightsquigarrow)$ , but we have found it clearer to have only one arrow for primitive functions plus the closure type  $\{\tau\}$ .

[Eisenberg and Peyton Jones 2017; Peyton Jones and Launchbury 1991], arity [Bolingbroke and Peyton Jones 2009; Downen et al. 2019; Marlow and Peyton Jones 2004], and mixed evaluation strategies [Downen and Ariola 2018] in intermediate languages. Specifically, we provide a type system in which the kind  $\kappa$  of a type  $\tau : \kappa$  classifies all the intensional properties needed to compile an expression of type  $\tau$ , namely its *representation*, *arity*, and *levity*.

### 3.1 Why Polymorphism Is a Problem

Consider this example of a polymorphic definition, in our intermediate language

$$\begin{aligned} poly &: \forall a. (\text{Int} \rightsquigarrow \text{Int} \rightsquigarrow a) \rightsquigarrow (a, a) \\ poly &= \Lambda a. \lambda(f : \text{Int} \rightsquigarrow \text{Int} \rightsquigarrow a). \text{let } g : \text{Int} \rightsquigarrow a = f \text{ 3 in } (g \text{ 5}, g \text{ 4}) \end{aligned}$$

What is the arity of  $f$  and  $g$ ? With “arities in the types,” we count the arrows, and answer 2 and 1 respectively. *But what if  $a$  were instantiated by  $\text{Bool} \rightsquigarrow \text{Bool}$ ?* Then suddenly the answers become 3 and 2 respectively. Yikes!

One solution is to monomorphize the entire program, creating type-specialized versions of each polymorphic function. This is patently unsatisfactory. First, it is a whole-program transformation. Second, some languages (including Haskell and OCaml) support polymorphic recursion, which makes static monomorphization impossible. Third, there is a risk of creating many copies of essentially the same function, many of which wastefully compile to identical machine code.

Instead, we retain the traditional type-erasure model: each polymorphic function is compiled to a single chunk of machine code that works the same no matter how its type is specialized. Under this model, our function  $poly$  above is problematic. To compile code we must know the arity of every function we call (because its arity determines its calling convention, and thus, what code to generate), but in  $poly$  we do not know a stable arity of  $f$  or  $g$  for every instance of  $a$ .

### 3.2 Nonuniform Representations and Polymorphism

Interestingly, this exact same problem has arisen before, in the context of *unboxed data types*. A contribution of this paper is to show that both can be solved with the same idea.

Nearly thirty years ago, GHC introduced the idea of distinguishing boxed and unboxed data types in its intermediate language [Peyton Jones and Launchbury 1991]. They introduced two distinct types for integers:  $\text{Int\#}$  for primitive, unboxed machine integers; and  $\text{Int}$  for boxed integers, represented as a pointer to a heap-allocated object and defined as:

$$\text{data Int} = \text{I\# Int\#}$$

The (sole) goal of distinguishing these two types is efficiency. If we had only boxed integers, then even simple addition would be forced to evaluate and unbox each argument, then box up the result. By making boxing and unboxing explicit we expose much more low-level information to the optimizer, and can eliminate lots of intermediate boxes. For example, assume we have  $\text{plus\#}$  and  $\text{minus\#}$  primitive operations, both of type  $\text{Int\#} \rightarrow \text{Int\#} \rightarrow \text{Int\#}$ , and consider:

$$\begin{array}{ll} \text{plus, minus} : \text{Int} \rightarrow \text{Int} \rightarrow \text{Int} & \text{sumFrom} : \text{Int} \rightarrow \text{Int} \\ \text{plus (I\# x) (I\# y)} = \text{I\# (plus\# x y)} & \text{sumFrom (I\# 0)} = \text{I\# 0} \\ \text{minus (I\# x) (I\# y)} = \text{I\# (minus\# x y)} & \text{sumFrom x} = \text{plus x (sumFrom (minus x (I\# 1)))} \end{array}$$

This definition is wasteful; each recursive step allocates several new boxes on the heap only to be immediately used by  $\text{plus}$ ,  $\text{minus}$ , and  $\text{sumFrom}$ . Instead, as Peyton Jones and Launchbury [1991] show, the recursive function can be optimized using the *worker/wrapper* transformation, like so:

$$\begin{array}{ll} \text{sumFrom} : \text{Int} \rightarrow \text{Int} & \text{sumFrom\#} : \text{Int\#} \rightarrow \text{Int\#} \\ \text{sumFrom (I\# x)} = \text{I\# (sumFrom\# x)} & \text{sumFrom\# 0} = 0 \\ & \text{sumFrom\# x} = \text{plus\# x (sumFrom\# (minus\# x 1))} \end{array}$$

Now, the recursion is done by the more efficient *sumFrom#* function which works directly on machine integers; no boxes are allocated or consumed, and so *sumFrom#* can be compiled with no intermediate allocations. *sumFrom* becomes a wrapper around *sumFrom#*, just handling the issues of boxing and unboxing. This is a huge gain in both time and space.

Notice the similarities between this work on unboxed data types and the earlier discussion about arity. In both cases, we make *IL* expose more primitive, but less convenient, types of values (primitive functions and unboxed integers respectively), along with a way to explicitly “box” them into a more convenient form (using *Clos* and *I#* respectively) and later “unbox” them (using *App* and pattern matching on *I#* respectively). Making these boxing and unboxing operations explicit in the syntax of *IL* programs makes them accessible to the optimizer.

Alas, unboxed types cause trouble with polymorphism. For example, consider the generic, higher-order binary application function:

$$\text{binapp} : \forall a b c. (a \rightarrow b \rightarrow c) \rightarrow a \rightarrow b \rightarrow c \qquad \text{binapp } f \ x \ y = f \ x \ y$$

To compile *binapp* to a single piece of code, we cannot expect the calls *binapp plus* and *binapp plus#* to both work; *binapp* would somehow have to handle arguments *x* and *y* with different representations (perhaps stored in different registers) to pass them along to *plus* versus *plus#*.

Notice that this is the *exact same problem* that we had with *poly* in Section 3.1: the code we compile for *binapp* depends on how the type variables *a* and *b* are instantiated.

### 3.3 A Stop-Gap Solution

Because of the difficulty with polymorphism, unboxed types were originally introduced with a draconian restriction on polymorphism: *polymorphic type variables (like  $a, b, c$  in *binapp*) cannot be instantiated with unboxed types (like *Int#*)* [Peyton Jones and Launchbury 1991]. The mechanism for enforcing the restriction was the *kind system*: *Int* has kind  $\star$  (the kind of ordinary data types), but *Int#* has a different kind  $\#$  (the kind of unboxed data). GHC then required that the kind of a quantified type variable *t* could be  $\star$  or, say,  $(\star \rightarrow \star)$ , but never  $\#$ .

Since this is the same problem as the one we encountered for arity, in function *poly* in Section 3.1, we might expect the same solution to work. And indeed that is the approach adopted by Downen et al. [2019]: type variables cannot range over arrow types  $\tau \rightsquigarrow \sigma$ , whose kind is different from  $\star$ .

While this approach has served GHC well for nearly three decades, it has two major inadequacies:

(1) *Too restrictive*. Consider the *error* function, which prints a message and prematurely ends the program. It can have the type  $\forall t. \text{String} \rightarrow t$ . Because *error* never returns a value, it doesn’t matter how *t* is represented; truly the same code can be used no matter if *t* is boxed or unboxed. But the draconian restriction on unboxed types rejects this kind of polymorphism. Instead, specialized versions are needed, like *errorInt#* : *String*  $\rightarrow$  *Int#*, even though the compiled code is identical. For *error* we would prefer to be able to say that the representation of *t* doesn’t matter.

(2) *Kind polymorphism*. In a language with kind polymorphism, we can write a term with type  $\forall k. \forall (a : k). \tau$ . Here, *a* can have *any* kind *k*, including  $\#$ . We have now lost our ability to prevent quantifiers over  $\#$ -types, and need a new solution to the problem of uncompletable polymorphism.<sup>3</sup>

<sup>3</sup>A particularly knowledgeable reader might be aware that GHC supported kind polymorphism and restricted quantification for a few years. This worked because a type like  $\forall k(a:k). a \rightarrow a$  is ill-kinded; the  $\rightarrow$  type former puts requirements on *a*. However, with *kind equalities* [Weirich et al. 2013], a type variable  $a : k$  can be cast to a more specific kind, causing chaos. In fact, this interaction originally spurred the development of Eisenberg and Peyton Jones [2017].



### 3.4 A Better Way: Look to the Kinds

Fortunately, both (1) and (2) can be solved, by using...more polymorphism. Let's start with representation polymorphism, using the approach of [Eisenberg and Peyton Jones \[2017\]](#).<sup>4</sup> We give `Int` and `Int#` these more refined kinds:

$$\text{Int} : \text{TYPE } \text{PtrR} \quad \text{Int\#} : \text{TYPE } \text{IntR}$$

where  $\text{TYPE} : \text{Rep} \rightarrow \star$ , and  $\text{PtrR}, \text{IntR} : \text{Rep}$ .<sup>5</sup> The idea is that, given  $\tau : \text{TYPE } r$ , values of type  $\tau$  have a runtime representation described by  $r$ . Now we can quantify over types represented by a heap pointer with  $\forall(a : \text{TYPE } \text{PtrR}).\tau$ . But we can *also* define representation-polymorphic functions—like `error` :  $\forall(r : \text{Rep})(a : \text{TYPE } r)$ . `String`  $\rightarrow a$ —solving problem (1).

What of problem (2)? Instead of limiting how type variables can be instantiated (which is incompatible with full kind polymorphism), we instead add side conditions to the typing rules for abstraction and application, excluding uncompileable programs like so:

$$\frac{\Gamma, x:\tau \vdash e : \sigma \quad \Gamma \vdash \tau \text{ mono-rep}}{\Gamma \vdash \lambda x:\tau. e : \tau \leadsto \sigma} \text{FUN-I} \quad \frac{\Gamma \vdash e : \tau \leadsto \sigma \quad \Gamma \vdash e' : \tau \quad \Gamma \vdash \tau \text{ mono-rep}}{\Gamma \vdash e e' : \sigma} \text{FUN-E}$$

The  $\tau$  **mono-rep** caveat says  $\tau$ 's kind must be representation-monomorphic; that is, it can be `TYPE PtrR`; or `TYPE IntR`; but *not* `TYPE r`. For example, we might try to give `binapp` this type

$$\text{binapp} : \forall(r_a, r_b, r_c : \text{Rep})(a : \text{TYPE } r_a)(b : \text{TYPE } r_b)(c : \text{TYPE } r_c).(a \rightarrow b \rightarrow c) \rightarrow a \rightarrow b \rightarrow c$$

$$\text{binapp } f \ x \ y = f \ x \ y$$

But that should be rejected because we cannot compile the call  $(f \ x \ y)$  without knowing how  $x$  and  $y$  are represented (and thus, how to retrieve them and pass them to  $f$ ). Indeed, it is rejected (by *both* FUN-I and FUN-E). On the other hand, perhaps surprisingly, this type for `binapp` is fine:

$$\text{binapp} : \forall(r_c : \text{Rep})(a : \text{TYPE } \text{PtrR})(b : \text{TYPE } \text{PtrR})(c : \text{TYPE } r_c).(a \rightarrow b \rightarrow c) \rightarrow a \rightarrow b \rightarrow c$$

Notice that `binapp` can be representation-polymorphic in the return type  $c$ . Because our compiler supports tail-call elimination,  $f$  is the one to return a value of type  $c$  to the caller, not `binapp`.

### 3.5 Arity Polymorphism

Arity and representations share problems (1) and (2), and thankfully, they share solutions, too. We add a second parameter to `TYPE` that describes the arity of the function, like this:<sup>6</sup>

$$\text{poly} : \forall(a : \text{TYPE } \text{PtrR } \text{Call}[2]).(\text{Int} \leadsto \text{Int} \leadsto a) \leadsto (a, a)$$

$$\text{poly} = \Lambda(a : \text{TYPE } \text{PtrR } \text{Call}[2]).\lambda(f : \text{Int} \leadsto \text{Int} \leadsto a). \text{let } g : \text{Int} \leadsto a = f \ 3 \text{ in } (g \ 5, g \ 4)$$

Here,  $a$ 's kind says that it only ranges over arity-2 types, of kind `TYPE PtrR Call[2]`, while  $f$ 's type  $\text{Int} \leadsto \text{Int} \leadsto a$  has kind `TYPE PtrR Call[4]`, hence  $f$  has arity 4 in total. Similarly,  $g : \text{Int} \leadsto a$  of kind `TYPE PtrR Call[3]` has arity 3. In short, *the calling convention of a function (how many arguments to pass simultaneously) is described by its kind*. Of course, this meant that we had to choose a particular arity for  $a$ , just as we had to choose a particular representation for the  $a$  and  $b$  in `binapp`'s type. We enforce this “particular arity” constraint not at the point of abstraction, but at the point of application. Here is the augmented FUN-E rule:

$$\frac{\Gamma \vdash e : \tau \leadsto \sigma \quad \Gamma \vdash e' : \tau \quad \Gamma \vdash \tau \text{ mono-rep} \quad \Gamma \vdash \tau \text{ mono-conv}}{\Gamma \vdash e e' : \sigma} \text{FUN-E}$$

where the new side condition  $\tau$  **mono-conv** ensures `TYPE`'s second parameter is statically known. Preparing an argument of function type is precisely the point at which a compiler must compile

<sup>4</sup>The paper *Levity Polymorphism* in fact describes *representation polymorphism*. For levity polymorphism, see Section 3.6.

<sup>5</sup>GHC classifies representations like `PtrR`, `IntR`, etc., as `RuntimeRep`; here, we shorten the name to just `Rep`.

<sup>6</sup>These kinds are somewhat simplified; the full story is in Section 4.

code for that argument that actually takes the number of arguments specified by its kind. Requiring the calling convention be monomorphic fixes that number statically; it can't be a type variable, say.

### 3.6 Evaluation Strategy and Levity

We have discussed *arity* and *representation*, two of our three design axes. But what of *levity*?

When compiling a function call  $f (1 + 1)$ , we must know the evaluation strategy to use: do we evaluate  $(1 + 1)$  *before* calling  $f$  or *after*, on-demand? Programming languages typically commit to a choice here, which is usually to evaluate  $(1 + 1)$  before the function call—the eager, call-by-value strategy. Haskell makes the opposite choice, implementing the lazy, call-by-need strategy;  $(1 + 1)$  is evaluated only when its value is needed, and only once. But regardless of a language's choice of evaluation strategy, some scenarios require the opposite. Programmers in eager languages sometimes use constructs like *Delay* and *Force* [Wadler et al. 1998] to embed lazy evaluation, and programmers in lazy languages manually introduce strictness (such as Haskell's *seq*) to force eager evaluation. So a compiler's intermediate language should support both lazy and eager evaluation, and ideally without favoring either over the other.

We can use types (or, more precisely, their kinds) to control evaluation order. For example, suppose we have types  $\text{Int}^L$  and  $\text{Int}^U$ , where the  $L$  stands for “Lifted” (the type has an extra bottom element,  $\perp$ , denoting divergent computation) and  $U$  for “Unlifted” (the type has no extra bottom element). Operationally, a value of lifted type must be represented as a pointer to a heap-allocated object because it may be an unevaluated thunk; a value of unlifted type may well still be represented by a pointer, but to the value itself, not a thunk. Variables of unlifted types cannot be bound to (divergent) computations, so they must be evaluated eagerly at binding-time; variables of lifted types may be bound to any unevaluated computation, so those bindings may be evaluated lazily. The types of an eager language (like OCaml) are all unlifted, whereas the types of a lazy language (like Haskell) are all lifted. But our  $\mathcal{IL}$  supports both, and hence can be a target for both.

Once again, we can track levity in the kinds, like this:

$$\text{Int}^L : \text{TYPE PtrR Eval}^L \quad \text{Int}^U : \text{TYPE PtrR Eval}^U \quad \text{Int\#} : \text{TYPE IntR Eval}^U$$

The levity of a type is all about its evaluation strategy. We already know primitive functions cannot be evaluated without calling them (Section 2.4), and thus a type has *either* a levity *or* an arity. We accordingly re-use the second component of the kind  $\text{TYPE}$ , which in Section 3.5 described the arity of a primitive function, using  $\text{Eval}^L$  and  $\text{Eval}^U$  to classify lifted and unlifted data types respectively. Reflecting this dual use, we describe the second argument of  $\text{TYPE}$  as the kind's *convention*, connoting “calling convention” for functions and “evaluation convention” for data.

Once again, we can be polymorphic in both levity (Section 4.4) and convention (Section 4.5). And we still keep the program compilable despite type erasure using suitable side conditions (Section 4.6). In the case of levity, the **mono-conv** premise in  $\text{FUN-E}$  (recall Section 3.5) is exactly what is needed to specify whether to compile this application using call-by-need or call-by-value.

### 3.7 From $\mathcal{IL}$ to $\mathcal{ML}$

Returning to arities, how can *poly* in Section 3.5 be compiled? In particular, since  $g$ 's kind specifies that it has arity 3, we must compile  $g$  to code that takes three arguments. So, before generating code, we  $\eta$ -expand  $g$ . But apparently we can't, at least not in the confines of  $\mathcal{IL}$ 's type system, because  $a$  is not an arrow type!

To have any hope of solving this problem, we must do more than say  $a$  is an arity-2 type. We must also spell out the representations of its two arguments, so that we can generate code for passing the  $\eta$ -expanded function arguments. In *poly*, the explicited calling convention looks like:

$$\begin{aligned} \text{poly} &= \Lambda(a : \text{TYPE PtrR Call}[\text{PtrR}, \text{IntR}]). \lambda(f : \text{Int} \rightsquigarrow \text{Int} \rightsquigarrow a). \\ &\quad \text{let } g : \text{Int} \rightsquigarrow a = f \ 3 \ \text{in } (g \ 5, g \ 4) \end{aligned}$$



$Kind \ni \kappa ::= \text{TYPE } \rho \ v$	
$Type \ni \tau, \sigma ::= t \mid T_p \mid T_d \mid \tau \rightsquigarrow \sigma \mid \forall \chi. \sigma$	$TyVar \ni \chi ::= t:\kappa \mid r \mid g \mid n$
$Representation \ni \rho ::= r \mid \text{PtrR} \mid \text{IntR} \mid \dots$	$PrimType \ni T_p ::= \text{Int\#} \mid \dots$
$Levity \ni \gamma ::= g \mid L \mid U$	$DataType \ni T_d ::= \text{Int}^\gamma \mid \dots$
$Convention \ni v ::= n \mid \text{Eval}^\gamma \mid \text{Call}[\alpha]$	$Arity \ni \alpha ::= \rho, \alpha \mid \varepsilon \mid \text{arity}(v)$
$Expr \ni e ::= x \mid c \mid I^\gamma e \mid \text{case } e \text{ of } I^\gamma x \rightarrow e' \mid \lambda x:\tau. e \mid e \ e' \mid \text{Clos}^\gamma e \mid \text{App } e \mid \Lambda \chi. e \mid e \ \phi$	
$Const \ni c ::= i \mid op$	$PrimOp \ni op ::= \text{error} \mid \dots$
$Answer \ni A ::= x \mid c \mid I^\gamma A \mid \text{Clos}^\gamma e \mid \Lambda \chi. A \mid A \ \phi$	$Erasable \ni \phi ::= \tau \mid \rho \mid \gamma \mid v$

Fig. 1. Syntax of  $\mathcal{IL}$ : An intermediate language with levity, representation, and arity polymorphism

The  $\text{Call}[\text{PtrR}, \text{IntR}]$  in  $a$ 's kind describes the representation of its two arguments.

Before code generation, we compile from  $\mathcal{IL}$  into a lower-level representation  $\mathcal{ML}$  (suggesting “machine language”).  $\mathcal{ML}$  is still statically typed, but much more coarsely:  $\mathcal{ML}$ 's types correspond to representations and calling conventions. So compiled polymorphic code might look like:

$$poly = \lambda(f:\text{PtrR}). \text{let } g : \text{PtrR} = \lambda(x:\text{PtrR}, y:\text{PtrR}, z:\text{IntR}). f(3, x, y, z) \\ \text{in } (\lambda(y:\text{PtrR}, z:\text{IntR}). g(5, y, z), \lambda(y:\text{PtrR}, z:\text{IntR}). g(4, y, z))$$

In  $\mathcal{ML}$  the functions are uncurried, and specify all their arguments and their representations. As such, every function is fully  $\eta$ -expanded, and every call is fully saturated. This representation is less convenient for the optimizer (for which we use  $\mathcal{IL}$ ), but is just right for the code generator. We discuss  $\mathcal{ML}$ , and the lowering transformation from  $\mathcal{IL}$  to  $\mathcal{ML}$ , in Section 6.

#### 4 THE INTERMEDIATE LANGUAGE ( $\mathcal{IL}$ )

Our intermediate language, which we call  $\mathcal{IL}$ , is an explicitly-typed  $\lambda$ -calculus based closely on System F, with its syntax (Fig. 1), formation of types (Fig. 2), and type system (Fig. 3). We cover the kind system of  $\mathcal{IL}$  in more detail in Sections 4.2 to 4.7.

The syntax for expressions  $e$  is given in Fig. 1, and includes the following forms.

- Variables  $x$ , and constants  $c$ .
- *Primitive integers* have type  $\text{Int\#} : \text{TYPE IntR Eval}^U$ , and are represented by a machine integer (hence  $\text{IntR}$  in the kind). There are numeric constants,  $i : \text{Int\#}$ , and primitive functions ( $op$ ) that operate over  $\text{Int\#}$ , such as  $\text{plus\#} : \text{Int\#} \rightsquigarrow \text{Int\#} \rightsquigarrow \text{Int\#}$ .
- *Boxed integers* have type  $\text{Int}^\gamma : \text{TYPE PtrR Eval}^\gamma$ , where the levity  $\gamma$  can be  $L, U$  or a levity variable  $g$ . A data-constructor application  $I^\gamma e : \text{Int}^\gamma$  allocates a box containing the value of  $e$  (see Section 3.6). Such  $I^\gamma$  boxes are unpacked by a pattern-matching  $\text{case}$ .
- *Primitive functions* are introduced and eliminated with the familiar forms  $\lambda x:\tau. e$  and  $e \ e'$ . As discussed in Section 2.2 we use a wavy arrow ( $\tau \rightsquigarrow \sigma$ ) to remind ourselves that in this  $\lambda$ -calculus functions have some unusual behavior: they enjoy unrestricted  $\eta$ -expansion. Multiple arguments, even though they are curried, are always passed simultaneously; e.g., if  $f : \text{Int\#} \rightsquigarrow \text{Int\#} \rightsquigarrow \text{Int\#} \rightsquigarrow \text{Int\#}$ , then  $f$  is always called with all three arguments.
- *Function closures* are introduced and eliminated with the forms  $\text{Clos}^\gamma e$  and  $\text{App } e$ , respectively. Similar to boxed integers,  $\text{Clos}^\gamma e$  allocates a new closure containing  $e$ , where  $\gamma$  records the levity of the computation responsible for that allocation, just as with  $I^\gamma$ .
- *Type abstraction*  $\Lambda \chi. e$  and *application*  $e \ \phi$ , have the same meaning as in System F. While their operational interpretation in  $\mathcal{IL}$  is analogous to  $\lambda x:\tau. e$  and  $e \ e'$ , these abstractions and applications are fully *erased* by compilation to a lower-level representation. The only new

Formation rules for kinds, levities, representations, conventions, arities:

$$\begin{array}{c}
\frac{\Gamma \vdash \rho \text{ rep} \quad \Gamma \vdash \nu \text{ conv}}{\Gamma \vdash \text{TYPE } \rho \nu \text{ kind}} \quad \Gamma \vdash \kappa \text{ kind} \\
\frac{}{\Gamma \vdash \rho \text{ rep}} \quad \frac{}{\Gamma, r \vdash r \text{ rep}} \quad \frac{}{\Gamma \vdash \text{PtrR rep}} \quad \frac{}{\Gamma \vdash \text{IntR rep}} \\
\frac{}{\Gamma \vdash \gamma \text{ lev}} \quad \frac{}{\Gamma, g \vdash g \text{ lev}} \quad \frac{}{\Gamma \vdash \text{L lev}} \quad \frac{}{\Gamma \vdash \text{U lev}} \\
\frac{}{\Gamma \vdash \nu \text{ conv}} \quad \frac{}{\Gamma, n \vdash n \text{ conv}} \quad \frac{\Gamma \vdash \gamma \text{ lev}}{\Gamma \vdash \text{Eval}^\gamma \text{ conv}} \quad \frac{\Gamma \vdash \alpha \text{ ari}}{\Gamma \vdash \text{Call}[\alpha] \text{ conv}} \\
\frac{\Gamma \vdash \rho \text{ rep} \quad \Gamma \vdash \alpha \text{ ari}}{\Gamma \vdash \rho, \alpha \text{ ari}} \quad \frac{}{\Gamma \vdash \varepsilon \text{ ari}} \quad \frac{\Gamma \vdash \nu \text{ conv}}{\Gamma \vdash \text{arity}(\nu) \text{ ari}} \\
\Gamma \vdash \alpha \text{ ari} \\
\text{Kinds of types } \Gamma \vdash \tau : \kappa \\
\frac{\Gamma \vdash \kappa \text{ kind}}{\Gamma, t : \kappa \vdash t : \kappa} \text{TVAR} \quad \frac{}{\Gamma \vdash \text{Int\#} : \text{TYPE IntR Eval}^{\text{U}} \text{Int\#}} \text{INT\#} \quad \frac{\Gamma \vdash \gamma \text{ lev}}{\Gamma \vdash \text{Int}^\gamma : \text{TYPE PtrR Eval}^\gamma \text{Int}} \text{INT} \\
\frac{\Gamma \vdash \tau_1 : \text{TYPE } \rho_1 \nu_1 \quad \Gamma \vdash \tau_2 : \text{TYPE } \rho_2 \nu_2}{\Gamma \vdash \tau_1 \rightsquigarrow \tau_2 : \text{TYPE PtrR Call}[\rho_1, \text{arity}(\nu_2)]} \text{FUN} \quad \frac{\Gamma \vdash \gamma \text{ lev} \quad \Gamma \vdash \tau : \text{TYPE } \rho \nu}{\Gamma \vdash \gamma\{\tau\} : \text{TYPE PtrR Eval}^\gamma \text{CLO}} \text{CLO} \\
\frac{\Gamma, \chi \vdash \sigma : \text{TYPE } \rho \nu \quad \Gamma \vdash \text{TYPE } \rho \nu \text{ kind}}{\Gamma \vdash \forall \chi. \sigma : \text{TYPE } \rho \nu} \text{FORALL} \quad \frac{\Gamma \vdash \tau : \kappa \quad \kappa = \kappa'}{\Gamma \vdash \tau : \kappa'} \text{K-CONV}
\end{array}$$

Reflexivity, transitivity, symmetry, and compatibility for  $\kappa = \kappa'$ , plus the following rules:

$$\text{arity}(\text{Eval}^\gamma) = \varepsilon \quad \text{arity}(\text{Call}[\alpha]) = \alpha$$

Monomorphism restrictions:

$$\frac{\Gamma \vdash \tau : \text{TYPE } \rho \nu \quad \vdash \rho \text{ rep}}{\Gamma \vdash \tau \text{ mono-rep}} \text{MONO-REP} \quad \frac{\Gamma \vdash \tau : \text{TYPE } \rho \nu \quad \vdash \nu \text{ conv}}{\Gamma \vdash \tau \text{ mono-conv}} \text{MONO-CONV}$$

Fig. 2. Kind and levity system of  $\mathcal{IL}$

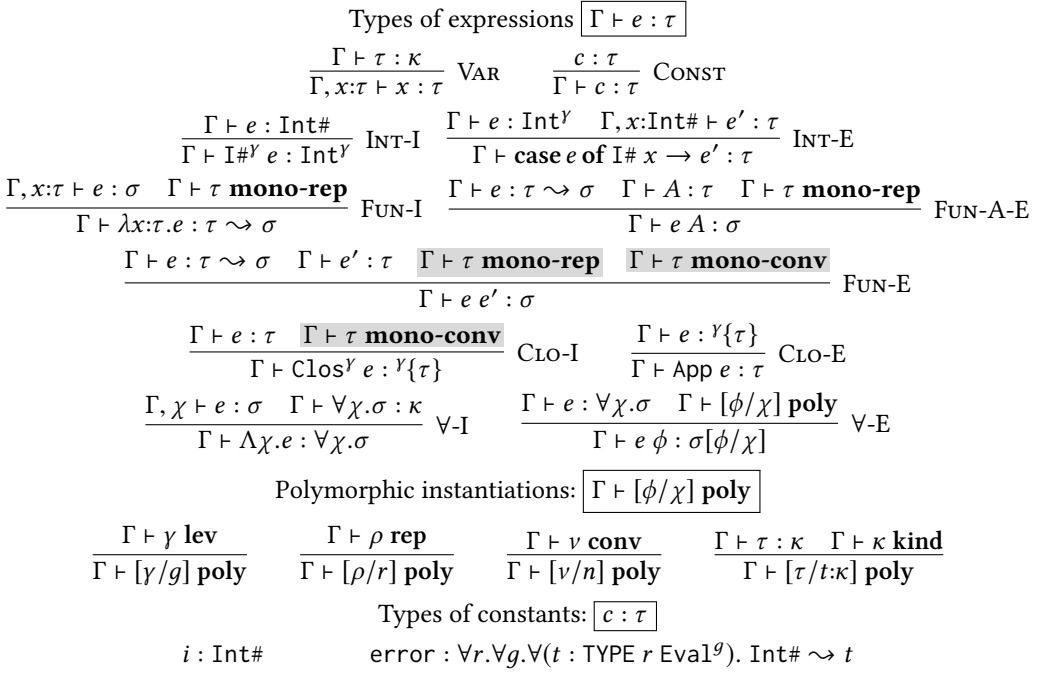
addition here is that binders  $\chi$  and arguments  $\phi$  range over four different sorts of type-level variables and arguments (all erasable). The expressiveness and restrictions on these different kinds of polymorphism are detailed in Sections 4.3 to 4.8. Type binders are annotated with their kinds, and other binders are distinguished by their naming convention; we use a Latin name for type-level variables and the corresponding Greek name for the syntactic category. For clarity, we sometimes state the sort of variable being quantified ( $r$ :Rep,  $g$ :Lev, or  $n$ :Conv).

A subset of these expressions are *Answers* (denoted by  $A$ , Fig. 1), meaning that they are possible results of evaluation. This classification accounts for the eventual *type erasure* mentioned above. The apparent redex  $(\Lambda \chi. A) \phi$  is an answer, because it is erased to  $A$  at compile-time. Answers include *closures*  $\text{Clos}^\gamma e$  as usual, but not primitive functions  $\lambda x : \tau. e$ ; recall that functions are *called* instead of *evaluated*, and thus cannot be answers of evaluation (Section 2.4).

The rules for well-formed types and kinds are given in Fig. 2, and include type variables  $t$ , primitive types  $T_p$  (of which we supply one,  $\text{Int\#}$ ), algebraic data types  $T_d$  (of which we supply one,  $\text{Int}$ ), polymorphic types  $\forall \chi. \sigma$ , primitive function types  $\tau \rightsquigarrow \sigma$  and a closure type  $\gamma\{\tau\}$ .

$\mathcal{IL}$ 's type system given in Fig. 3 ensures the following guarantees for well-behaved programs:

- *Static Compilation*: Every well-typed program can be compiled to a lower-level representation for a monomorphic machine, which models details such as specialized registers (for integers versus pointers) and function calls with multiple arguments (Theorem 3).
- *Type Safety*: If a well-typed program is equal to a number (as per Section 4.9), then its compiled code computes that number (Theorem 4). In particular, executing well-typed programs never

Fig. 3. Type system of  $\mathcal{IL}$ 

gets stuck, because the correct kind of register is used to store each value, and primitive functions are always called with the correct number and kind of arguments.

Static compilation requires some typing rules in Fig. 3 to explicitly refer to the *kinds of types*, as described in Fig. 2, that can be assigned to certain expressions. This shows up in the occasional  $\tau$  **mono-rep** and  $\tau$  **mono-conv** side conditions, which we elaborate in Section 4.6. Intuitively,  $\tau$  **mono-rep** captures the fact that  $\tau$  has a *statically-known, monomorphic representation*, and  $\tau$  **mono-conv** says that  $\tau$  has a *statically-known, monomorphic convention*.

#### 4.1 Simplifying Assumptions

To maintain a minimal presentation of  $\mathcal{IL}$ , we make many simplifying assumptions that reduce its number of features to a small representative core illustrating our main objectives. A realistic implementation will include more features, which can either be added as an extension to  $\mathcal{IL}$ , or encoded in terms of the features shown here. Our simplifying assumptions are as follows:

- Higher kinds (e.g.,  $\star \rightarrow \star$ ) are not included, but are a standard extension orthogonal to  $\mathcal{IL}$ .
- $\text{Int}\#$  is the only exemplary primitive unboxed type, which introduces the only non-pointer representation  $\text{IntR}$ . Other atomic unboxed types and representations—say, for floating-point numbers, characters, arrays, *etc.*—can be added straightforwardly, as can primitive operations on these types, like  $\text{plus}\#$  and  $\text{minus}\#$  of type  $\text{Int}\# \rightsquigarrow \text{Int}\# \rightsquigarrow \text{Int}\#$ . Unboxed tuple types  $(\# \tau_1, \dots, \tau_n \#)$  are an extension that introduces a compound representation roughly dual to the calling convention of functions:  $\text{Call}[\rho_1, \dots, \rho_n]$  describes the representations of values that a function needs to consume, while  $\text{Tuple}[\rho_1, \dots, \rho_n]$  [Eisenberg and Peyton Jones 2017] describes the representations of values that a tuple has within it.
- $\text{Int}^\gamma$  serves as the only example of an algebraic data type, which happens to be parameterized by a levity  $\gamma$  specifying whether the result of the constructor is evaluated eagerly (U) or

lazily (L). In general, there should be a way to declare new user-defined algebraic data types. These need not be levity polymorphic (*i.e.*, have a  $\gamma$  parameter), but  $\mathcal{IL}$  makes it possible to combine levity polymorphism with user-defined data types; see Section 4.8.

- There is no built-in **let**-binding construct but, as usual, a non-recursive **let** can be regarded as shorthand for  $\lambda$  and application:  $\text{let } x:\tau = e \text{ in } e' \triangleq (\lambda x:\tau. e') e$ . The typing rules for **let**-bindings can be derived from the rules for primitive function types  $\tau \rightsquigarrow \sigma$ . As such, **let**-bindings inherit similar side conditions on the type of the bound variable; see Section 4.6.
- **error** is the only source of computational effects in  $\mathcal{IL}$  as presented here. Recursive **let**-bindings, which are essential for practical functional programming, can be added straightforwardly, as can other computational effects, such as printing and state as in OCaml, with additional primitive operations.
- Because **error** is the only side-effect in  $\mathcal{IL}$ , call-by-name and call-by-need evaluation always give the same result [Ariola and Felleisen 1997]. So in  $\mathcal{IL}$ , we interpret “lifted” (L) as call-by-need for operational concerns, and as call-by-name for equational reasoning. To accommodate richer side-effects, the choice should instead be made explicit. This can be done directly in  $\mathcal{IL}$  by further dividing L into separate levities for call-by-name (CBN) and call-by-need (Need) evaluation. For example, evaluating  $\text{let } x:\text{Int}^\gamma = (\text{print "bye "}; \text{I}\#^\gamma 0) \text{ in } (\text{print "hi "}; x)$  prints: “bye hi ” with  $\gamma = \text{U}$ , “hi bye bye ” with  $\gamma = \text{CBN}$ , and “hi bye ” with  $\gamma = \text{Need}$ .

## 4.2 Kinds, Representations, Levities, and Conventions

A kind  $\kappa$  has the form  $\text{TYPE } \rho \, v$ , where  $\rho$  describes the *representation* of the type, and  $v$  describes its *convention*, that is, what operations are allowed on that type. Suppose  $x : \tau : \text{TYPE } \rho \, v$ ; that is,  $x$  is a term variable of type  $\tau$ , whose kind is  $\text{TYPE } \rho \, v$ . The representation  $\rho$  specifies the runtime representation of the value of  $x$ . Referring to Fig. 1,  $\rho$  can be:

- $\text{PtrR}$ , meaning that  $x$  is a pointer into the garbage-collected heap.
- $\text{IntR}$ , meaning that  $x$  is a machine integer (not a pointer).
- $r$ , a representation variable bound by a  $\forall$ ; that is, we support representation polymorphism [Eisenberg and Peyton Jones 2017].

The convention  $v$  describes the allowed operations on  $x$ , *i.e.*, how it can be consumed.  $v$  can be:

- $\text{Eval}^\text{U}$ , meaning that  $x$  cannot be bound to a computation like  $\perp$  (hence U for “Unlifted”). This kind is used for primitive values, and heap pointers that point directly to the value itself (such as a heap-allocated array).
- $\text{Eval}^\text{L}$ , meaning that  $x$  may be bound to a computation like  $\perp$  (hence L for “Lifted”). This kind is used for thunks, which might need evaluation to get its value, and might diverge doing so.
- $\text{Eval}^g$ , where  $g$  is a levity variable bound by  $\forall$ ; that is, we support levity polymorphism.
- $\text{Call}[\alpha]$ , meaning that  $x$  is a primitive function (not a thunk) with an *arity* described by  $\alpha$ . The arity of an  $\mathcal{IL}$  primitive function might be a fixed non-empty list  $\rho_1, \dots, \rho_m$ , so  $x$  takes precisely  $m \geq 1$  arguments,<sup>7</sup> with representations given by  $\rho_i$ . Otherwise, the arity is  $\rho_1, \dots, \rho_m, \text{arity}(v)$ , meaning that  $x$  takes *at least*  $m$  arguments with the listed representations, followed by possibly some more arguments given by the arity of  $v$ .
- $n$ , a convention variable bound by  $\forall$ ; that is, we support convention polymorphism.

The kind  $\text{TYPE } \text{PtrR } \text{Eval}^\text{L}$  expresses the uniform representation of a value in a lazy language: a pointer to a lifted (*i.e.*, possibly a thunk) object stored in the heap. Because this kind is so common,

<sup>7</sup>There is no type in  $\mathcal{IL}$  with the convention  $\text{Call}[]$ , but it can easily arise in extensions of  $\mathcal{IL}$ . In practice, unboxed tuple arguments are passed simultaneously in several registers. So the type  $(\# \text{Bool}, \text{Int}\#, \text{String}\#) \rightsquigarrow \text{Int}^\text{L}$  can be given the kind  $\text{TYPE } \text{PtrR } \text{Call}[\text{PtrR}, \text{IntR}, \text{PtrR}]$ . The nullary case of unboxed tuple arguments,  $(\# \#) \rightsquigarrow \text{Int}^\text{L}$ , can then be given the kind  $\text{TYPE } \text{PtrR } \text{Call}[]$ .

we often abbreviate it to  $\star$  for the default kind. In a call-by-value language we would instead define the default kind  $\star$  as  $\text{TYPE PtrR Eval}^U$ , and  $\text{Eval}^L$  would be used sparingly, if at all.

### 4.3 Calling Conventions in Kinds

We track the arity of a function type in its kind, as described in Sections 3.5 and 3.7. For example:

$$\begin{aligned} \text{Int}\# \rightsquigarrow \text{Int}\# & : \text{TYPE PtrR Call}[\text{IntR}] \\ \text{Int}\# \rightsquigarrow \text{Int}\# \rightsquigarrow \text{Int}\# & : \text{TYPE PtrR Call}[\text{IntR}, \text{IntR}] \\ \text{Int}\# \rightsquigarrow \text{Int}^L \rightsquigarrow \text{Int}\# & : \text{TYPE PtrR Call}[\text{IntR}, \text{PtrR}] \end{aligned}$$

The convention of each of these primitive function types has the form  $\text{Call}[\alpha]$ , where the arity given by  $\alpha$  describes the arguments needed to fully call functions of that type.

The formation rule for  $(\tau_1 \rightsquigarrow \tau_2)$  keeps track of these arities. Looking at rule FUN in Fig. 2, we see the kind of  $(\tau_1 \rightsquigarrow \tau_2)$  has a calling convention of  $\text{Call}[\rho_1, \text{arity}(v_2)]$ . Two special cases are:

$$\frac{\Gamma \vdash \tau_1 : \text{TYPE } \rho_1 \vee \quad \Gamma \vdash \tau_2 : \text{TYPE } \rho' \text{ Call}[\rho_2, \dots, \rho_m]}{\Gamma \vdash \tau_1 \rightsquigarrow \tau_2 : \text{TYPE PtrR Call}[\rho_1, \rho_2, \dots, \rho_m]} \quad \frac{\Gamma \vdash \tau_1 : \text{TYPE } \rho_1 \vee \quad \Gamma \vdash \tau_2 : \text{TYPE } \rho' \text{ Eval}^Y}{\Gamma \vdash \tau_1 \rightsquigarrow \tau_2 : \text{TYPE PtrR Call}[\rho_1]}$$

The representation of the first argument ( $\rho_1$ ) is that of  $\tau_1$ . The rest of the arguments come from the convention  $v_2$  of  $\tau_2$ , via the type-level operation *arity*, as defined in Fig. 2. It returns the arity  $\alpha$  of  $\text{Call}[\alpha]$ ; and the empty arity in the case of  $\text{Eval}^Y$ . But  $v_2$  might also be a variable  $n$ , and then *arity*( $n$ ) is stuck; that is why *arity*( $v$ ) is part of the syntax of  $\alpha$  in Fig. 1. Rule K-CONV allows calls to *arity* to be calculated whenever desired.<sup>8</sup>

### 4.4 Polymorphism in Levity and Representation

We are used to polymorphism over types, but we can gainfully employ polymorphism over levities, representations, and conventions, which is extremely useful in practice. For example, levity polymorphism lets us write some functions that work uniformly over both strict and lazy values. Adding two boxed integers can be defined as

$$\text{plus} : \forall g_1. \forall g_2. \forall g_3. \text{Int}^{g_1} \rightsquigarrow \text{Int}^{g_2} \rightsquigarrow \text{Int}^{g_3} \quad \text{plus} (\text{I}\# x) (\text{I}\# y) = \text{I}\# (\text{plus}\# x y)$$

which is short-hand for the following definition in  $\mathcal{IL}$ :

$$\begin{aligned} \text{plus} & : \forall g_1. \forall g_2. \forall g_3. \text{Int}^{g_1} \rightsquigarrow \text{Int}^{g_2} \rightsquigarrow \text{Int}^{g_3} \\ \text{plus} & = \Lambda g_1. \Lambda g_2. \Lambda g_3. \lambda (x' : \text{Int}^{g_1}). \lambda (y' : \text{Int}^{g_2}). \text{case } x' \text{ of } \text{I}\# x \rightarrow \text{case } y' \text{ of } \text{I}\# y \rightarrow \text{I}\#^{g_3} (\text{plus}\# x y) \end{aligned}$$

Notice that in this definition, the levities  $g_i$  of the argument and return types are statically unknown, so we must be able to pattern-match on and return values with unknown levities. Specifically, rule INT-E in Fig. 3 allows a case-expression to scrutinize an integer of arbitrary levity  $\gamma$ . Operationally, a levity-polymorphic **case** has to test the scrutinee to see if it is a thunk (in case the levity variable is instantiated to L), and if so evaluate it. In essence, we can interpret a **case** on an unknown levity as a lifted one because a **case** is always strict and, if it happens that  $g$  is U, the branch for handling a thunk is simply dead code.<sup>9</sup>

Similarly, suppose we had a primitive type of arrays,  $\text{Array}\#^Y$ , with kinding rule

$$\frac{\Gamma \vdash \gamma \text{ lev} \quad \Gamma \vdash \tau : \text{TYPE PtrR } v}{\Gamma \vdash \text{Array}\#^Y \tau : \text{TYPE PtrR Eval}^Y}$$

<sup>8</sup>Alternatively, we could require *arity*( $v$ ) to be fully calculated in the FUN kinding rule. This would let us remove *arity*( $v$ ) from the grammar of arities, but also forces an additional restriction on the formation of types and expressions, specifically FUN and FUN-I, to rule out *arity*( $n$ ). Such a restriction comes with the cost of breaking a pleasant property of  $\mathcal{IL}$ : except for the  $\forall$  quantifier, any type made from well-kinded types is itself well-kinded.

<sup>9</sup>We assume that the concrete, runtime representation of values (not thunks) is the same for both eager and lazy evaluation. This is true in GHC and seems likely in other systems supporting laziness, but it may not hold in some systems.

From a representation point of view, an  $\text{Array}^\gamma$  is represented by a pointer and contains pointers. The array itself can be lifted or unlifted, and (independently) can contain lifted or unlifted values. For example, the type  $\text{Array}^\#{}^L(\text{Array}^\#{}^U \text{Int}^U)$  is a lifted array of pointers, each of which points directly to an array of pointers to (boxed) integers. The ability to exclude the possibility of intermediate thunks in this data structure is very valuable in high-performance code, as a recent spate of GHC proposals shows [Eisenberg 2019; Graf 2020; Martin 2019a,b,c; Theriault 2019].

#### 4.5 Polymorphism in Convention

We may also be polymorphic in conventions. Consider the reverse-apply function

$$\text{revapp } x \ f = f \ x$$

For now, suppose that it returns  $\text{Int}\#$ . What type should *revapp* have? Here are two possibilities:

- (1)  $\text{revapp} : \forall (t : \text{TYPE PtrR Eval}^L). t \rightsquigarrow (t \rightsquigarrow \text{Int}\#) \rightsquigarrow \text{Int}\#$
- (2)  $\text{revapp} : \forall (t : \text{TYPE IntR Eval}^U). t \rightsquigarrow (t \rightsquigarrow \text{Int}\#) \rightsquigarrow \text{Int}\#$

We want to compile a function like *revapp* to a single block of efficient machine code. To do so, we *must know the representation of  $x$* , because we have to generate instructions to move  $x$  around. If  $x$  is represented by an integer, it will be passed in one sort of register; if a float, in another; if a pointer then yet another.<sup>10</sup> So we can choose (1) or (2), but not both.

On the other hand, consider these other possible types:

- (3)  $\text{revapp} : \forall (t : \text{TYPE PtrR Eval}^U). t \rightsquigarrow (t \rightsquigarrow \text{Int}\#) \rightsquigarrow \text{Int}\#$
- (4)  $\text{revapp} : \forall (t : \text{TYPE PtrR Call[IntR]}). t \rightsquigarrow (t \rightsquigarrow \text{Int}\#) \rightsquigarrow \text{Int}\#$

Since we are simply moving  $x$  around, but not otherwise acting upon it, we *can* simultaneously allow (1), (3), and (4). That is, we can be completely polymorphic in its convention  $n$ , thus:

- (4)  $\text{revapp} : \forall (n : \text{Conv}) (t : \text{TYPE PtrR } n). t \rightsquigarrow (t \rightsquigarrow \text{Int}\#) \rightsquigarrow \text{Int}\#$

What about the return type of  $f$ ? The code for *revapp* does not manipulate  $f$ 's return value at all (it does not even move it around, thanks to tail-call elimination), so it can be completely polymorphic in the representation or  $t_2$ , giving this type:

- (5)  $\text{revapp} : \forall (n : \text{Conv}) (r : \text{Rep}) (g : \text{Lev}) (t_1 : \text{TYPE PtrR } n) (t_2 : \text{TYPE } r \text{ Eval}^g). t_1 \rightsquigarrow (t_1 \rightsquigarrow t_2) \rightsquigarrow t_2$

But notice that, unlike the argument type  $t_1$ , *revapp* cannot be polymorphic in the convention of  $t_2$ , because a compiler needs to statically know the arity of the  $\lambda$ -bound argument  $f$  to generate code for  $(f \ x)$  inside *revapp* (Section 4.6). It can, however, be evaluated with any levity.

#### 4.6 Restrictions on Polymorphism

Of course we have the usual restrictions on polymorphism,<sup>11</sup> but  $\mathcal{IL}$ 's polymorphism introduces some new issues. We have already seen how unrestricted polymorphism is incompatible with efficient static code generation<sup>12</sup> in Section 4.5, where we cannot allow *revapp*'s type argument  $t_1$  to have a representation-polymorphic kind. A second restriction is demonstrated by:

$$\text{twice } f \ x = f \ (f \ x)$$

<sup>10</sup>Even if pointers occupy the same sort of register as integers, they are treated quite differently by the garbage collector, so the code generator treats them differently.

<sup>11</sup>E.g., every free variable (of every sort) that appears anywhere in the type checking judgment  $\Gamma \vdash e : \tau$  must be bound by  $\Gamma$ .

<sup>12</sup>Runtime code generation would allow the system to clone fresh code for each representationally-distinct instantiation of a function. But this is a pretty big hammer: only .NET does this. To keep things simple, we assume static code generation.



Should  $(f\ x)$  be eagerly or lazily evaluated? If it has a lifted type, then we can build a thunk for it, and pass that thunk to  $f$ . Otherwise, we must evaluate it before the call—remember, variables of unlifted types are always bound to values, never thunks. So we can give *twice* either of these types:

- (1)  $twice : \forall(t:\text{TYPE PtrR Eval}^L). (t \rightsquigarrow t) \rightsquigarrow t \rightsquigarrow t$
- (2)  $twice : \forall(t:\text{TYPE PtrR Eval}^U). (t \rightsquigarrow t) \rightsquigarrow t \rightsquigarrow t$

but we must choose: unlike *revapp*, *twice* cannot be polymorphic in  $t$ 's convention.

The restrictions on polymorphism used to ensure static compilability are embodied in the shaded premises in the type system of Fig. 3. The judgment  $\Gamma \vdash \tau$  **mono-rep**, defined in Fig. 2 checks that the representation of  $\tau$  is monomorphic; that is, that it mentions no variables. This is ensured by the empty context in the second premise of rule **MONO-REP**. There is an equivalent judgment  $\Gamma \vdash \tau$  **mono-conv** for conventions. Now returning to Fig. 3 we see the shaded premises:

- Rule **FUN-E**: for a general application, the argument type must be monomorphic in both the representation (so that we know how to store it while passing), and convention (so that we know when to evaluate it, or for first-class primitive function arguments, how to create it).
- Rule **FUN-A-E**: in the special case where the argument syntactically has the form of an answer, we can allow the argument type to be levity-polymorphic, since in this case there is no distinction between lazy or eager. If the application is lazy, then the argument does not need to be evaluated anyway. If application is eager, then the argument is already a value (or a variable that must be bound to a value), and again, does not need to be evaluated.
- Rule **CLO-I**: this rule boxes a primitive function so, just as in **FUN-E**, we must know that function's arity statically, so we can compile code for it that starts by taking the correct number of arguments.

We can justify these restrictions intuitively, but how do we know that these are the “right” restrictions? To answer that question we will show, in Section 6, how to compile  $\mathcal{IL}$  into our lower level  $\mathcal{ML}$ . In the translation from  $\mathcal{IL}$  to  $\mathcal{ML}$  in Fig. 8, we need exactly the shaded monomorphic restrictions of Fig. 3. If any of these restrictions were removed, then there would be expressions that are well-typed, yet uncompileable.

Our rules also include two additional, unshaded, monomorphism restrictions, in the **FUN-I** and **FUN-A-E** rules. These restrictions enforce an extra invariant on the environment  $\Gamma$ : *every variable in  $\Gamma$  has a monomorphic representation*. Besides making intuitive sense, this invariant could be necessary in a compiler accounting for more low-level details like storing free variables in a closure; doing so certainly requires knowing their representation. However, perhaps shockingly, the compilation scheme we give in Section 6 *does not* require any monomorphism restrictions in **FUN-I** and **FUN-A-E**: they could be deleted and yet all closed, well-typed expressions could still be compiled. This example suggests different compilers might need different restrictions on polymorphism. And from the reverse standpoint, other compilation schemes might allow for new and more adventurous possibilities for levity, representation, and convention polymorphism.

As an example of why different restrictions are needed for the **FUN-E** and **FUN-A-E** rules, recall the encoding of **let**-bindings from Section 4.1, which gives us the typing rules

$$\frac{\Gamma \vdash e : \tau \quad \Gamma, x : \tau \vdash e' : \sigma \quad \Gamma \vdash \tau \text{ mono-rep} \quad \Gamma \vdash \tau \text{ mono-conv}}{\Gamma \vdash \text{let } x : \tau = e \text{ in } e' : \sigma} \text{LET} \quad \frac{\Gamma \vdash A : \tau \quad \Gamma, x : \tau \vdash e : \sigma \quad \Gamma \vdash \tau \text{ mono-rep}}{\Gamma \vdash \text{let } x : \tau = A \text{ in } e : \sigma} \text{LET-A}$$

derived by composing **FUN-I** with **FUN-E** or **FUN-A-E**, respectively. That is, in general a **let** can bind a variables to expressions with both a monomorphic representation and convention; but convention-polymorphic answers can also be bound. Neither typing rule is more general than the other: **LET-A** only applies to some bound expressions (pre-evaluated answers), whereas **LET** only applies to bindings of some types (ones with statically-known conventions).

Now, consider these **let**-bindings, which follow the restrictions of either **LET** or **LET-A** above:

- |  |  |
|--|--|
| (1) $\text{let } x : \forall g. \text{Int}^g$<br>$\quad x = \Lambda g. \text{I}\#^g 0 \text{ in } \dots$                       | (3) $\text{let } x : \forall n. \forall (t : \text{TYPE PtrR } n). t$<br>$\quad x = y \text{ in } \dots$   |
| (2) $\text{let } x : \text{Int}^L$<br>$\quad x = \text{plus } L \ L \ L \ (\text{I}\#^L 1) (\text{I}\#^L 2) \text{ in } \dots$ | (4) $\text{let } x : \forall g. \text{Int}^g \rightsquigarrow \text{Int}^g$<br>$\quad x = \Lambda g. \lambda (y : \text{Int}^g). f \ (f' \ y) \text{ in } \dots$ |

Both (1) and (3) bind  $x$  to syntactically manifest answers with convention-polymorphic types ( $\forall g. \text{Int}^g : \text{TYPE PtrR Eval}^g$  and  $\forall n. \forall (t : \text{TYPE PtrR } n). t : \text{TYPE PtrR } n$ , respectively), which is well-typed by **LET-A**. In contrast, (2) and (3) bind  $x$  to non-answer expressions, so  $x$  needs a convention-monomorphic type (here,  $\text{Int}^L : \text{TYPE PtrR Eval}^L$  and  $\forall g. \text{Int}^g \rightsquigarrow \text{Int}^g : \text{TYPE PtrR Call}[\text{PtrR}]$ , respectively) in order to be well-typed by the **LET** rule.

#### 4.7 The FORALL Rule

The polymorphic quantifier  $\forall t : \kappa. \sigma$  has no impact on a type's kind: it just inherits the kind of  $\sigma$  (rule **FORALL** in Fig. 2). Intuitively, this is because these quantifiers will be totally *erased* by compilation, and have no impact on the final runtime code; the  $\forall$  is “invisible” to the lower-level machine.

However, now that variables may appear in kinds, we must be careful they do not escape their scope. For example, the following type is not well-kinded:

$$\forall r. \forall (t : \text{TYPE } r \text{ Eval}^L). t \rightsquigarrow t \quad :? \quad \text{TYPE PtrR Call}[r]$$

Here the representation  $r$  of the first parameter escapes in the calling convention of this primitive function type, because  $r$  is meant to be local to the type itself. This nonsense is prevented by the second premise of rule **FORALL** in Fig. 2 and the second premise of rule  $\forall$ -I in Fig. 3.

#### 4.8 User-Defined Types and Code Reuse

As we have seen, polymorphism over these type descriptors—representation, convention, and levity—increases the opportunities for code reuse. We can even have *data types* that are polymorphic over these descriptors, although this is beyond the scope of the  $\mathcal{LL}$  we formally describe here. For example definition of boxed integers (given in Section 3.2) might be generalized to

**data**  $\text{Int } (g : \text{Lev}) : \text{TYPE PtrR Eval}^g$  **where**  
 $\text{I}\# : \text{Int}\# \rightsquigarrow \text{Int } g$

Now, the  $\text{Int}$  type is parameterized by a chosen levity  $(g : \text{Lev})$ , which determines whether or not the boxed integers are evaluated eagerly or lazily.

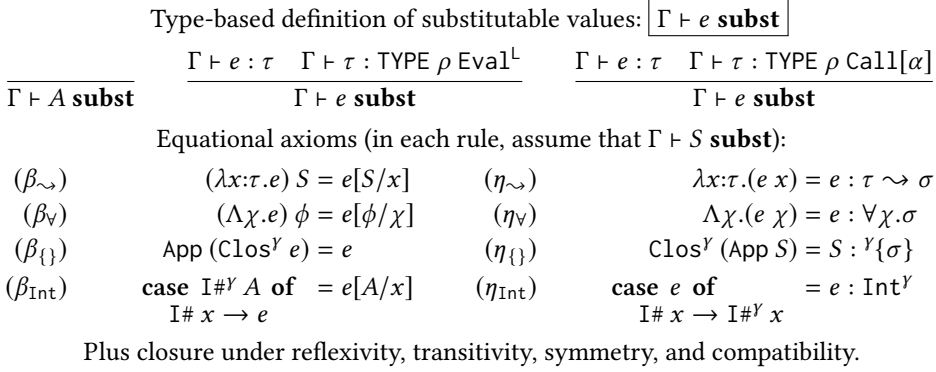
Polymorphism over a constructor's arguments and result can be combined within a single definition. For example, here is a further generalized definition of lists whose spine can be either eagerly or lazily evaluated ( $g$ ), and containing elements of any arity or evaluation strategy ( $n$ ):

**data**  $\text{List } (g : \text{Lev}) (n : \text{Conv}) (t : \text{TYPE PtrR } n) : \text{TYPE PtrR Eval}^g$  **where**  
 $\text{Nil} : \text{List } g \ n \ t$   
 $\text{Cons} : t \rightsquigarrow \text{List } g \ n \ t \rightsquigarrow \text{List } g \ n \ t$

Despite restrictions on polymorphism, we can define some levity-polymorphic functions over this type. For example, we could write the following polymorphic definition which is capable of summing up a list of integers with any combination of levities:

$\text{sum} \quad \quad \quad : \forall g_1 \ g_2 \ g_3. \text{List } g_1 \ \text{Eval}^{g_2} (\text{Int } g_2) \rightsquigarrow \text{Int } g_3$   
 $\text{sum Nil} \quad \quad \quad = \text{I}\# 0$   
 $\text{sum } (\text{Cons } (\text{I}\# \ x) \ xs) = \text{case sum } xs \text{ of } \text{I}\# \ y \rightarrow \text{I}\# (\text{plus}\# \ x \ y)$

$\text{sum}$ 's caller chooses the levity of the list's spine ( $g_1$ ), the list's elements ( $g_2$ ), and the result ( $g_3$ ). This is possible because  $\text{sum}$  is completely strict anyway; if it is given an evaluated list or an unevaluated

Fig. 4. Equational theory of  $\mathcal{IL}$ 

thunk, the entire thing will be added together before a value is returned. Therefore, the same code can be used for any combination of eager or lazy evaluation.

Other functions, such as the standard definition for mapping over a list:

$$\text{map } f \text{ Nil} = \text{Nil} \qquad \text{map } f (\text{Cons } x \text{ xs}) = \text{Cons } (f \ x) (\text{map } f \text{ xs})$$

cannot be so levity polymorphic. That's because the order of evaluation in  $\text{Cons } (f \ x) (\text{map } f \text{ xs})$  depends on the levity of  $(f \ x)$  and the recursive call  $(\text{map } f \text{ xs})$ . We can give this type assignment to  $\text{map}$  that specifies its result should be lazy, like in Haskell:

$$\text{map} : \forall g \ n \ (a : \text{TYPE } \text{PtrR } n) \ (b : \text{TYPE } \text{PtrR } \text{Eval}^L) . (a \rightsquigarrow b) \rightsquigarrow \text{List } g \ n \ a \rightsquigarrow \text{List } L \ \text{Eval}^L \ b$$

Note that the input list can store values of any convention at all, and like with  $\text{sum}$ , it can be either spine-strict or spine-lazy;  $\text{map}$  is strict in its second argument either way. Other evaluation orders for  $\text{map}$  can be specified by replacing one or both of  $L$  in  $(\text{List } L \ \text{Eval}^L \ b)$  with  $U$ . Although the definition of  $\text{map}$  appears the same in  $\mathcal{IL}$ , this change will compile to very different machine code.

#### 4.9 Equational Theory

The equational theory for  $\mathcal{IL}$ , as defined in Fig. 4, gives us a framework to reason about equivalent  $\mathcal{IL}$  expressions. We use this as the basis for correctness of compiling a high-level language to  $\mathcal{IL}$  (Theorem 2) and further on to a low-level language (Theorem 3). The rules for function closures (namely  $\beta_{\{\}}$  and  $\eta_{\{\}}$ ) and boxed integers ( $\beta_{\text{Int}}$  and  $\eta_{\text{Int}}$ ) are unsurprising, as are the rules for erasable abstractions ( $\beta_{\forall}$  and  $\eta_{\forall}$ ). More distinctive is  $\eta_{\rightsquigarrow}$ , which (as discussed in Section 2) allows unrestricted  $\eta$  in either direction for *any* expression of a primitive function type  $\tau \rightsquigarrow \sigma$ .

That leaves the reduction rule  $\beta_{\rightsquigarrow}$ . As is usual in a call-by-value  $\lambda$ -calculus, only some expressions—called *substitutable values*, denoted by the metavariable  $S$ —can be passed to a primitive function and substituted for its formal parameter by the  $\beta_{\rightsquigarrow}$  rule. The  $\beta_{\rightsquigarrow}$  rule only fires if the argument of the application is substitutable. Unlike most systems—using a purely syntactic definition of substitutable values— $\mathcal{IL}$  identifies these expressions by their *type* and *kind*, as defined by the rules for the  $\Gamma \vdash e \text{ subst}$  judgment. Using the type and kind of the argument to define substitutability, rather than only its syntax, allows us to integrate several different evaluation orders (call-by-value, call-by-name, etc.) within the same language.

Consider the rules for the  $\Gamma \vdash e \text{ subst}$  judgment. Firstly, we designate all *answers*  $A$  to be *substitutable*. Notice that, up to type erasure, each answer is considered a value in the call-by-value  $\lambda$ -calculus. In other cases, the criteria for substitutability of a typed term  $e : \tau$  depends on the convention of  $\tau$  rather than the syntax of  $e$ . For example, in a call-by-name setting, every expression can be substituted for a variable. So in  $\mathcal{IL}$ , *all* lazily-evaluated arguments—which have

the convention  $\text{Eval}^L$ —are substitutable values, and hence allow  $\beta_{\sim}$  to fire. In contrast, the only values in call-by-value languages are answers ( $A$ ), so variables of type  $\text{Eval}^U$  can only be substituted with answers. More generally, answers are always substitutable in all of the evaluation strategies we are interested in here, so we can say something more:  $A$  is substitutable for *all* conventions. This extra step is helpful in case we are dealing with an expression—like  $x$  or  $I\#^g i$ —which has an unknown convention but will inevitably be substitutable in any case.

For example, consider the different evaluation orders of a function call with lifted or unlifted arguments. On the one hand, we can express a lazy call to *plus* (Section 4.4) such as

$$(\lambda x:\text{Int}^L.e) (\text{plus } L (I\#^L 1) (I\#^L 2)) =_{\beta_{\sim}} e[\text{plus } L (I\#^L 1) (I\#^L 2)/x]$$

which substitutes the unevaluated argument for the parameter  $x$  right away according to  $\beta_{\sim}$ . This is possible because the argument has the type  $\text{Int}^L : \text{TYPE PtrR Eval}^L$ , and so it is substitutable by virtue of its type. On the other hand, the corresponding eager call would be

$$(\lambda x:\text{Int}^U.e) (\text{plus } U (I\#^U 1) (I\#^U 2)) = (\lambda x:\text{Int}^U.e) (I\#^U 3) =_{\beta_{\sim}} e[I\#^U 3/x]$$

Here, we cannot apply the  $\beta_{\sim}$  directly as before, because the argument is not substitutable: it has the type  $\text{Int}^U : \text{TYPE PtrR Eval}^U$  but is not an answer. Instead we must first evaluate the argument; the result is the answer  $(I\#^U 3)$ , and  $\beta_{\sim}$  can now fire.

## 5 COMPILATION TO $\mathcal{IL}$ FROM A HIGHER LEVEL

$\mathcal{IL}$  is, by design, a fairly low-level language making fine distinctions about representation, calling conventions, evaluation orders, and so on. This makes it a target for *both* eager *and* lazy languages. To see how, we now give translations for call-by-name and call-by-value System F into  $\mathcal{IL}$ .

### 5.1 Call-by-Name System F to $\mathcal{IL}$

To translate call-by-name System F into  $\mathcal{IL}$ , we begin by picking a single “uniform”  $\mathcal{IL}$  kind  $\star$  that captures all the types of the source language, namely  $\star = \text{TYPE PtrR Eval}^L$ . Each source-language type  $\tau$  translates to an  $\mathcal{IL}$  type  $\llbracket \tau \rrbracket$  of kind  $\star$ ; that is, a pointer to a lifted value, perhaps a thunk.

Fig. 5 gives this type translation. To get the correct call-by-name semantics for numbers, we use the boxed integer type  $\text{Int}^L$ , which happily has the correct kind. However, even though the primitive function type  $\llbracket \tau \rrbracket \rightsquigarrow \llbracket \sigma \rrbracket$  has the correct call-by-name semantics, it has the wrong kind  $\text{TYPE PtrR Call}[\text{PtrR}]$ , so the translation coerces  $\text{Call}[\text{PtrR}]$  to  $\text{Eval}^L$  with a closure type. Polymorphic type abstraction is unchanged, only clarified that bound type variable ranges over  $\star$ .

To compile expressions, we only need to expand out the additions prescribed by the translation of types. Numeric constants need to be boxed, functions and their calls need the explicit coercions to and from closures, and bound type variables are annotated with their uniform kind.

### 5.2 Call-by-Value System F to $\mathcal{IL}$

We can compile a call-by-value version of system F using virtually the same procedure as above. Again, we need to decide on a uniform kind that is suitable for each source-level type, which is  $\star = \text{TYPE PtrR Eval}^U$  for call-by-value evaluation. As before, we can compile source-level types following the invariant that  $\llbracket \tau \rrbracket : \star$ , as again shown in Fig. 5.

Note that we still compile integers to a boxed type, so that all values are represented uniformly by a pointer, but this time we make it unlifted to reflect the call-by-value semantics. Function types are still wrapped in a closure, as in the call-by-name case, but this time unlifted ones.

The translation of polymorphism is more complex, however, due a mismatch with the semantics of call-by-value System F. With call-by-value evaluation, the abstraction  $\Lambda t.\perp$  is a value, even though  $\perp$  diverges, whereas in call-by-name they would be  $\eta$ -equivalent. We need to make sure

$$\begin{array}{c}
\text{Call-by-name where } \star = \text{TYPE PtrR Eval}^L \\
\begin{array}{llll}
\llbracket \text{Int} \rrbracket \triangleq \text{Int}^L & \llbracket t \rrbracket \triangleq t & \llbracket x \rrbracket \triangleq x & \llbracket i \rrbracket \triangleq \text{I}^\#^L i \\
\llbracket \tau \rightarrow \sigma \rrbracket \triangleq {}^L\{\llbracket \tau \rrbracket \leadsto \llbracket \sigma \rrbracket\} & \llbracket \lambda x:\tau.e \rrbracket \triangleq \text{Clos}^L(\lambda x:\llbracket \tau \rrbracket.\llbracket e \rrbracket) & \llbracket e \ e' \rrbracket \triangleq (\text{App } \llbracket e \rrbracket) \llbracket e' \rrbracket & \\
\llbracket \forall t.\tau \rrbracket \triangleq \forall t:\star.\llbracket \tau \rrbracket & \llbracket \Lambda t.e \rrbracket \triangleq \Lambda t:\star.\llbracket e \rrbracket & \llbracket e \ \tau \rrbracket \triangleq \llbracket e \rrbracket \llbracket \tau \rrbracket & 
\end{array} \\
\\
\text{Call-by-value where } \star = \text{TYPE PtrR Eval}^U \\
\begin{array}{llll}
\llbracket \text{Int} \rrbracket \triangleq \text{Int}^U & \llbracket t \rrbracket \triangleq t & \llbracket x \rrbracket \triangleq x & \llbracket i \rrbracket \triangleq \text{I}^\#^U i \\
\llbracket \tau \rightarrow \sigma \rrbracket \triangleq {}^U\{\llbracket \tau \rrbracket \leadsto \llbracket \sigma \rrbracket\} & \llbracket \lambda x:\tau.e \rrbracket \triangleq \text{Clos}^U(\lambda x:\llbracket \tau \rrbracket.\llbracket e \rrbracket) & \llbracket e \ e' \rrbracket \triangleq (\text{App } \llbracket e \rrbracket) \llbracket e' \rrbracket & \\
\llbracket \forall t.\tau \rrbracket \triangleq {}^U\{\forall t:\star.\llbracket \tau \rrbracket\} & \llbracket \Lambda t.e \rrbracket \triangleq \text{Clos}^U(\Lambda t:\star.\llbracket e \rrbracket) & \llbracket e \ \tau \rrbracket \triangleq (\text{App } \llbracket e \rrbracket) \llbracket \tau \rrbracket & 
\end{array}
\end{array}$$

Fig. 5. Compiling call-by-name and call-by-value System F to  $\mathcal{IL}$ 

Call-by-name substitutable values:

$$V ::= e$$

Equational axioms (for both call-by-name and -value definitions of  $V$ ):

$$(\beta_{\rightarrow}) \quad (\lambda x:\tau.e) \ V = e[V/x]$$

$$(\beta_{\forall}) \quad (\Lambda t.e) \ \tau = e[\tau/t]$$

$$(name) \quad (\lambda x:\tau.e \ x) \ e' = e \ e'$$

Call-by-value substitutable values:

$$V ::= x \mid \lambda x:\tau.e \mid \Lambda t.e$$

$$(\eta_{\rightarrow}) \quad \lambda x:\tau.(V \ x) = V : \tau \rightarrow \sigma$$

$$(\eta_{\forall}) \quad \Lambda t.(V \ t) = V : \forall t.\sigma$$

Fig. 6. Equational theory of System F; call-by-value and call-by-name

that this abstraction is still a value even after the polymorphic  $\Lambda$  is erased. For that reason, we must introduce the additional call-by-value closure which is preserved to runtime.

The compilation of call-by-value expressions is nearly the same as call-by-name expressions. Besides swapping  $L$  for  $U$ , the only difference is in for polymorphic abstractions and instantiations. These have an extra closure that signifies call-by-value evaluation and, more importantly, makes sure that the value  $\Lambda t.e$  in System F compiles to a value, namely  $\text{Clos}^U(\Lambda t:\star.\llbracket e \rrbracket)$ , which is essential when side effects or non-termination enters the picture. For example, evaluating  $\Lambda t.\text{error } t$  returns immediately in call-by-value System F, but the corresponding  $\Lambda t:\star.\text{error PtrR } U \ t \ i$  causes an error in  $\mathcal{IL}$ . Intuitively, the explicit closure and application serve to codify the standard definition of type erasure for call-by-value System F, which traditionally erases  $\Lambda t.e$  into  $\lambda().e$ .

### 5.3 Correctness of Source-to- $\mathcal{IL}$ Compilation

Compiling call-by-name and call-by-value System F into  $\mathcal{IL}$  are both correct: type checking and equalities are preserved by translation. We assume the standard type system for System F and semantics in Fig. 6. Note that the *name* axiom, which gives a name to the argument of a function, corresponds to a right-to-left evaluation order for the call-by-value semantics, and is a consequence of  $\beta_{\rightarrow}$  in the call-by-name semantics. The preservation of types is straightforward, where the compilation of a typing environment  $\llbracket \Gamma \rrbracket$  is defined pointwise.

**Theorem 1** (Type Preservation). *If  $\Gamma \vdash e : \tau$  is derivable then so is  $\llbracket \Gamma \rrbracket \vdash \llbracket e \rrbracket : \llbracket \tau \rrbracket$ .*

**Theorem 2** (Soundness and Completeness).  *$\Gamma \vdash e = e' : \tau$  in call-by-name System F if and only if  $\llbracket \Gamma \rrbracket \vdash \llbracket e \rrbracket = \llbracket e' \rrbracket : \llbracket \tau \rrbracket$  via call-by-name compilation, and likewise for call-by-value.*

## 6 COMPILATION FROM $\mathcal{IL}$ TO A LOWER LEVEL

In order to illustrate how our levity-polymorphic intermediate language  $\mathcal{IL}$  might be compiled to a more conventional, lower-level language, we introduce an abstract machine that supports

$PrimRep \ni \pi ::= PtrR \mid IntR$	$Strictness \ni \psi ::= L \mid U$	$KnownConv \ni \eta ::= Eval^\psi \mid Call[\overline{\pi}]$
$Expr \ni e ::= W \mid W(\overline{a}) \mid App\ e(\overline{a}) \mid \text{case } e \text{ of } I\#(x_{IntR}) \rightarrow e' \mid \text{let } x_\pi^\psi = e \text{ in } e'$		
$Reference \ni R ::= I\#(a) \mid Clos^n W \mid \lambda(\overline{x_\pi}).e$		
$Value \ni V ::= c \mid I\#(a) \mid Clos^n W$	$WHNF \ni W ::= R \mid a$	$Arg \ni a ::= c \mid x_\pi$
$StackCont \ni K ::= \varepsilon \mid App(\overline{a}); K \mid \text{case } I\#(x_{IntR}) \rightarrow e; K \mid \text{let } x_\pi \text{ in } e; K \mid \text{set } x; K$		
$Store \ni H ::= \varepsilon \mid [x := R]H \mid [x := \text{memo } e]H$		
$MachineState \ni m ::= \langle e \mid K \mid H \rangle \mid \text{error}(n)$	$Final \ni Fin ::= \langle V \mid \varepsilon \mid H \rangle \mid \text{error}(n)$	
$(PshApp)$	$\langle App\ e(\overline{a}) \mid K \mid H \rangle \mapsto \langle e \mid App(\overline{a}); K \mid H \rangle$	
$(PshCase)$	$\langle \text{case } e \text{ of } I\#(x_{IntR}) \rightarrow e' \mid K \mid H \rangle \mapsto \langle e \mid \text{case } I\#(x_{IntR}) \rightarrow e'; K \mid H \rangle$	
$(PshLet)$	$\langle \text{let } x_\pi^U = e \text{ in } e' \mid K \mid H \rangle \mapsto \langle e \mid \text{let } x_\pi \text{ in } e'; K \mid H \rangle$	
$(LAlloc)$	$\langle \text{let } x_{PtrR}^L = e' \text{ in } e \mid K \mid H \rangle \mapsto \langle e[y_{PtrR}/x_{PtrR}] \mid K \mid [y := \text{memo } e']H \rangle$	
$(Call)$	$\langle (\lambda(\overline{x_\pi}).e)(\overline{a}) \mid K \mid H \rangle \mapsto \langle \overline{e[a/x_\pi]} \mid K \mid H \rangle$	
$(Apply)$	$\langle Clos^n W \mid App(\overline{a}); K \mid H \rangle \mapsto \langle W(\overline{a}) \mid K \mid H \rangle$	(if $n =  \overline{a} $ )
$(Unbox)$	$\langle I\#(a) \mid \text{case } I\#(x_{IntR}) \rightarrow e; K \mid H \rangle \mapsto \langle e[a/x_{IntR}] \mid K \mid H \rangle$	
$(Move)$	$\langle c \mid \text{let } x_\pi \text{ in } e; K \mid H \rangle \mapsto \langle e[c/x_\pi] \mid K \mid H \rangle$	
$(SAlloc)$	$\langle R \mid \text{let } x_{PtrR} \text{ in } e; K \mid H \rangle \mapsto \langle e[y_{PtrR}/x_{PtrR}] \mid K \mid [y := R]H \rangle$	
$(Fun)$	$\langle y_{PtrR}(\overline{a}) \mid K \mid [y := R]H \rangle \mapsto \langle R(\overline{a}) \mid K \mid [y := R]H \rangle$	
$(Look)$	$\langle y_{PtrR} \mid K \mid [y := R]H \rangle \mapsto \langle R \mid K \mid [y := R]H \rangle$	
$(Force)$	$\langle y_{PtrR} \mid K \mid [y := \text{memo } e]H \rangle \mapsto \langle e \mid \text{set } y; K \mid [y := \text{memo } e]H \rangle$	
$(Memo)$	$\langle R \mid \text{set } y; K \mid [y := \text{memo } e]H \rangle \mapsto \langle R \mid K \mid [y := R]H \rangle$	
$(Error)$	$\langle \text{error}(n) \mid K \mid H \rangle \mapsto \text{error}(n)$	

Fig. 7. Syntax and semantics for  $\mathcal{ML}$ 

non-uniform representations and function arities, both of which must be monomorphic. The idea is to model a realistic machine architecture with multiple basic kinds of data representations (e.g., pointers vs. integers), and where primitive functions are passed multiple arguments but are otherwise first order (though primitive function pointers may be passed as arguments and invoked). The syntax of this language, which we call  $\mathcal{ML}$ , is given in Fig. 7, along with its abstract machine.<sup>13</sup>

The syntax of  $\mathcal{ML}$  is restricted from the more  $\lambda$ -calculus-inspired  $\mathcal{IL}$  in several ways:

- Expressions follow the A-normal form (ANF) convention [Sabry and Felleisen 1993]: all arguments  $a$  are either variables or constants. To support ANF,  $\mathcal{ML}$  has a **let** construct.
- Primitive functions  $(\lambda(\overline{y_\pi}).e)$  and calls  $(W(\overline{a}))$  have an explicit arity and pass multiple arguments at once, but cannot be partially applied.
- An applied function cannot be an arbitrary expression; it must be a weak head-normal form, namely a reference to a  $\lambda$ , variable, or constant. Every application  $W(\overline{a})$  can be resolved in at most two steps: lookup  $W$  if it's a variable, then apply  $W$  if it's a  $\lambda$  or constant (like error).

As such, it is impossible to chain several calls in a row. For example,  $f(1)(2)$  is not a legal expression in  $\mathcal{ML}$ . If  $f$  is an arity 2 primitive function, it must be called as  $f(1, 2)$ ; if it is an arity-1 primitive function returning a closure of an arity-1 primitive function, it must be called as  $(App\ f(1))(2)$ . The

<sup>13</sup>For aficionados of GHC,  $\mathcal{IL}$  is like GHC's Core language, while  $\mathcal{ML}$  is like the STG language [Peyton Jones 1992].



number of arguments passed at once is explicitly fixed in each  $\lambda$ -abstraction and call site. In other words,  $\mathcal{ML}$  does not support polymorphism of function arity.

$\mathcal{ML}$ 's syntax includes annotations that make the semantically important information in types explicit, so the syntax of programs makes it clear how they are executed. In particular, each variable is annotated with its representation  $\pi$ , which must be either a pointer ( $\text{PtrR}$ ) or integer ( $\text{IntR}$ ), corresponding to an assignment to an appropriate register. In other words, all variables are permanently assigned a representation—intuitively, stored in either an integer or address register. By design,  $\mathcal{ML}$  does not support polymorphism over these representations because the different types of machine registers are distinct, and the choice made is fixed in the code.

Additionally, each **let** binding is annotated as either eager ( $\text{U}$ ) or lazy ( $\text{L}$ ). This controls whether the right-hand side of the **let** is evaluated first before being bound to the variable, or bound first and evaluated later as needed. Again, this decision about evaluation order must be statically chosen for each **let**, so  $\mathcal{ML}$  does not support polymorphism over evaluation order.

### 6.1 The Semantics of $\mathcal{ML}$

Executing an  $\mathcal{ML}$  program involves a machine configuration of the form  $\langle e \mid K \mid H \rangle$ , where  $e$  is the expression being evaluated,  $K$  is the *continuation* or *call stack* of evaluation, and  $H$  is the heap for storing allocated memory. Heaps may contain both values ( $[x := V]H$ ) or unevaluated thunks ( $[x := \text{memo } e]H$ ). Both call stacks and heaps are conventional, except that a stack may contain an application of many arguments in a single stack frame, like  $\text{App}(\bar{a}); K$ . The other cases of stack frames include a **case**, a strict **let** binding, and a **set** construct to memoize thunk evaluation.

Many of the steps of the machine are also conventional, including those for pushing stack frames ( $\text{PshApp}$ ,  $\text{PshCase}$ ,  $\text{PshLet}$ ) and allocating memory ( $\text{LAlloc}$ ,  $\text{SAlloc}$ ), but note that we include cases for *both* lazy bindings ( $\text{LAlloc}$ ) and strict ones ( $\text{PshLet}$ ,  $\text{SAlloc}$ ). Next we have the rules for performing interesting reductions. *Apply* resolves the application of a closure by extracting the primitive function it contains, and *Call* calls a primitive function directly. Note that *Apply* can check that the number of arguments matches the arity of the closure at runtime (and potentially respond appropriately if they do not match, as we do later in Section 7). Instead, *Call* is merely undefined when the arguments don't match the bound parameters, representing a type or memory unsafe error. In addition, we have *Move* for moving a constant into an appropriate variable (corresponding to a register) and *Unbox* for extracting the contents of a boxed integer. Finally, we have the rules for handling pointer variables at runtime. *Fun* expects a function pointer to map to a value. For other pointers, we have to check if it is evaluated, to either *Look* up values or else *Force* thunks.<sup>14</sup> When a forced thunk returns a value, it is *Memoized* to share the result on future uses.

### 6.2 Compilation

Compilation from  $\mathcal{IL}$  to the low-level machine language  $\mathcal{ML}$  is given in Fig. 8. The top-level translation is  $\mathcal{E}_\nu \llbracket e \rrbracket_\theta^\Gamma$ , which compiles a typed expression  $\Gamma \vdash e : \tau$  given  $\tau$ 's convention is  $\nu$ . The environment  $\theta$  is a mapping from  $\mathcal{IL}$  variables to  $\mathcal{ML}$  arguments (either constants or representation-annotated variables) written as  $[a_1/x_1] \dots [a_n/x_n]$ .

A key part in understanding the compilation in Fig. 8 is to remember the distinction between calling and evaluating. In our system, only expressions with types like  $\text{Int}^\#$ ,  $\text{Int}^\vee$ , and  ${}^\vee\{\tau\}$  can be evaluated. In contrast, expressions with types like  $\tau \rightsquigarrow \sigma$  can only be called. Implementing this distinction is the main role of  $\mathcal{E}_\nu \llbracket e \rrbracket_\theta^\Gamma$ , which takes into account the convention  $\nu$  of  $e$ : if it is  $\text{Eval}^\vee$

<sup>14</sup>Note that this uniform check on pointers  $y_{\text{PtrR}}$  is needed to support levity polymorphism for types like  $\text{Int}^\mathcal{G}$  and  ${}^\mathcal{G}\{\tau\}$ . In a more practical compiler, we could have specialized code that avoids a check when it is statically known, due to type checking that  $y_{\text{PtrR}}$  must be unlifted, so that the *Look* step always applies without a dynamic check. Thus, a language which is call-by-value by default does not have to pay the runtime penalty for thunks unless they are actually being used.

**In the following, all equations are tried left-to-right, top-to-bottom.**

Top-level eta expansion:

$$\mathcal{E}_V \llbracket A \rrbracket_\theta^\Gamma \triangleq \mathcal{A} \llbracket A \rrbracket_\theta^\Gamma \quad \mathcal{E}_{\text{Eval}^\psi} \llbracket e \rrbracket_\theta^\Gamma \triangleq C \llbracket e \rrbracket_\theta^\Gamma(\varepsilon) \quad \mathcal{E}_{\text{Call}[\bar{\pi}]} \llbracket e \rrbracket_\theta^\Gamma \triangleq \lambda(\bar{x}_\pi). C \llbracket e \rrbracket_\theta^\Gamma(\bar{x}_\pi)$$

Constants and variables:

$$\begin{aligned} C \llbracket i \rrbracket_\theta^\Gamma(\varepsilon) &\triangleq i & C \llbracket x \rrbracket_\theta^{\Gamma, x:\tau}(\varepsilon) &\triangleq \theta(x) & (\text{if } \Gamma \vdash \tau \xrightarrow{\text{conv}} \text{Eval}^\psi) \\ C \llbracket \text{error} \rrbracket_\theta^\Gamma(a) &\triangleq \text{error}(a) & C \llbracket x \rrbracket_\theta^{\Gamma, x:\tau}(\bar{a}) &\triangleq (\theta(x))(\bar{a}) & (\text{if } \Gamma \vdash \tau \xrightarrow{\text{conv}} \text{Call}[\bar{\pi}]) \end{aligned}$$

Applications (the following equations are tried top-to-bottom):

$$\begin{aligned} C \llbracket \text{App } e \rrbracket_\theta^\Gamma(\bar{a}) &\triangleq \text{App } (C \llbracket e \rrbracket_\theta^\Gamma(\varepsilon))(\bar{a}) \\ C \llbracket e \phi \rrbracket_\theta^\Gamma(\bar{a}) &\triangleq C \llbracket e \rrbracket_\theta^\Gamma(\bar{a}) \\ C \llbracket e A \rrbracket_\theta^\Gamma(\bar{a}) &\triangleq C \llbracket e \rrbracket_\theta^\Gamma(\mathcal{A} \llbracket A \rrbracket_\theta^\Gamma, \bar{a}) & (\text{if } \mathcal{A} \llbracket A \rrbracket_\theta^\Gamma = x_\pi \text{ or } \mathcal{A} \llbracket A \rrbracket_\theta^\Gamma = c) \\ C \llbracket e A \rrbracket_\theta^\Gamma(\bar{a}) &\triangleq \text{let } x_{\text{PtrR}}^U = \mathcal{A} \llbracket A \rrbracket_\theta^\Gamma \text{ in } C \llbracket e \rrbracket_\theta^\Gamma(x_{\text{PtrR}}, \bar{a}) \\ C \llbracket e e' \rrbracket_\theta^\Gamma(\bar{a}) &\triangleq \text{let } x_{\pi}^{\text{lev}(\eta)} = \mathcal{E}_\eta \llbracket e' \rrbracket_\theta^\Gamma \text{ in } C \llbracket e \rrbracket_\theta^\Gamma(x_\pi, \bar{a}) & (\text{if } \Gamma \vdash e':\tau, \Gamma \vdash \tau \xrightarrow{\text{rep}} \pi, \text{ and } \Gamma \vdash \tau \xrightarrow{\text{conv}} \eta) \end{aligned}$$

Boxing and unboxing:

$$\begin{aligned} C \llbracket I\#^Y A \rrbracket_\theta^\Gamma(\varepsilon) &\triangleq \mathcal{A} \llbracket I\#^Y A \rrbracket_\theta^\Gamma \\ C \llbracket I\#^Y e \rrbracket_\theta^\Gamma(\varepsilon) &\triangleq \text{let } x_{\text{IntR}}^U = C \llbracket e \rrbracket_\theta^\Gamma(\varepsilon) \text{ in } I\#(x_{\text{IntR}}) \\ C \llbracket \text{case } e' \text{ of } I\# x \rightarrow e \rrbracket_\theta^\Gamma(\bar{a}) &\triangleq \text{case } C \llbracket e' \rrbracket_\theta^\Gamma(\varepsilon) \text{ of } I\#(x_{\text{IntR}}) \rightarrow C \llbracket e \rrbracket_{[x_{\text{IntR}}/x]_\theta}^{\Gamma, x:\text{Int}\#}(\bar{a}) \end{aligned}$$

Abstractions:

$$\begin{aligned} C \llbracket \lambda x:\sigma. e \rrbracket_\theta^\Gamma(a', \bar{a}) &\triangleq C \llbracket e \rrbracket_{[a'/x]_\theta}^{\Gamma, x:\sigma}(\bar{a}) \\ C \llbracket \Lambda \chi. e \rrbracket_\theta^\Gamma(\bar{a}) &\triangleq C \llbracket e \rrbracket_\theta^{\Gamma, \chi}(\bar{a}) \\ C \llbracket \text{Clos}^Y e \rrbracket_\theta^\Gamma(\varepsilon) &\triangleq \mathcal{A} \llbracket \text{Clos}^Y e \rrbracket_\theta^\Gamma \end{aligned}$$

Type erasure of answers:

$$\begin{aligned} \mathcal{A} \llbracket c \rrbracket_\theta^\Gamma &\triangleq c & \mathcal{A} \llbracket x \rrbracket_\theta^\Gamma &\triangleq \theta(x) & \mathcal{A} \llbracket I\#^Y A \rrbracket_\theta^\Gamma &\triangleq I\#(\mathcal{A} \llbracket A \rrbracket_\theta^\Gamma) \\ \mathcal{A} \llbracket A \phi \rrbracket_\theta^\Gamma &\triangleq \mathcal{A} \llbracket A \rrbracket_\theta^\Gamma & \mathcal{A} \llbracket \Lambda \chi. A \rrbracket_\theta^\Gamma &\triangleq \mathcal{A} \llbracket A \rrbracket_\theta^\Gamma \\ \mathcal{A} \llbracket \text{Clos}^Y e \rrbracket_\theta^\Gamma &\triangleq \text{Clos}^{\text{arity}(\eta)} \mathcal{E}_{\text{Call}[\text{arity}(\eta)]} \llbracket e \rrbracket_\theta^\Gamma & (\text{if } \Gamma \vdash e:\tau \text{ and } \Gamma \vdash \tau \xrightarrow{\text{conv}} \eta) \end{aligned}$$

Calculating known representations and conventions, and the levity of a known convention:

$$\frac{\Gamma \vdash \tau : \text{TYPE } \pi \quad \pi \in \text{PrimRep}}{\Gamma \vdash \tau \xrightarrow{\text{rep}} \pi} \quad \frac{\Gamma \vdash \tau : \text{TYPE } \rho \quad \eta \in \text{KnownConv}}{\Gamma \vdash \tau \xrightarrow{\text{conv}} \eta} \quad \begin{aligned} \text{lev}(\text{Eval}^\psi) &= \psi \\ \text{lev}(\text{Call}[\bar{\pi}]) &= \text{U} \end{aligned}$$

Fig. 8. Compiling  $\mathcal{IL}$  to  $\mathcal{ML}$

then we can evaluate the result of  $e$  directly by the main compilation translation, otherwise if it is  $\text{Call}[\bar{\pi}]$  then  $e$  must be called (not evaluated). To make sure that the definition and call sites of a primitive function match, we always fully  $\eta$ -expand these expressions when they are defined: either on the right-hand side of **lets** or in the body of **Closures**. Because of  $\eta$ -expansion, this step of compilation is only defined when the calling convention is statically known (e.g., it is not a variable  $n$  or partially-defined like  $\text{Call}[r_1, r_2, \text{arity}(n)]$ ). In any case, we next move to the main work-horse of compilation,  $C \llbracket e \rrbracket_\theta^\Gamma(\bar{a})$ , that produces  $\mathcal{ML}$  code to evaluate the result of  $e$  applied to the arguments  $(\bar{a})$ . Again, there are invariants to this translation that we will enumerate shortly.

For  $C \llbracket e \rrbracket_\theta^\Gamma(\bar{a})$ , literal constants are just passed through, but compiling a call to **error** assumes precisely one argument. What if the user has written a partial application of **error**? Such partial

applications are *always*  $\eta$ -expanded to be fully saturated, satisfying the requirement here. Compiling a variable  $x$  looks it up in the environment  $\theta$ , but there is different  $\mathcal{ML}$  code for *evaluating*  $x$  versus *calling* it, even when there are no arguments. In other words, a primitive function variable  $x$  with the empty calling convention  $\text{Call}[]$  compiles as  $C[x]_{[y_{\text{PtrR}}/x]\theta}^\Gamma(\varepsilon) = y_{\text{PtrR}}()$ —a nullary function call—but a variable  $x$  with the convention  $\text{Eval}^\gamma$  compiles as to the pointer lookup  $C[x]_{[y_{\text{PtrR}}/x]\theta}^\Gamma(\varepsilon) = y_{\text{PtrR}}$ .

Closure applications are compiled straightforwardly, and erasable arguments  $\phi$  and binders  $\Lambda\chi.e$  are simply dropped.  $\mathcal{IL}$  answers can be compiled outright to a  $\mathcal{ML}$  WHNF via  $\mathcal{A}[A]_\theta^\Gamma$ . Compiling an application to an answer  $A$  depends on the nature of that answer:

- If  $\mathcal{A}[A]_\theta^\Gamma$  is a variable or constant (after type erasure), then it can be passed directly.
- Otherwise, name the argument with a **let** (respecting the  $\Lambda$ -normal form) and pass it by reference. In this case, the compiled argument will always have the form  $\text{I\#}(a)$  or  $\text{Clos}^n W$ , meaning the **let**-binding will always be represented as a pointer into the heap.<sup>15</sup>

In the variable case, we do not need to track the levity or representation of the argument, because these decisions have already been made by the context, when the variable definition itself was compiled. Crucially, we did not have to look up any information in the typing environment to compile answer arguments; this is why no highlighted premises are needed in rule  $\text{FUN-A-E}$ .

In the case of an application  $e\ e'$  to an arbitrary argument that needs to be computed, corresponding to  $\text{FUN-E}$ , we always generate a **let** similar to the second case for  $e\ A$ . However, for  $\text{FUN-E}$ , we need to determine the representation, convention, *and* levity of the binding, which could truly be anything. This corresponds to the highlighted side conditions in  $\text{FUN-E}$ .

### 6.3 Correctness of Compilation

Notice how the same polymorphism restrictions used in the typing rules also appear during compilation. Even though the defined compilation translation is partial (not every syntactically valid expression can be compiled), all well-typed  $\mathcal{IL}$  expressions with a known convention have a defined compilation to  $\mathcal{ML}$ . In particular,  $\mathcal{E}_\eta[e]$  is well-defined for any closed expression  $\vdash e : \tau : \text{TYPE } \rho\ \eta$ , where the syntax of known conventions  $\eta$  is given in Fig. 7.

In fact, we allow for a little more levity polymorphism during compilation:  $\mathcal{E}_{\text{Eval}^g}[e]$ , for a polymorphic levity  $g$ , is also allowed. That's because the generated code will be executed exactly when expression is evaluated: in other words, when a computation is forced, there is no difference between eager (U) or lazy (L). This added flexibility is essential for compiling levity polymorphic expressions appearing in strict contexts, such as in the discriminant of a **case** or first argument of **App**. Although implicit, the  $C$  translation assumes that evaluation of the compiled expression is being forced. In contrast, the  $\mathcal{A}$  translation does not assume this, because it is used in contexts that do *not* force the expression. This small difference is how we are able to pass variables of *any* convention (eager, lazy, or primitive functions) without erroneously introducing extra strictness.

During compilation, we occasionally need to know representation ( $\pi$ ) or convention ( $\eta$ ) of a sub-expression. This appears in Fig. 8 as highlighted side conditions  $\Gamma \vdash \tau \overset{\text{rep}}{\rightsquigarrow} \pi$  and  $\Gamma \vdash \tau \overset{\text{conv}}{\rightsquigarrow} \eta$ , respectively. In general, compilation could fail if the representation or convention in the kind of  $\tau$  are partially unknown—that is, contains free variables. But any closed representation has the form  $\pi$  and any closed convention is equivalent to a  $\eta$ , as defined by the syntax of  $\mathcal{ML}$  in Fig. 7. The places where this requirement appears corresponds exactly to the highlighted monomorphism restrictions in Fig. 3.

**Theorem 3** (Closed Compilation). *If  $\vdash e : \tau$  and  $\vdash \tau : \text{TYPE } \rho\ \nu$  then  $\mathcal{E}_\nu[e]$  is defined.*

<sup>15</sup>A more feature-rich language may allow for representations other than just a single pointer, in this case. Even then, answers compile to values with syntactically manifest representations, so no additional typing information is needed here.

This theorem just states when compilation succeeds in generating  $\mathcal{ML}$  code from a closed  $\mathcal{IL}$  expression. We should also expect that compilation preserves the behavior of that  $\mathcal{IL}$  expression as well. In other words, if an expression is equal to some answer in  $\mathcal{IL}$  (as per Fig. 4), then executing the compiled code should give the same answer in  $\mathcal{ML}$ . But we are not interested in evaluating primitive functions directly—they are called, not evaluated!—so answers will be of some Evaluatable types like (un)boxed integers, which are simple enough values to line up on the nose.

**Theorem 4** (Soundness and Completeness).

- (1) For any  $\vdash e : \text{Int}\#$ ,  $\vdash e = i : \text{Int}\#$  if and only if  $\langle \mathcal{E}_{\text{Eval}^u} \llbracket e \rrbracket \mid \varepsilon \mid \varepsilon \rangle \mapsto^* \langle i \mid \varepsilon \mid H \rangle$ .
- (2) For any  $\vdash e : \text{Int}^Y$ ,  $\vdash e = \text{I}\#^Y i : \text{Int}^Y$  if and only if  $\langle \mathcal{E}_{\text{Eval}^Y} \llbracket e \rrbracket \mid \varepsilon \mid \varepsilon \rangle \mapsto^* \langle \text{I}\#(i) \mid \varepsilon \mid H \rangle$ .

## 7 DYNAMIC ARITY

Consider the following program written in Haskell, where *exp* is some expensive function:<sup>16</sup>

<pre>data T = MkT (Int → Int → Int) t1 = MkT (λx. let z = exp x in λy. z + y) t2 = MkT plus</pre>	<pre>appT1 : T → Int → Int appT1 (MkT g) = g 1  appT2 : T → Int appT2 (MkT g) = g 1 2</pre>
---	---

In terms of the informal notion of arity in the source language (Section 2), we can say that *t1* stores an arity 1 closure and *t2* stores an arity 2 closure. Likewise, *appT1* performs an arity 1 application and *appT2* performs an arity 2 application. This can be compiled to  $\mathcal{IL}$  (extended with data type declarations) similar to the translation in Fig. 5, which formalizes the arities like so:

<pre>data T = MkT<sup>L</sup> {Int ∼<sup>L</sup> {Int ∼<sup>L</sup> Int}} t1 = MkT (Clos<sup>L</sup> λx. let z = exp x in Clos<sup>L</sup> λy. z + y) t2 = MkT (Clos<sup>L</sup> λx. Clos<sup>L</sup> λy. plus x y)</pre>	<pre>appT1 : T ∼<sup>L</sup> {Int ∼<sup>L</sup> Int} appT1 (MkT g) = App g 1  appT2 : T ∼<sup>L</sup> Int appT2 (MkT g) = App (App g 1) 2</pre>
---	---

Notice how *appT2* (*MkTg*) does not apply *g* to both arguments at once; instead, it must evaluate *g* applied to 1 first, then apply the returned closure to 2 in a separate step. The two-step application process is *mandated* by the type of the *MkT* constructor, even though most of the time *MkT* will ultimately be used store a closure capable of accepting both arguments at once. This arity demotion can be seen in *t1*, where the binary *plus* function is wrapped up in a chain of two unary closures, as required by *MkT*. As a result, the call *appT1 t2* must pass 1 and 2 separately to *plus*, losing the opportunity for the faster binary calling convention that seemed possible in the source program.

We can attempt improve the performance of *t2* and *appT2* with an alternate translation:

<pre>data T' = MkT' {Int ∼<sup>L</sup> Int ∼<sup>L</sup> Int} t1' = MkT' (Clos<sup>L</sup> λx. λy. let z = exp x in z + y) t2' = MkT' (Clos<sup>L</sup> plus)</pre>	<pre>appT1' : T' ∼<sup>L</sup> Int ∼<sup>L</sup> Int appT1' (MkT' g) x = App g 1 x  appT2' : T' ∼<sup>L</sup> Int appT2' (MkT' g) = App g 1 2</pre>
---	---

Now, *MkT'* contains a closure of a single binary function. This way, *appT2' t2'* steps to the single binary application *plus* 1 2, passing both arguments at once to *plus*. But the type of *MkT'* makes it impossible to contain closures that memoize work when applied to only one argument—an unfortunate, unintended ramification of this “optimization.” Unlike before, the closure inside *t1'* recomputes *exp x* every time the second argument is given, instead of memoizing the result once the first argument is provided. Here, we make this fact syntactically explicit by  $\eta$ -expanding the definitions of *t1'* and *appT1'*, but because  $\sim$ -types are fully extensional, the same recomputation

<sup>16</sup>The exact same example applies to OCaml as well by replacing  $\text{L}$  with  $\text{U}$  in translations into  $\mathcal{IL}$  that follow.

will occur no matter what. For example,  $\text{map } (\text{appT1}' \ t1') [1..1000]$  recomputes  $\text{exp } 1$  for all 1000 elements of the list, whereas the previous  $\text{map } (\text{appT1 } t1) [1..1000]$  only computes  $\text{exp } 1$  once.

The root of the problem is that the data type definition of  $T$  forces us to pick *one* type for the constructor  $\text{MkT}$ . Because  $\mathcal{IL}$  statically links this type to an arity, we are thus forced to pick *one* arity for  $\text{MkT}$  closures that is used throughout the entire program.  $\text{MkT}$  may often be used to hold essentially binary function closures, but the computational effect of unary application of  $\text{MkT}$  closures can still be crucial in certain corners of a program—either for asymptotic complexity in Haskell or for side effects in OCaml.

Rather than enforcing *statically* that every call is exactly saturated, perhaps we could perform a *dynamic* check. Given a binary application, we could interrogate the function value to check its arity, and behave differently depending on whether that arity is 1 (an over-saturated call), 2 (exactly saturated, the fast case), or greater than 2 (unsaturated). This is, in fact, what GHC does today, and corresponds to the “unknown” function calls of Marlow and Peyton Jones [2004] which inspect the arities of closure values at runtime to choose the best calling convention. But how can  $\mathcal{IL}$ —with its statically-tracked notion of arity—accommodate dynamic arity checking?

The key is to allow for the primitive functions contained in closures to have a different arity than their call site, thus requiring a dynamic check on all Applications of closures. Applying too few arguments creates a partial application, and applying too many is broken down into several steps. More formally, the syntax of  $\mathcal{ML}$  can be extended with partial applications  $\text{Clos}^n f(\bar{a})$ , where  $f$  has been applied to arguments  $\bar{a}$  so far and  $n$  is the number of remaining arguments expected before  $f$  can be called. Now, consider these extra rules for dynamic handling a runtime arity mismatch:

$$\begin{aligned}
 (\text{Apply}) \quad & \langle \text{Clos}^n W(\bar{a}) \mid \text{App}(\bar{a}'); K \mid H \rangle \mapsto \langle W(\bar{a}, \bar{a}') \mid K \mid H \rangle & (\text{if } |\bar{a}'| = n) \\
 (\text{PApp}) \quad & \langle \text{Clos}^n W(\bar{a}) \mid \text{App}(\bar{a}'); K \mid H \rangle \mapsto \langle \text{Clos}^{n-|\bar{a}'|} W(\bar{a}, \bar{a}') \mid K \mid H \rangle & (\text{if } |\bar{a}'| < n) \\
 (\text{OApp}) \quad & \langle \text{Clos}^n W(\bar{a}) \mid \text{App}(\bar{a}', \bar{a}''); K \mid H \rangle \mapsto \langle W(\bar{a}, \bar{a}') \mid \text{App}(\bar{a}''); K \mid H \rangle & (\text{if } |\bar{a}'| = n, |\bar{a}''| > 0)
 \end{aligned}$$

In practice, these extra rules for partial- and over-application lets us treat closure types like  ${}^L\{\text{Int} \rightsquigarrow {}^L\{\text{Int} \rightsquigarrow \text{Int}\}\}$  and  ${}^L\{\text{Int} \rightsquigarrow \text{Int} \rightsquigarrow \text{Int}\}$  as the same, without endangering type safety due to arity mismatch. For example, this would eliminate the difference between the  $T$  and  $T'$  data types from before. With these dynamic arity checks, it is safe to call  $\text{appT1 } t2'$ : this results in a partial application because the caller ( $\text{appT1}$ ) only provides one argument to the binary closure ( $t2'$ ). Likewise, it is safe to call  $\text{appT2}' \ t1$ : this results in an over application where the caller ( $\text{appT2}'$ ) wants to pass two arguments to a unary closure ( $t1$ ). Yet, in the cases where the optimal arities do match (like  $\text{appT1 } t1$  and  $\text{appT2}' \ t2'$ ), the fastest calling convention is used at runtime. Therefore, the types  $T$  and  $T'$  can be used interchangeably, in some sense, as long as arities are checked at runtime. More generally, dynamic arity checks lets us safely equate these two closure types:

$${}^V\{\bar{\tau} \rightsquigarrow {}^V\{\sigma\}\} = {}^V\{\bar{\tau} \rightsquigarrow \sigma\}$$

This could be formalized in  $\mathcal{IL}$  as a type equality or a type-safe coercion [Breitner et al. 2016]: both types are represented identically at runtime as closure objects with some runtime arity count.

In the end, both known and unknown calls of Marlow and Peyton Jones [2004] can be captured in the intermediate language. The arity of a type  $\tau \rightsquigarrow \sigma$  is known statically by its kind, and the program must provide the right number of arguments and binders. However, types like  ${}^V\{\tau \rightsquigarrow \sigma\}$  and  ${}^V\{\tau \rightsquigarrow {}^V\{\sigma\}\}$ , which both store additional arity information at runtime, can be freely interchanged at compile time, as long as the arities are checked at runtime. To be clear, this extension has a trade-off: the closures described here are subject to extra dynamic checks. It is possible that an implementation would want to have both statically checked closures and dynamically checked ones. We can accommodate both by simply having two different closure types (with their own

Clos and App). Then, an optimizing compiler, or an expert user, can select the one with the best performance for a particular part of a program.

## 8 RELATED WORK

The system presented in this paper is the culmination of several independent lines of work on expressing performance issues directly in an intermediate language. The underlying theme is to capture the low-level details of calling conventions as features of a higher-level functional language.

### 8.1 Representation and Levity in the Kinds

The idea of distinguishing (un)boxed and (un)lifted types goes back to Peyton Jones and Launchbury [1991]. That distinction has been static until recent work added levity polymorphism to the mix [Eisenberg and Peyton Jones 2017], and shown that its utility is greater than expected (see Section 7 of that work). However, Eisenberg and Peyton Jones [2017] conflates *levity polymorphism* and *representation polymorphism*. Our contribution separates the two completely, with applications that are polymorphic in one but not the other. One of our main requirements is to generate only *one* piece of code for every polymorphic definition. Certain definitions that must be rejected, because compilation would depend on a choice made at runtime. An alternative approach by Dunfield [2015] accepts more uses of levity polymorphism at the cost of generating *different* code for each choice—an exponential blowup of code size in practice—which we avoid.

### 8.2 Optimizing Curried Functions

Previous work established methods for optimizing curried function calls dynamically at runtime, avoiding the overhead of naively calling  $((f\ 1)\ 2)\ 3$  by passing one argument at a time. In practice,  $f$  will often expect all three arguments before doing any interesting work, so those calls should be fused when possible. Fusing can be done by pushing many arguments on the stack at once (the push/enter model) [Krivine 2007; Leroy 1990] or by evaluating the arity of closures (the eval/apply model) [Marlow and Peyton Jones 2004]. In this work, we capture this dynamic type of optimization within the syntax and types of programs, as described in Section 7.

### 8.3 Function Arity in Types

While there is performance to be gained by dynamically optimizing curried function calls at runtime, it is even better to optimize statically at compile time. Of course, this is easy to do when the compiler can find the definition of the called function [Marlow and Peyton Jones 2004]. This scheme is easily thwarted by higher-order functions, so a less syntactic approach—like one based on types—can be beneficial. Uncurrying—representing a function  $a \rightarrow b \rightarrow c \rightarrow d$  as  $(a, b, c) \rightarrow d$ —is an obvious place to start, and has been investigated before [Bolingbroke and Peyton Jones 2009; Dargaye and Leroy 2009; Hannan and Hicks 1998]. However, when polymorphism is brought into the picture, type quantification is irreparably fused with multi-arity functions; see [Downen et al. 2019, Section 8.1].

Following Downen et al. [2019],  $\mathcal{IL}$  instead retains the curried form of function types. However,  $\mathcal{IL}$  goes significantly beyond that work by supporting type polymorphism over arrow types (Section 4.3), and polymorphism over levities (Section 4.4) and conventions (Section 4.5). Another difference is that Downen et al. [2019] had two function arrows,  $(\tau \xrightarrow{\gamma} \sigma)$  and  $(\tau \rightsquigarrow \sigma)$ , whereas  $\mathcal{IL}$  has just one arrow  $(\tau \rightsquigarrow \sigma)$ , plus the closure type  ${}^\gamma\{\tau\}$ . The two are inter-convertible: we showed how to translate  $\tau \xrightarrow{\gamma} \sigma$  to  $\mathcal{IL}$  for either  $\gamma$  in Fig. 5, and in the other direction we have  ${}^\gamma\{\tau\} = () \xrightarrow{\gamma} \tau$  with  $(\text{Clos}^\gamma e) = \lambda^\gamma().e$  and  $(\text{App } e) = e ()$ . Note that, to make the analogy operationally exact, the unit type  $()$  should be an unboxed, empty tuple (*i.e.*, represented as 0 arguments at runtime). The approach here has a greater economy of concepts, and a nice correspondence with  $\text{Int}\#$  and  $\text{Int}^\gamma$ . However, two function arrows might be better for a practical compiler.



## 8.4 The Glasgow Haskell Compiler

GHC already implements a rich kind system, including polymorphism over types, kinds, and representations. Indeed GHC goes further: instead of a stratified zoo of different things (types, kinds, representations, *etc.*) as in Fig. 1, they are all *types* [Weirich et al. 2013] kept separate by their *kinds*. This is a fantastic simplification, immediately allowing polymorphism over all these conceptually-different things. This does, however, make it hard *not* to have polymorphism! Returning to Section 4.3, it would be hard to prevent instantiation of a forall-quantified type variable with an arrow type, requiring a restriction like “this quantified variable can have any kind *other than* Call.” So GHC’s infrastructure strongly encourages fully-fledged polymorphism.

## 8.5 Logical Foundations

The  $\mathcal{IL}$  language is not an ad-hoc collection of design compromises driven by only performance considerations. Rather, it grows directly from principled foundations in logic.

Previous work on unboxed types and extensional functions shares the observation that *lifting*—in the sense of denotational semantics—corresponds to a *mismatch* between machine primitives and the semantics of a programming language. Unlifted types can be implemented more directly—and therefore more efficiently—in a machine. But the cause of lifting depends on the type: unlifted integers need to be call-by-value whereas unlifted curried functions need to be call-by-name. The first reconciliation was achieved in call-by-push-value [Levy 2001], which avoids all lifting unless explicitly requested. As such, this paper can be seen as a practical extension of this foundation.

The same connection between types and evaluation is also tied to *focusing and polarity* [Andreoli 1992; Laurent 2002] in proof search, which corresponds to pattern-matching in functional programming [Zeilberger 2008, 2009] and semantics and computation [Munch-Maccagnoni 2009, 2013]. Recently, these mixed evaluation strategy languages have been extended with practical features like call-by-need evaluation [Downen and Ariola 2018; McDermott and Mycroft 2019] to model shared computation. Of note, the types in  $\mathcal{IL}$  used for boxing and indirection correspond exactly to the “polarity shifts” of Downen and Ariola [2018] to and from call-by-need. In particular, the boxed integer type corresponds to an “up shift” ( $\text{Int} = \uparrow\text{Int}\#$ ) and the function closures to a “down shift” ( $\{\tau \rightsquigarrow \sigma\} = \downarrow(\tau \rightsquigarrow \sigma)$ ). For the sake of usability,  $\mathcal{IL}$  performs other implicit polarity conversions of types based on their context. For example, closing over a non-function type like  $\text{Int}\#$  implicitly shifts it to a “nullary function” (there written  $\uparrow\text{Int}$ ), expressed by the encoding  $\{\text{Int}\# \} = \downarrow\uparrow\text{Int}\#$ .

## 9 CONCLUSION

This paper illustrates a cohesive system for including low-level details—specifically *representation*, *levity*, and *arity*—inside a higher-level intermediate language. Not only does this let the language express intensional properties of programs, it also lets programs *abstract* over these details when they do not impact compilation. Parts of this work have been implemented already in the Glasgow Haskell Compiler, and we intend to further implement the entirety of kinds as calling conventions. The story presented here takes an explicitly typed intermediate language and—through type-driven elaboration—compiles it to an untyped target language. As future work, it could be enlightening to consider how types might be preserved by compilation by giving a sufficiently expressive type system for the lower-level language. Since the main objective is to capture performance in the intermediate language, we would also like to characterize the cost of computation in its semantics.

## ACKNOWLEDGMENTS

The authors would like to thank Norman Ramsey and the anonymous reviewers, whose feedback was invaluable in improving the presentation of this paper. The material is based upon work supported by the National Science Foundation under Grant No. 1719158.

## REFERENCES

- Jean-Marc Andreoli. 1992. Logic Programming with Focusing Proofs in Linear Logic. *Journal of Logic and Computation* 2, 3 (1992), 297–347. <https://doi.org/10.1093/logcom/2.3.297>
- Zena M. Ariola and Matthias Felleisen. 1997. The Call-By-Need Lambda Calculus. *Journal of Functional Programming* 7, 3 (May 1997), 265–301. <https://doi.org/10.1017/S0956796897002724>
- Maximilian C. Bolingbroke and Simon L. Peyton Jones. 2009. Types Are Calling Conventions. In *Proceedings of the 2nd ACM SIGPLAN Symposium on Haskell (Haskell '09)*. ACM, 1–12.
- Joachim Breitner. 2014. Call Arity. In *Trends in Functional Programming - 15th International Symposium, TFP 2014, Soesterberg, The Netherlands, May 26-28, 2014. Revised Selected Papers*. 34–50.
- Joachim Breitner, Richard A. Eisenberg, Simon Peyton Jones, and Stephanie Weirich. 2016. Safe zero-cost coercions for Haskell. *Journal of Functional Programming* 26 (2016), e15. <https://doi.org/10.1017/S0956796816000150>
- Zaynah Dargaye and Xavier Leroy. 2009. A verified framework for higher-order uncurrying optimizations. *Higher-Order and Symbolic Computation* 22, 3 (2009), 199–231.
- Paul Downen and Zena M. Ariola. 2018. Beyond Polarity: Towards a Multi-Discipline Intermediate Language with Sharing. In *27th EACSL Annual Conference on Computer Science Logic, CSL 2018, September 4-7, 2018, Birmingham, UK*. 21:1–21:23.
- Paul Downen, Zachary Sullivan, Zena M. Ariola, and Simon Peyton Jones. 2019. Making a Faster Curry with Extensional Types. In *Proceedings of the 12th ACM SIGPLAN International Symposium on Haskell (Haskell 2019)*. Association for Computing Machinery, New York, NY, USA, 58–70. <https://doi.org/10.1145/3331545.3342594>
- Joshua Dunfield. 2015. Elaborating evaluation-order polymorphism. In *Proceedings of the 20th ACM SIGPLAN International Conference on Functional Programming, ICFP 2015, Vancouver, BC, Canada, September 1-3, 2015*. 256–268.
- Richard Eisenberg. 2019. GHC Proposal 29: revised levity polymorphism. <https://github.com/ghc-proposals/ghc-proposals/blob/master/proposals/0029-levity-polymorphism.rst>
- Richard A. Eisenberg and Simon Peyton Jones. 2017. Levity polymorphism. In *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2017, Barcelona, Spain, June 18-23, 2017*. 525–539.
- Sebastian Graf. 2020. GHC Proposal 265: unlifted data types. <https://github.com/ghc-proposals/ghc-proposals/blob/master/proposals/0265-unlifted-datatypes.rst>
- John Hannan and Patrick Hicks. 1998. Higher-Order unCurrying. In *POPL '98, Proceedings of the 25th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, San Diego, CA, USA, January 19-21, 1998*. 1–11.
- Jean-Louis Krivine. 2007. A Call-By-Name Lambda-Calculus Machine. *Higher-Order and Symbolic Computation* 20, 3 (2007), 199–207.
- Olivier Laurent. 2002. *Étude de la polarisation en logique*. Ph.D. Dissertation. Université de la Méditerranée - Aix-Marseille II.
- Xavier Leroy. 1990. *The ZINC experiment: an economical implementation of the ML language*. Technical report 117. INRIA.
- Paul Blain Levy. 2001. *Call-By-Push-Value*. Ph.D. Dissertation. Queen Mary and Westfield College, University of London.
- Simon Marlow and Simon L. Peyton Jones. 2004. Making a fast curry: push/enter vs. eval/apply for higher-order languages. In *Proceedings of the Ninth ACM SIGPLAN International Conference on Functional Programming, ICFP 2004, Snow Bird, UT, USA, September 19-21, 2004*. ACM, 4–15.
- Andrew Martin. 2019a. GHC Proposal 112: unlifted arrays. <https://github.com/ghc-proposals/ghc-proposals/blob/master/proposals/0112-unlifted-array.rst>
- Andrew Martin. 2019b. GHC Proposal 203: pointer rep. <https://github.com/ghc-proposals/ghc-proposals/blob/master/proposals/0203-pointer-rep.rst>
- Andrew Martin. 2019c. GHC Proposal 98: unlifted newtypes. <https://github.com/ghc-proposals/ghc-proposals/blob/master/proposals/0098-unlifted-newtypes.rst>
- Dylan McDermott and Alan Mycroft. 2019. Extended Call-by-Push-Value: Reasoning About Effectful Programs and Evaluation Order. In *Programming Languages and Systems*, Luís Caires (Ed.). Springer International Publishing, Cham, 235–262.
- Guillaume Munch-Maccagnoni. 2009. Focalisation and Classical Realisability. In *Computer Science Logic: 23rd international Workshop, CSL 2009, 18th Annual Conference of the EACSL (CSL 2009)*, Erich Grädel and Reinhard Kahle (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 409–423.
- Guillaume Munch-Maccagnoni. 2013. *Syntax and Models of a non-Associative Composition of Programs and Proofs*. Ph.D. Dissertation. Université Paris Diderot.
- Simon L. Peyton Jones. 1992. Implementing Lazy Functional Languages on Stock Hardware: The Spineless Tagless G-machine. *Journal of Functional Programming* 2, 2 (1992), 127–202.
- Simon L. Peyton Jones and John Launchbury. 1991. Unboxed Values As First Class Citizens in a Non-Strict Functional Language. In *Proceedings of the 5th ACM Conference on Functional Programming Languages and Computer Architecture*. Springer-Verlag, London, UK, UK, 636–666.

- Amr Sabry and Matthias Felleisen. 1993. Reasoning About Programs in Continuation-Passing Style. *Lisp and Symbolic Computation* 6, 3-4 (Nov. 1993), 289–360.
- Amr Sabry and Philip Wadler. 1997. A Reflection on Call-by-Value. *ACM Transactions on Programming Languages and Systems* 19, 6 (1997), 916–941.
- Alex Theriault. 2019. GHC Proposal 209: levity-polymorphic Lift. <https://github.com/ghc-proposals/ghc-proposals/blob/master/proposals/0209-levity-polymorphic-lift.rst>
- Philip Wadler, Walid Taha, and David MacQueen. 1998. How to add laziness to a strict language without even being odd. In *Proceedings of the Standard ML Workshop*.
- Stephanie Weirich, Justin Hsu, and Richard A. Eisenberg. 2013. System FC with Explicit Kind Equality. In *International Conference on Functional Programming (ICFP '13)*. ACM.
- Noam Zeilberger. 2008. On the Unity of Duality. *Annals of Pure and Applied Logic* 153, 1 (2008), 660–96.
- Noam Zeilberger. 2009. *The Logical Basis of Evaluation Order and Pattern-Matching*. Ph.D. Dissertation. Carnegie Mellon University.