

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 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 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.

Draft March 2020. We would appreciate any comments or feedback to the email addresses below.

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

1 INTRODUCTION

Polymorphism supports re-use by allowing one chunk of executable to work on 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 re-use if we could specialize to, say, “any type that is 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*:

$$twice\ f\ x = f\ (f\ x)$$

We would like to be able to control 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 the type system.

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 (Section 2):

- *Representation*. How is this argument represented at runtime?
- *Levity*. What is the evaluation strategy of this argument (call-by-value or call-by-need)?
- *Calling convention*. For functions, how many arguments does this function take before its code can be executed (its arity)?

Indeed, it turns out that 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 simple, uniform language, but the compiler needs a more expressive intermediate language so that it can express low-level representation choices, and expose that code to the optimizer. For example, the programmer might work exclusively with boxed integer values, of type `Int` say, but the intermediate language

Authors’ addresses: Paul Downen; Zena M. Ariola, University of Oregon, 1477 E. 13th Ave. Eugene, Oregon, 97403, USA, pdownen@cs.uoregon.edu, ariola@cs.uoregon.edu; Simon Peyton Jones, Microsoft Research, 21 Station Rd. Cambridge, CB1 2FB, UK, simonpj@microsoft.com; Richard A. Eisenberg, Bryn Mawr College, 101 N. Merion Ave. Bryn Mawr, PA, 19010, USA, Tweag I/O, Cambridge, UK, rae@richarde.dev.

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].

In this paper we build directly on several earlier works that track different representations [Eisenberg and Peyton Jones 2017] and function arities [Downen et al. 2019] within a type and kind system, but here we bring them together into a single framework, more powerful and more precise than any of its predecessors. Specifically, we make these contributions:

- We introduce a polymorphic intermediate language that statically captures *calling conventions in kinds* (Section 3), and has polymorphism over the *representation*, *levity*, and *arity* of types.
- We show how to compile our polymorphic intermediate language to a more conventional lower-level representation (Section 4) that has different representations of values (e.g., pointers versus integers) and multi-arity functions, but is *not* polymorphic. Crucially, compilation is driven by kinds and uses static restrictions to ensure that polymorphic code can always be compiled to monomorphic code, without sacrificing type erasure or duplicating code.
- Even though we statically track calling conventions in kinds, there might still be opportunities for improving a call at runtime. We describe a small extension to our intermediate language to allow for dynamic checks on the arities of closures at runtime, so that we can use the best arity available during execution (Section 5).
- 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 6).
- We discuss how having kinds as calling conventions can be used to generalize user-defined data types (Section 7) for expressing more efficient code in a high-level language and enriching a source language with both lazy and eager data types.
- We provide evidence of correctness for the full compilation process (Theorems 2 and 4) for call-by-name and call-by-value evaluation, thus showing that our intermediate language is equally well-suited for representing *both eager and lazy* functional languages.

We discuss related work and conclude in Sections 8 and 9.

2 SYSTEM F IS IMPRACTICAL AS AN INTERMEDIATE LANGUAGE

When developing an intermediate language for functional programs, System F serves as a wonderful starting point that solves the challenge of polymorphism. However, it does not address several concerns involving the efficiency of those programs. Imagine a baseline of System F: everything is represented uniformly by a pointer to support polymorphism; function arguments are passed one-at-a-time to support currying; and the evaluation strategy is fixed as call-by-value or call-by-name (the particular choice is not important) to support either laziness or strictness, but not both. Before going into our main contribution—a type and kind system that enables efficient calling conventions—in the next section, let us unpack these individual issues in more detail.

2.1 Polymorphism

Many optimizing compilers use *type erasure* to handle polymorphic functions, wherein typing information is not present at runtime. Consider the lowly identity function $id : t \rightarrow t$. Its type tells us that id works over *any* type t . Yet even a function as simple as id must do *some* work. It must know where to find its argument, where to put its result, and then copy from the first spot to the second. Type erasure imposes the constraint that the code generated for id must be able to do all of this for *any* type t : type erasure means that we cannot inspect the particular choice of t , so we must have one sequence of instructions that works for any t .

Many types are represented by pointers, suggesting that we can compile *id* to expect one pointer as input and produce one pointer as output. Unfortunately, this simplification prevents specializing *id* to work with other types of runtime values like `Int#`—the type describing machine integers—which may have a different size or use different registers than pointers.

We thus arrive at a dilemma: either somehow restrict *id* not to be able to work with `Int#`, or ban types like `Int#` from our language altogether. We say that the latter is impractical, because constant, unavoidable indirection via pointers is too slow for industrial programming. Effectively, a naïve implementation of System F makes this latter, impractical choice. Instead, a real-world programming language must somehow restrict *id*. We do so by classifying `Int#` as having a different *kind* than types represented by pointers [Eisenberg and Peyton Jones 2017].

This second, more practical, choice is similar to the move to System F_ω , an extension of System F with many different kinds of types for expressing type operators like `Array`. In System F_ω , we are forced to state the kind of types a type variable like *t* ranges over. For example, the identity function is elaborated as $\forall t:\star. t \rightarrow t$, where the kind \star must somehow specify that *t* is represented as a pointer. A type like `Int#`, then, must have a kind which is different from \star .

2.2 Higher-order functions

Functional programming languages get significant mileage out of *currying*: a function *plus* of type `Int → Int → Int` is regarded as a function of one argument (of type `Int`) that returns a function of one argument (another `Int`), that finally returns an `Int`. Such a function can be *partially applied*. For example *map* (*plus* 1) is a function that increments every element in a list; here both *map* and *plus* are partially applied. Currying is very convenient for programmers, and can lead to measurable reduction in code size [Arvind and Ekanadham 1988].

But we do not want to *implement* currying by making every multi-argument function return a chain of (heap-allocated) function closures; that would be unacceptably inefficient. So compilers like GHC and OCamlC go to great lengths to group together adjacent lambdas, and pass the arguments “all at once” to an uncurried variant of the function [Leroy 1990; Marlow and Peyton Jones 2004]. The number of arguments simultaneously expected by a function before doing work is called its *arity*. When calling a (let-bound) function whose definition is in scope, GHC can ensure that the right number of arguments are passed to match its arity and thus producing efficient code.

But calling unknown (lambda-bound) functions whose definitions cannot be found is harder. Consider the third clause of *zipWith*:

$$\begin{aligned} \text{zipWith } f \ [] \ ys &= [] \\ \text{zipWith } f \ xs \ [] &= [] \\ \text{zipWith } f \ (x : xs') \ (y : ys') &= (f \ x \ y) : (\text{zipWith } xs' \ ys') \end{aligned}$$

The arity of the function *f* passed to *zipWith* might be 1 or 2. It might even be 3 or 0 (if it is an unevaluated thunk), depending on the caller. So the code for *zipWith* is forced to account for all these possibilities in the call $(f \ x \ y)$. But we could generate *much* better code if we statically knew *f*’s arity and, building on earlier work [Downen et al. 2019], that is what we do here.

2.3 Evaluation strategy

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 *later, on-demand*? The designer of a particular version of System F makes this choice in its semantics. Real-world programming languages also have committed to a choice here. This choice is frequently 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, for example, to make a control-flow

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

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

decision, and only once. However, regardless of a language’s choice of evaluation strategy, some scenarios demand making the opposite choice. Programmers in an eager language sometimes use constructs like *Delay* and *Force* [Wadler et al. 1998] to embed a lazy evaluation, and programmers in a lazy language use strict bindings or primitives (such as Haskell’s *seq*) to force eager evaluation. We thus want a language that supports both lazy and eager evaluation, giving the programmer a convenient means to choose between the two.

With the terminology of “liftedness”—whether a variable of a type can be bound to a computation, or only to values—the types of an eager language (like OCaml) are all unlifted, whereas the types of a lazy language (like Haskell) are all lifted. We can thus call the change between lazy and eager computation a change of liftedness, or *levity*. A function that can work with both lifted and unlifted types is *levity polymorphic*.¹ In order to support levity polymorphism in our language, we once again leverage the kind system, distinguishing between the two modes in the kinds of types.

3 OUR KEY CONTRIBUTION: THE INTERMEDIATE LANGUAGE (\mathcal{IL})

Previous research has suggested that *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 that combines several different previous works that focused on addressing representations [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.

Why *kinds*, rather than *types*? Because of polymorphism. A type can tell us important runtime details (such as the representations of values or the arity of functions), but polymorphism means that types might be statically unknown inside of a definition. Yet we still want to compile polymorphic definitions, and in such a way that the same code can be reused for every instantiation. Therefore, we encode just enough intensional information within kinds to express the low-level details needed to compile polymorphic code. Think of these kinds as a coarser-grained descriptions, that leave more subtle issues like safety to the finer-grain world of types. Better still, making details like

¹Previous work [Eisenberg and Peyton Jones 2017] is titled *Levity Polymorphism*. Yet the paper does not, in our opinion, deliver exactly on the promise in the title. Instead, it describes *representation polymorphism*, choosing not to distinguish between levity and runtime representation. As a consequence, levity polymorphism is only possible (and in fact, mandatory) as part of representation polymorphism.

$$\begin{array}{c}
\text{Types of expressions } \boxed{\Gamma \vdash e : \tau} \\
\\
\frac{}{\Gamma, x:\tau \vdash x : \tau} \text{VAR} \quad \frac{c : \tau}{\Gamma \vdash c : \tau} \\
\\
\frac{\Gamma \vdash e : \text{Int}\#}{\Gamma \vdash \text{I}\#^\gamma e : \text{Int}^\gamma} \text{INT-I} \quad \frac{\Gamma \vdash e : \text{Int}^\gamma \quad \Gamma, x:\text{Int}\# \vdash e' : \tau}{\Gamma \vdash \text{case } e \text{ of } \text{I}\# x \rightarrow e' : \tau} \text{INT-E} \\
\\
\frac{\Gamma, x:\tau \vdash e : \sigma \quad \Gamma \vdash \tau \text{ mono-rep}}{\Gamma \vdash \lambda x:\tau. e : \tau \rightsquigarrow \sigma} \text{LAM-I} \quad \frac{\Gamma \vdash e : \tau \rightsquigarrow \sigma \quad \Gamma \vdash \tau \text{ mono-rep} \quad \Gamma \vdash \tau \text{ mono-conv}}{\Gamma \vdash e e' : \sigma} \text{LAM-E} \\
\\
\frac{\Gamma \vdash e : \tau \rightsquigarrow \sigma \quad \Gamma \vdash P : \tau \quad \Gamma \vdash \tau \text{ mono-rep}}{\Gamma \vdash e P : \sigma} \text{LAM-P-E} \\
\\
\frac{\Gamma \vdash e : \tau \quad \Gamma \vdash \tau \text{ mono-conv}}{\Gamma \vdash \text{Clos}^\gamma e : \gamma\{\tau\}} \text{CLO-I} \quad \frac{\Gamma \vdash e : \gamma\{\tau\}}{\Gamma \vdash \text{App } e : \tau} \text{CLO-E} \\
\\
\frac{\Gamma, \chi \vdash e : \sigma \quad \Gamma \vdash \forall \chi. \sigma : \kappa}{\Gamma \vdash \lambda \chi. e : \forall \chi. \sigma} \forall_I \quad \frac{\Gamma \vdash e : \forall \chi. \sigma \quad \Gamma \vdash [\phi/\chi] \text{ poly}}{\Gamma \vdash e \phi : \sigma[\phi/\chi]} \forall_E \\
\\
\text{Polymorphic instantiations: } \boxed{\Gamma \vdash [\phi/\chi] \text{ poly}} \\
\\
\frac{\Gamma \vdash \gamma \text{ lev}}{\Gamma \vdash [\gamma/g] \text{ poly}} \quad \frac{\Gamma \vdash \rho \text{ rep}}{\Gamma \vdash [\rho/r] \text{ poly}} \quad \frac{\Gamma \vdash \nu \text{ conv}}{\Gamma \vdash [\nu/v] \text{ poly}} \quad \frac{\Gamma \vdash \tau : \kappa \quad \Gamma \vdash \kappa \text{ kind}}{\Gamma \vdash [\tau/t:\kappa] \text{ poly}} \\
\\
\text{Types of constants: } \boxed{c : \tau} \\
\\
n : \text{Int}\# \quad \text{error} : \forall r. \forall g. \forall t:\text{TYPE } r \text{ Eval}^g. \text{Int}\# \rightsquigarrow t
\end{array}$$

Fig. 2. Type system of \mathcal{IL}

representation, arity, and levity explicit lets us express when they do *not* need to be known—analogue to polymorphic code that is independent of a type—thereby improving code reuse.

The goal of this paper is thus to provide a type system in which the kind κ of a type $\tau:\kappa$ classifies both the *representation* and the *calling convention* of values of type τ . To give ourselves a solid foundation, we focus on the compiler’s *explicitly-typed intermediate language*, rather than the source language. Our intermediate language, which we call \mathcal{IL} , is an explicitly-typed λ -calculus based closely on System F, so much of its syntax (Fig. 1) and typing rules (Figs. 2 and 3) should be familiar. In particular the syntax for expressions e includes:

- Variables x , constants c .
- Two types for integer values, Int^γ (boxed, where γ refers to the levity, see Section 3.1) and $\text{Int}\#$ (unboxed), together with expression forms for explicitly boxing ($\text{I}\#^\gamma e$) and unboxing (**case**); see Section 3.2.
- Functions whose arity is statically known, introduced and eliminated with the familiar looking forms $\lambda x:\tau. e$ and $e_1 e_2$. In spite of apparent familiarity, the notational change in the type to a wavy arrow ($\tau_1 \rightsquigarrow \tau_2$), rather than the more conventional ($\tau_1 \rightarrow \tau_2$), signals some rather subtle, but crucial, semantic differences that we detail in Section 3.3.
- Functions of unknown arity, introduced and eliminated with the new forms $\text{Clos}^\gamma e$ and $\text{App } e$, respectively, are also detailed in Section 3.3.

Kinds of types $\boxed{\Gamma \vdash \tau : \kappa}$	
$\frac{\Gamma \vdash \kappa \text{ kind}}{\Gamma, t : \kappa \vdash t : \kappa} \text{T}_{\text{VAR}}$	$\frac{\Gamma \vdash \gamma \text{ lev}}{\Gamma \vdash \text{Int}^\gamma : \text{TYPE PtrR Eval}^\gamma \text{Int}} \text{Int}^\#$
$\frac{\Gamma \vdash \gamma \text{ lev} \quad \Gamma \vdash \tau : \text{TYPE } \rho \ v}{\Gamma \vdash \gamma\{\tau\} : \text{TYPE PtrR Eval}^\gamma} \text{T}_{\text{CL}}$	$\frac{\Gamma, \chi \vdash \sigma : \text{TYPE } \rho \ v \quad \Gamma \vdash \text{TYPE } \rho \ v \text{ kind}}{\Gamma \vdash \forall \chi. \sigma : \text{TYPE } \rho \ v} \text{FORALL}$
$\frac{\Gamma \vdash \tau_1 : \text{TYPE } \rho_1 \ v_1 \quad \Gamma \vdash \tau_2 : \text{TYPE } \rho_2 \ v_2}{\Gamma \vdash \tau_1 \rightsquigarrow \tau_2 : \text{TYPE PtrR Call}[\rho_1, \text{arity}(v_2)]} \text{ARROW}$	$\frac{\Gamma \vdash \tau : \kappa \quad \kappa = \kappa'}{\Gamma \vdash \tau : \kappa'} \text{K-CONV}$
Reflexivity, transitivity, symmetry, and compatibility for $\boxed{\kappa = \kappa'}$, plus the following rules:	
$\text{arity}(\text{Eval}^\gamma) = \varepsilon \quad \text{arity}(\text{Call}[\alpha]) = \alpha$	
Monomorphism restrictions:	
$\frac{\Gamma \vdash \tau : \text{TYPE } \rho \ v \quad \vdash \rho \text{ rep}}{\Gamma \vdash \tau \text{ mono-rep}}$	$\frac{\Gamma \vdash \tau : \text{TYPE } \rho \ v \quad \vdash v \text{ conv}}{\Gamma \vdash \tau \text{ mono-conv}}$
Formation rules for kinds, levities, representations, calling conventions, arities	
$\boxed{\Gamma \vdash \kappa \text{ kind}}$	$\frac{\Gamma \vdash \rho \text{ rep} \quad \Gamma \vdash v \text{ conv}}{\Gamma \vdash \text{TYPE } \rho \ v \text{ kind}}$
$\boxed{\Gamma \vdash \gamma \text{ lev}}$	$\overline{\Gamma \vdash L \text{ lev}} \quad \overline{\Gamma \vdash U \text{ lev}} \quad \overline{\Gamma, g \vdash g \text{ lev}}$
$\boxed{\Gamma \vdash \rho \text{ rep}}$	$\overline{\Gamma \vdash \text{PtrR rep}} \quad \overline{\Gamma \vdash \text{IntR rep}} \quad \overline{\Gamma, r \vdash r \text{ rep}}$
$\boxed{\Gamma \vdash v \text{ conv}}$	$\frac{\Gamma \vdash \gamma \text{ lev}}{\Gamma \vdash \text{Eval}^\gamma \text{ conv}} \quad \frac{\Gamma \vdash \alpha \text{ ari}}{\Gamma \vdash \text{Call}[\alpha] \text{ conv}} \quad \overline{\Gamma, v \vdash v \text{ conv}}$
$\boxed{\Gamma \vdash \alpha \text{ ari}}$	$\overline{\Gamma \vdash \varepsilon \text{ ari}} \quad \frac{\Gamma \vdash \rho \text{ rep} \quad \Gamma \vdash \alpha \text{ ari}}{\Gamma \vdash \rho, \alpha \text{ ari}} \quad \frac{\Gamma \vdash v \text{ conv}}{\Gamma \vdash \text{arity}(v) \text{ ari}}$

Fig. 3. Kind and levity system of \mathcal{IL}

- Type abstraction $\lambda \chi. e$ and application $e \phi$. The only unusual thing here is that binders χ and arguments ϕ range over four different sorts of (erasable) variables and arguments respectively. We use a consistent naming convention, with a Latin-font name for variables (t, g, v , and r) and the corresponding Greek-font name for the syntactic category (τ, γ, v , and ρ respectively).

Note that a subset of these expressions are *Passive* (denoted by P), meaning that they do not require evaluation when passing them to a function or binding them to a variable. Of note, this classification accounts for eventual *type erasure* (which is why $\lambda \chi. e$ is passive only when e is and $P \phi$ is passive because ϕ arguments do not survive compilation). Passive expressions include *closures* $\text{Clos}^\gamma e$ as usual, but not *functions* $\lambda \chi : \tau. e$ because functions are *called* instead of *evaluated*. Types τ, σ 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, Int), polymorphic types $\forall \chi. \sigma$, function types $\tau \rightsquigarrow \sigma$ and a closure type $\gamma\{\tau\}$.

3.1 Kinds, representations, calling conventions, and levities

A kind κ has the form $\text{TYPE } \rho \ v$, where ρ describes the *representation* of the type, and v describes its *calling convention*, that is, what operations are allowed on that type. We do not support higher kinds, nor kind polymorphism; adding either is entirely straightforward, but introduces distracting details. Here we focus on the essentials of the intensional details in kinds.

Suppose $x : \tau : \text{TYPE } \rho \ v$; that is, x is a term variable of type τ , whose kind is $\text{TYPE } \rho \ v$. We now describe what ρ and v mean. The *representation* ρ describes how the value of x is represented at runtime. Referring to Fig. 1, ρ can be:

- PtrR , meaning that x is represented by a pointer into the garbage-collected heap.
- IntR , meaning that x is represented by a machine integer (not a pointer). In reality we would have many such kinds, for integers of different widths, for floating point values, and so on.
- r , a representation variable, which can be bound by a \forall ; that is, we support representation polymorphism [Eisenberg and Peyton Jones 2017].

The *convention* v describes how x may be consumed. Referring to Fig. 1, we see that v can be:

- Eval^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.
- Eval^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.
- Eval^g , where g is a levity variable; that is, we support levity polymorphism.
- $\text{Call}[\alpha]$, meaning that x is a function (not a thunk) with an *arity* described by α . The arity of a function is either a fixed list ρ_1, \dots, ρ_n , in which case x takes precisely n arguments, whose representations are given by ρ_i . Otherwise, the arity will have the form $\rho_1, \dots, \rho_n, \text{arity}(v)$, meaning that x takes *at least* n arguments, followed by possibly some more arguments given by the arity of another calling convention v . Call is a completely new concept: see Section 3.4.
- v , a calling-convention variable, which can be bound by a \forall ; that is, we support polymorphism over calling conventions.

The kind $\text{TYPE } \text{PtrR } \text{Eval}^L$ expresses the uniform representation of a value in a lazy language, as a pointer to a lifted (i.e. possibly a thunk) object. Because this kind is so common, 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 $\text{TYPE } \text{PtrR } \text{Eval}^U$ and Eval^L would be used sparingly, if at all.

3.2 Boxed and unboxed data types

Nearly thirty years ago, GHC introduced the idea of distinguishing boxed and unboxed data types in its intermediate language [Peyton Jones and Launchbury 1991]. We adopt this idea, but as a special case in a more general framework. Specifically, we have:

- A primitive type $\text{Int}\#$, of kind $\text{TYPE } \text{IntR } \text{Eval}^U$. The IntR says that $\text{Int}\#$ is represented by a machine integer; while the Eval^U says that it is unlifted (cannot be a thunk). The type $\text{Int}\#$ comes with literal constants n , and primitive operations op over it. We specify just one such primitive type, but in reality there would be many more, with different representations (machine integers of various widths, floating point numbers of various widths, etc).
- A boxed type Int^γ , of kind $\text{TYPE } \text{PtrR } \text{Eval}^\gamma$. This kind specifies that an Int^γ is represented by a pointer to a heap-allocated object (hence “boxed”), and that it may be lifted (if γ is L) or unlifted (if γ is U).

In reality there would be a way to declare new user-defined algebraic data types, and these types need not be levity polymorphic (i.e., have a γ argument). The sole goal of distinguishing two types is efficiency. If we had only boxed integers, then even simple addition would be forced to evaluate

and unbox each argument, and box up the result. By making these operations explicit we expose much more to the optimizer, and can eliminate lots of intermediate boxes. For example, consider:²

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

where `plus#` and `minus#` are primitive operations (extending *PrimOp* from Fig. 1), both of type $\text{Int\#} \rightsquigarrow \text{Int\#} \rightsquigarrow \text{Int\#}$.

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

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

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 allocation. `sumFrom` becomes just a wrapper around `sumFrom#` which handles all of the issues of boxing and unboxing. Instead of allocating several extraneous boxes at each step of the loop as before, the optimized code will only ever unbox the given number once at the start and then allocate the final box at the end. This is a huge gain in both time and space!

We add one new refinement, however: we can distinguish Int^l from Int^u , thus reusing the same function—and eventually generating the same code—for both strict and lazy languages. We will return to this point after having introduced levity polymorphism (see section Section 3.5).

3.3 Lambdas and closures

Consider these two function definitions in, say, OCaml:

$$f_1 = \lambda x. \lambda y. \text{let } z = \text{print } x \text{ in plus\# } x \ y \quad f_2 = \lambda x. \text{let } z = \text{print } x \text{ in } \lambda y. \text{plus\# } x \ y$$

These two functions are η -equivalent to one another, and yet, they are not the same! The side effect of printing `x` reveals the difference when the functions are partially applied. For example, `let z = f1 10 in 20` will just return 20, but `let z = f2 10 in 20` will first print 10 and then return 20. For this reason, the OCaml compiler cannot freely η -convert between `f1` and `f2`, due to side effects and eagerness changing the result of a partial application.

Even in a pure, lazy language like Haskell, a similar issue arises once we think about operational concerns. Consider now these two Haskell function definitions [Downen et al. 2019]:

$$\begin{array}{l} f_1, f_2 : \text{Int\#} \rightarrow \text{Int\#} \rightarrow \text{Int\#} \\ f_1 = \lambda x. \lambda y. \text{let } z = h \ x \ x \ \text{in } e \ y \ z \\ f_2 = \lambda x. \text{let } z = h \ x \ x \ \text{in } \lambda y. e \ y \ z \end{array}$$

Suppose we have a call `(f1 10 20)`. Currying suggests that we should first call `(f1 10)`, returning a (heap-allocated) function closure, and then call that closure with 20. This is unacceptably inefficient in practice; instead, we want to pass 10 and 20 *simultaneously*, with no intermediate function closure. Doing so is fine, because, `f1` can do no work until it is applied to two arguments.

But `f2` is quite different: we *should* pass the arguments one at a time. Imagine we had the call `map (f2 10) xs`: then we want to share the computation of `(h 10 10)` among all the elements of `xs`. Unlike `f1`, function `f2` can do (potentially expensive) work when applied to one argument.

²Here and elsewhere we reduce notational clutter by using syntactic sugar such as pattern matching.

In this case, we could η -expand f_2 anyway, to get f_1 . Now calls to f_2 would be more efficient—but at the cost of an asymptotic slowdown when used with *map*. We also saw previously that if we have side effects, like in OCaml, then η -expansion changes *semantics*, not just *efficiency*. Even Haskell has the *seq* operator, which also makes η -expansion semantically unsound. In short, unrestricted η -expansion is not allowed in any realistic source language.

In Haskell the distinction between f_1 and f_2 is not apparent in their types; indeed they both have the same type. The idea of Downen et al. [2019] is to express arity (*i.e.*, the ability to do computationally-safe η -expansion) in the types. In particular, in \mathcal{IL} *any expression e whatsoever of type $\tau_1 \rightsquigarrow \tau_2$ can safely be η -expanded to $\lambda x. e\ x$* .

What, then, of currying? If f_2 is eta-expandable to f_1 without changing the computational cost, how can we achieve the sharing desired by the Haskell programmer above? We can do it like this:

$$\begin{aligned} f_3 &:: \text{Int}\# \rightsquigarrow {}^L\{\text{Int}\# \rightsquigarrow \text{Int}\#\} \\ f_3 &= \lambda x. \text{let } z = h\ x\ x \text{ in Clos}^L(\lambda y. e\ y\ z) \end{aligned}$$

Here the “return a function closure” part of currying has become fully explicit in the expression $\text{Clos}^L e$, and is reflected in f_3 ’s type by ${}^L\{\text{Int}\# \rightsquigarrow \text{Int}\#\}$.

How can f_3 be called? Clearly we cannot write $(f_3\ 10\ 20)$ because the application $(e_1\ e_2)$ requires that $e_1 : \tau_1 \rightsquigarrow \tau_2$. Instead, Clos comes with its dual elimination form, App , and we call f_3 by writing $(\text{App}\ (f_3\ 10)\ 20)$. This call directly expresses the idea of applying f_3 to one argument 10, and then applying the resulting function closure to the second argument 20. The typing rules for Clos and App are just as you would expect: see rules CLO-I and CLO-E in Fig. 2.

3.4 Kinds are calling conventions

We motivated the idea that the *type* of a function could express its *arity*; just count the arrows. For example, $\text{Int}\# \rightsquigarrow \text{Int}\# \rightsquigarrow \text{Int}\#$ is the type of an arity 2 function. But suppose we have

$$id : \forall(t:\text{TYPE PtrR } \nu). t \rightsquigarrow t$$

What is *id*’s arity? You might reasonably answer “one,” but what about when t is instantiated to $(\text{Int}\# \rightsquigarrow \text{Int}\#)$? That type application presumably has type $(\text{Int}\# \rightsquigarrow \text{Int}\#) \rightsquigarrow \text{Int}\# \rightsquigarrow \text{Int}\#$, which appears to have arity 2. Suddenly it is not clear how many arguments *id* expects.

One way to address this problem, used by [Downen et al. 2019], is to *use the kind system to prevent instantiating type variable with an arrow type*. But that is a very drastic restriction: it means you cannot have, say, a list of functions; instead you would be forced to wrap them in Clos .

We take a more powerful approach here: instead of counting the arrows in the function’s type (which sometimes fails in the presence of type polymorphism), we simply look at the type’s kind. That kind expresses the function’s arity (and argument representations). 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 calling-convention component in these kinds is $\text{Call}[\alpha]$, where the arity list α describes the number and representation of arguments expected by the function. We need the latter so that, during compilation (Section 4), we can eta-expand a function to its full arity.

The kinding rule for $(\tau_1 \rightsquigarrow \tau_2)$ must track these arities. Looking at rule ARROW in Fig. 3 we see that the kind of $(\tau_1 \rightsquigarrow \tau_2)$ has a calling convention of $\text{Call}[\rho_1, \text{arity}(\nu_2)]$. The first argument has representation ρ_1 , the representation of τ_1 . The rest of the arguments come from the calling convention of τ_2 , via the type-level function *arity*, whose behavior is defined in Fig. 3. It returns the arity of $\text{Call}[\alpha]$; and the empty arity in the case of Eval^L since no arguments are needed

to evaluate a non-function. But v_2 might also be a variable v , and then $\text{arity}(v)$ is stuck; that is why $\text{arity}(v)$ is part of the syntax of α in Fig. 1. Rule $\kappa\text{-CONV}$ allows calls to arity to be calculated whenever desired.³

Here at last we see the key contribution of this paper: the calling convention of a function is expressed in the kind of its type, rather than in the type itself. What forces us to use kinds rather than types? Answer: the desire for abstraction, and specifically polymorphism. In a monomorphic system, “types are calling conventions” works just fine.

3.5 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, consider the error function—a motivating example of [Eisenberg and Peyton Jones 2017]—which takes an error code of type $\text{Int}\#$ and halts the program with some message. It can be given this type:

$$\text{error} : \forall r\ g. \forall (t : \text{TYPE } r \text{ Eval}^g). \text{Int}\# \rightsquigarrow t$$

Here, *error* is polymorphic in both the representation r and the levity g of the returned value because, in fact, it never returns a value. Lacking this polymorphism, we would be forced to define a whole family of monomorphic versions of the *error* function, one for each return representation and levity, which would be extremely painful—especially since they compile to the same code.

Using polymorphism over levity also allows us to write some functions that work uniformly over both strict and lazy values. For example, adding two boxed integers can be defined thus

$$\begin{aligned} \text{plus} &: \forall g_1. \forall g_2. \forall g_3. \text{Int}^{g_1} \rightsquigarrow \text{Int}^{g_2} \rightsquigarrow \text{Int}^{g_3} \\ \text{plus } g (\text{I}\# x) (\text{I}\# y) &= \text{I}\#^{g_3} (\text{plus}\# x y) \end{aligned}$$

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 Int^{g_i} are statically unknown, so we must be able to pattern-match on and return values with unknown levities. Specifically, rule CASE in Fig. 2 allows a case-expression to scrutinize an integer of arbitrary levity γ . Operationally, in a polymorphic situation the case-expression has to test the scrutinee to see if it is a thunk (in case the 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 case for handling a thunk is simply dead code.⁴

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

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

From a representation point of view, an $\text{Array}\#^Y$ 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 $\text{Array}\#^{\text{L}} (\text{Array}\#^{\text{U}} \text{Int}^{\text{U}})$ is a lifted array of pointers, each of which points directly

³Alternatively, we could require $\text{arity}(v)$ to be fully calculated in the ARROW kinding rule. This would let us remove $\text{arity}(v)$ from the grammar of arities, but also forces an additional restriction on the formation of types and expressions, specifically ARROW and LAM-I , to rule out $\text{arity}(v)$. The cost of such a restriction is to break the existing property that, except for the \forall quantifier, any type made from well-kinded types is itself well-kinded.

⁴We assume here that the concrete, run-time representation of evaluated lifted values is the same as the representation of unlifted values. This is true in GHC and seems likely in other systems that support laziness, but it is conceivably an invalid assumption in some systems.

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].

3.6 Polymorphism in calling convention

We may also be polymorphic in calling 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 another⁵. 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 v , thus:

$$(4) \quad \text{revapp} : \forall v. \forall (t:\text{TYPE PtrR } v). t \rightsquigarrow (t \rightsquigarrow \text{Int}\#) \rightsquigarrow \text{Int}\#$$

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), so we can be completely polymorphic in its representation, thus:

$$(5) \quad \text{revapp} : \forall v \ r \ g. \forall (t_1:\text{TYPE PtrR } v) (t_2:\text{TYPE } r \text{ Eval}^g). t_1 \rightsquigarrow (t_1 \rightsquigarrow t_2) \rightsquigarrow t_2$$

But notice that, unlike the argument type t_1 , *revapp* cannot be polymorphic in the calling convention of t_2 , as we discussed in Section 3.4. It can, however, be evaluated with any levity.

3.7 Restrictions on polymorphism

Of course we have the usual restrictions on polymorphism,⁶ but the polymorphism \mathcal{IL} introduces some new issues. We have already seen how unrestricted polymorphism is incompatible with efficient static code generation⁷ in Section 3.6, where we cannot allow *revapp*'s type argument t_1 to have a representation-polymorphic kind. A second restriction is demonstrated by this function:

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

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, unlifted types are always values, never thunks. So we can give *twice* either of these types:

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

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

⁵Even 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.

⁶For example, every free variable that appears anywhere in the type checking judgment $\Gamma \vdash e : \tau$ must be bound by Γ .

⁷Runtime 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 a static code generation.

The restrictions on polymorphism in our language are embodied in the shaded premises in the type system of Fig. 2. The judgment $\Gamma \vdash \tau$ **mono-rep**, defined in Fig. 3 checks that the representation of τ 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. 2 we see the shaded premises:

- Rule **LAM-E**: for a general application, the argument type must be monomorphic in both the representation (so that we know how to pass it to the function), and convention (so that we know when to evaluate it, or for first-class functions as arguments, how to define it).
- Rule **CLO-I**: this rule constructs a closure of an unknown arity from an expression with a known arity (as we elaborated in Section 3.3). As such, it needs to know the arity of the contained expression in order to define the closure.

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 4, how to compile \mathcal{IL} into a lower level “machine language” \mathcal{ML} . In the translation from \mathcal{IL} to \mathcal{ML} in Fig. 6, we need exactly the shaded monomorphic restrictions of Fig. 2. If any of these restrictions were removed, then there would be expressions that are well-typed, yet un-compilable.

Our rules also include two additional, unshaded, monomorphism restrictions, in the **LAM-I** and **LAM-P-E** rules. These restrictions enforce an extra invariant on the environment Γ : *every variable in Γ 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 4 *does not* require any monomorphism restrictions in **LAM-I** and **LAM-P-E**: they could be deleted and yet all closed, well-typed expressions could still be compiled. This example shows how 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.

3.8 The FORALL rule

The polymorphic quantifier $\forall t:\kappa.\sigma$ has no impact on a type’s kind: it just inherits the kind of σ (rule **FORALL** in Fig. 3). Intuitively, this is because these quantifiers will be totally *erased* by compilation, and have no impact on the final run-time code. Since kinds are meant to reflect the representations and actions that occur during run-time, the \forall is “invisible” to the lower-level machine.

However, now that variables may appear in kinds, we must be careful to avoid such a variable escaping its 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 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. 3 and the second premise of rule \forall_I in Fig. 2.

3.9 Let bindings

\mathcal{IL} does not have a **let**-binding construct but, as usual, a non-recursive **let** can be regarded as shorthand for lambda and application: **let** $x:\tau = e$ **in** $e' \triangleq (\lambda x:\tau. e') e$. The typing rule follows by composing **LAM-I** and **LAM-E**:

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

Type-based passable expressions (<i>i.e.</i> , substitutable values): $\boxed{\Gamma \vdash e \text{ pass}}$			
$\Gamma \vdash P \text{ pass}$	$\frac{\Gamma \vdash e : \tau \quad \Gamma \vdash \tau : \text{TYPE } \rho \text{ Eval}^L}{\Gamma \vdash e \text{ pass}}$	$\frac{\Gamma \vdash e : \tau \quad \Gamma \vdash \tau : \text{TYPE } \rho \text{ Call}[\alpha]}{\Gamma \vdash e \text{ pass}}$	
Equational axioms (in each rule, assume that $\Gamma \vdash S \text{ pass}$):			
(β_{\sim})	$(\lambda x:\tau.e) S = e[S/x]$	(η_{\sim})	$\lambda x:\tau.(e x) = e : \tau \rightsquigarrow \sigma$
(β_{\forall})	$(\lambda \chi.e) \phi = e[\phi/\chi]$	(η_{\forall})	$\lambda \chi.(e \chi) = e : \forall \chi.\sigma$
$(\beta_{\{\}})$	$\text{App}(\text{Clos}^Y e) = e$	$(\eta_{\{\}})$	$\text{Clos}^Y(\text{App } S) = S : {}^Y\{\sigma\}$
(β_{Int})	$\text{case } I\#^Y P \text{ of } = e[P/x]$	(η_{Int})	$\text{case } e \text{ of } = e : \text{Int}^Y$
	$I\# x \rightarrow e$		$I\# x \rightarrow I\#^Y x$
Plus closure under reflexivity, transitivity, symmetry, and compatibility.			

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

That is, the right-hand side of a **let** must have statically known representation and calling convention. We can add a similar derived typing rule **LET-P**, by composing **LAM-I** and **LAM-P-E**, for the special case when the right hand side is passive so its convention does not need to be known.

3.10 Equational theory

The equational theory for \mathcal{IL} is defined in Fig. 4. It gives us a framework to reason about equality in \mathcal{IL} , and ultimately about the correctness of compiling \mathcal{IL} to a low-level language (Theorem 1) as well as for compiling a high-level language to \mathcal{IL} (Theorem 4). The rules for **Clos/App** (namely $\beta_{\{\}}$ and $\eta_{\{\}}$) and **case**, $I\#$ (β_{Int} and η_{Int}) are unsurprising, as are the rules for erasable abstractions (β_{\forall} and η_{\forall}). More distinctive is η_{\sim} , which (as discussed in Section 3.3) allows unrestricted η in either direction for any expression of a function type.

That leaves the reduction rule β_{\sim} . As is usual in a call-by-value λ -calculus, only some expressions—called *passable* or *substitutable* expressions, and here denoted by the metavariable S —can be passed to a function and substituted for its formal parameter by the β_{\sim} rule. Only if the argument of a β -redex is passable does the β_{\sim} rule fire. Unlike most systems, which use a syntactic definition of the expressions that can be passed or substituted, \mathcal{IL} instead identifies these expressions by their type and kind, as defined by the rules for the $\Gamma \vdash e \text{ pass}$ judgment. Using the type and kind of the argument to define passability, rather than only its syntax, allows us to integrate several different evaluation orders (call-by-value, call-by-name, *etc.*) within the same language.

Let us examine closely the $\Gamma \vdash e \text{ pass}$ judgment. Firstly, we designate all *passive* expressions P to be *passable*. These expression forms include variables (which would always be substituted for passable expressions), constants, and applications of $I\#$ and **Clos**. Notice that, thus far, these are all considered values in the call-by-value λ -calculus. The grammar for P also looks through abstractions $\lambda \chi.P$ and applications $P \phi$, which are erased anyway during compilation, and thus have no impact on the runtime behavior of a program.

Notice how the notion of passability depends specifically on the calling convention for the type of that expression. For example, in a call-by-name setting, every expression can be substituted for a variable. So in \mathcal{IL} , *all* arguments that are to be lazily evaluated—which have the convention Eval^L —are considered passable, and hence allow β_{\sim} to fire. In contrast, only expressions that are

passive (P) are values in call-by-value languages. So we can only substitute a variable with an expression of the convention Eval^U —which should be strictly evaluated—when it is passive. More generally, passive expressions are always substitutable in all of the evaluation strategies we are interested in here, so we can say something more: a passive P is passable for *all* calling conventions. This extra step is helpful in case we are dealing with an expression—like x or $I^\#^g n$ —which has an unknown calling convention but will inevitably be passable 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 3.5) 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 β_{\sim} . This is possible because the argument has the type $\text{Int}^L : \text{TYPE PtrR Eval}^L$, and so it is passable. 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 β_{\sim} directly as before, because the argument is not passable: it has the type $\text{Int}^U : \text{TYPE PtrR Eval}^U$, and it is not passive. Instead we must first evaluate the argument; now it becomes passive, and β_{\sim} can fire.

Of course the substitution done by β_{\sim} is woefully inefficient: worse than just recursing down an expression, it irreparably *duplicates* the work of delayed computations, leading to an asymptotic slowdown in many cases. A better semantics cares about sharing work, which we show next in Section 4 when compiling to a lower-level representation. But we could still reason about complexity and sharing directly in \mathcal{IL} using a call-by-need semantics, just like the λ -calculus [Ariola and Felleisen 1997]. For more details about the operational semantics of \mathcal{IL} , see Appendix A.⁸

4 COMPILATION FROM \mathcal{IL} TO A LOWER LEVEL

4.1 Syntax and semantics of a machine language \mathcal{ML}

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 non-uniform representations and higher-arity functions, both of which must be monomorphic. The idea is to model a realistic machine architecture that supports multiple basic kinds of data representations (e.g., pointers vs. integers), and where functions are passed multiple arguments but are otherwise first order (though function pointers may be passed as arguments and invoked). The syntax of this language, which we call \mathcal{ML} , is given in Fig. 5, along with its abstract machine⁹.

The syntax of \mathcal{ML} is restricted from the more λ -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.
- Functions $(\lambda(\overline{y}).e)$ and their calls $(P(\overline{a}))$ pass many arguments at once, explicitly modeling multi-arity functions.
- The function components of an application form cannot be an arbitrary expression; it must be a passive expression, which is either a value V or a variable. This way, every application $P(\overline{a})$ can be resolved in at most two steps: lookup P if it is a variable, and then apply P if it is a λ -abstraction or a constant (like error).

Note that, as such, it is impossible to chain separate calls in a row. For example, $f(1)(2)$ is not a legal expression in \mathcal{ML} ; instead, if f is an arity 2 function, it must be called as $f(1, 2)$, while if it is

⁸All appendices are included as anonymized supplementary material.

⁹For aficionados of GHC, \mathcal{IL} is like Core language, while \mathcal{ML} is like the STG language [Peyton Jones 1992].

$\pi \in \text{PrimRep} ::= \text{PtrR} \mid \text{IntR}$	$\psi \in \text{Strictness} ::= \text{L} \mid \text{U}$	$\eta \in \text{KnownCC} ::= \text{Eval}^\psi \mid \text{Call}[\bar{\pi}]$
$e \in \text{Expr} ::= P \mid P(\bar{a}) \mid \text{App } e(\bar{a}) \mid \text{case } e \text{ of } \text{I\#}(x_{\text{IntR}}) \rightarrow e' \mid \text{let } x_\pi^\psi = e \text{ in } e'$		
$R \in \text{Reference} ::= \text{I\#}(a) \mid \text{Clos}^n P \mid \lambda(\bar{x}_\pi).e$		
$V \in \text{Value} ::= c \mid \text{I\#}(a) \mid \text{Clos}^n P$	$P \in \text{Passive} ::= R \mid x_\pi$	$a \in \text{Arg} ::= c \mid x_\pi$
$K \in \text{StackCont} ::= \varepsilon \mid \text{App}(\bar{a}); K \mid \text{case } \text{I\#}(x_{\text{IntR}}) \rightarrow e; K \mid \text{let } x_\pi \text{ in } e; K \mid \text{set } x; K$		
$H \in \text{Store} ::= \varepsilon \mid [x := R]H \mid [x := \text{memo } e]H$		
$m \in \text{MachineState} ::= \langle e \mid K \mid H \rangle \mid \text{error}(n)$	$A \in \text{Answer} ::= \langle V \mid \varepsilon \mid H \rangle \mid \text{error}(n)$	
(PshApp)	$\langle \text{App } e(\bar{a}) \mid K \mid H \rangle \mapsto \langle e \mid \text{App}(\bar{a}); K \mid H \rangle$	
(PshCase)	$\langle \text{case } e \text{ of } \text{I\#}(x_{\text{IntR}}) \rightarrow e' \mid K \mid H \rangle \mapsto \langle e \mid \text{case } \text{I\#}(x_{\text{IntR}}) \rightarrow e'; K \mid H \rangle$	
(PshLet)	$\langle \text{let } x_\pi^{\text{U}} = e \text{ in } e' \mid K \mid H \rangle \mapsto \langle e \mid \text{let } x_\pi \text{ in } e'; K \mid H \rangle$	
(LAlloc)	$\langle \text{let } x_{\text{PtrR}}^{\text{L}} = e' \text{ in } e \mid K \mid H \rangle \mapsto \langle e[y_{\text{PtrR}}/x_{\text{PtrR}}] \mid K \mid [y := \text{memo } e']H \rangle$	
(Call)	$\langle (\lambda(\bar{x}_\pi).e)(\bar{a}) \mid K \mid H \rangle \mapsto \langle e[\bar{a}/\bar{x}_\pi] \mid K \mid H \rangle$	
(Apply)	$\langle \text{Clos}^n P \mid \text{App}(\bar{a}); K \mid H \rangle \mapsto \langle P(\bar{a}) \mid K \mid H \rangle$	(if $n = \bar{a} $)
(Unbox)	$\langle \text{I\#}(a) \mid \text{case } \text{I\#}(x_{\text{IntR}}) \rightarrow e; K \mid H \rangle \mapsto \langle e[a/x_{\text{IntR}}] \mid K \mid H \rangle$	
(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)	$\langle R \mid \text{let } x_{\text{PtrR}} \text{ in } e; K \mid H \rangle \mapsto \langle e[y_{\text{PtrR}}/x_{\text{PtrR}}] \mid K \mid [y := R]H \rangle$	
(Fun)	$\langle y_{\text{PtrR}}(\bar{a}) \mid K \mid [y := R]H \rangle \mapsto \langle R(\bar{a}) \mid K \mid [y := R]H \rangle$	
(Look)	$\langle y_{\text{PtrR}} \mid K \mid [y := R]H \rangle \mapsto \langle R \mid K \mid [y := R]H \rangle$	
(Force)	$\langle y_{\text{PtrR}} \mid K \mid [y := \text{memo } e]H \rangle \mapsto \langle e \mid \text{set } y; K \mid [y := \text{memo } e]H \rangle$	
(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. 5. Syntax and semantics for \mathcal{ML}

an arity-1 function returning a closure of an arity-1 function it must be called as $(\text{App } f(1))(2)$. The number of arguments passed at once is explicitly fixed in each λ -abstraction defining a function and call site which uses it. In other words, \mathcal{ML} does not support polymorphism of function arity.

The syntax of \mathcal{ML} includes annotations that make explicit the semantically important aspect of types, which are implicit in \mathcal{IL} , so that it is clear from the syntax how to execute programs. In particular, each variable is annotated with its representation π , which must be either a pointer (PtrR) or integer (IntR). In other words, all variables are permanently assigned a representation—intuitively, specifying how they are stored in either an integer register or a pointer on the heap. By design, \mathcal{ML} does not support polymorphism over these representations because the machine's pointer registers and floating-point registers are distinct.

Additionally, each **let** binding is annotated as either eager (U) or lazy (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 these levities.

4.2 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 which 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 (PshApp , PshCase , PshLet) and allocating memory (LAlloc , VAlloc), but note that we include cases for *both* lazy bindings (LAlloc) and strict ones (PshLet , VAlloc).

Next we have the rules for performing interesting reductions. *Apply* resolves the application of a closure by extracting the function it contains, and *Call* calls a known 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 5). 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 called function to map to a value. For other pointers, we have to check whether or not it is evaluated, and merely *Look up* values or else *Force* thunks.¹⁰ When a forced thunk returns a value, it is *Memoized* to share the result on future uses of that pointer.

4.3 Compilation

The compilation transformation from \mathcal{IL} to the low-level machine language \mathcal{ML} is given in Fig. 6. The top-level compilation function is $\mathcal{E}_v \llbracket e \rrbracket_\theta^\Gamma$, which compiles a typed expression $\Gamma \vdash e : \tau$ where τ 's convention is v . The environment θ 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 function is to remember the distinction between calling and evaluating. In our system, only expressions with types like $\text{Int}\#$, 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 $\mathcal{E}_v \llbracket e \rrbracket_\theta^\Gamma$, which takes into account the calling convention v of e . If it is Eval^\vee , 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 function match, we always fully η -expand function expressions when they are defined: either on the right-hand side of `lets` or in the body of `Closures`. Because of this dependency, this step of compilation is only defined when the calling convention is statically known (e.g., it is not a variable v or partially-defined like $\text{Call}[r_1, r_2, \text{arity}(v)]$). In any case, we next move to the main work-horse of compilation, $\mathcal{C} \llbracket e \rrbracket_\theta^\Gamma(\bar{a})$, 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.

Even compiling constants and variables, at the top of the figure, brings in some complications. Literal constants are simply passed through, but when compiling a call to error, we see that

¹⁰Note that this uniform check on pointers y_{ptr} is needed to support levity polymorphism for types like Int^g and ${}^g\{\tau\}$. 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_{ptr} must be unlifted, so that the *Look* step always applies. This means that 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:

$$\mathcal{E}_v \llbracket P \rrbracket_\theta^\Gamma \triangleq \mathcal{P} \llbracket P \rrbracket_\theta^\Gamma \quad \mathcal{E}_{\text{Eval}^\gamma} \llbracket e \rrbracket_\theta^\Gamma \triangleq C \llbracket e \rrbracket_\theta^\Gamma(\varepsilon) \quad \mathcal{E}_{\text{Call}[\overline{\pi}]} \llbracket e \rrbracket_\theta^\Gamma \triangleq \lambda(\overline{x_\pi}). C \llbracket e \rrbracket_\theta^\Gamma(\overline{x_\pi})$$

Constants:

$$C \llbracket n \rrbracket_\theta^\Gamma(\varepsilon) \triangleq n \quad C \llbracket \text{error} \rrbracket_\theta^\Gamma(a) \triangleq \text{error}(a)$$

Variables:

$$C \llbracket x \rrbracket_\theta^{\Gamma, x:\tau}(\varepsilon) \triangleq \theta(x) \quad (\text{if } \Gamma \vdash \tau \rightsquigarrow^{\text{conv}} \text{Eval}^\psi)$$

$$C \llbracket x \rrbracket_\theta^{\Gamma, x:\tau}(\overline{a}) \triangleq (\theta(x))(\overline{a}) \quad (\text{if } \Gamma \vdash \tau \rightsquigarrow^{\text{conv}} \text{Call}[\overline{\pi}])$$

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

$$C \llbracket \text{App } e \rrbracket_\theta^\Gamma(\overline{a}) \triangleq \text{App } (C \llbracket e \rrbracket_\theta^\Gamma(\varepsilon))(\overline{a})$$

$$C \llbracket e \phi \rrbracket_\theta^\Gamma(\overline{a}) \triangleq C \llbracket e \rrbracket_\theta^\Gamma(\overline{a})$$

$$C \llbracket e P \rrbracket_\theta^\Gamma(\overline{a}) \triangleq C \llbracket e \rrbracket_\theta^\Gamma(\mathcal{P} \llbracket P \rrbracket_\theta^\Gamma, \overline{a}) \quad (\text{if } \mathcal{P} \llbracket P \rrbracket_\theta^\Gamma = x_\pi \text{ or } \mathcal{P} \llbracket P \rrbracket_\theta^\Gamma = c)$$

$$C \llbracket e P \rrbracket_\theta^\Gamma(\overline{a}) \triangleq \text{let } x_{\text{PtrR}}^U = \mathcal{P} \llbracket P \rrbracket_\theta^\Gamma \text{ in } C \llbracket e \rrbracket_\theta^\Gamma(x_{\text{PtrR}}, \overline{a})$$

$$C \llbracket e e' \rrbracket_\theta^\Gamma(\overline{a}) \triangleq \text{let } x_{\overline{\pi}}^{\text{lev}(\eta)} = \mathcal{E}_\eta \llbracket e' \rrbracket_\theta^\Gamma \text{ in } C \llbracket e \rrbracket_\theta^\Gamma(x_{\overline{\pi}}, \overline{a}) \quad (\text{if } \Gamma \vdash e':\tau, \Gamma \vdash \tau \rightsquigarrow^{\text{rep}} \pi, \text{ and } \Gamma \vdash \tau \rightsquigarrow^{\text{conv}} \eta)$$

Boxing and unboxing:

$$C \llbracket \text{I}\#^\gamma P \rrbracket_\theta^\Gamma(\varepsilon) \triangleq \mathcal{P} \llbracket \text{I}\#^\gamma P \rrbracket_\theta^\Gamma$$

$$C \llbracket \text{I}\#^\gamma e \rrbracket_\theta^\Gamma(\varepsilon) \triangleq \text{let } x_{\text{IntR}}^U = C \llbracket e \rrbracket_\theta^\Gamma(\varepsilon) \text{ in } \text{I}\#(x_{\text{IntR}})$$

$$C \llbracket \text{case } e' \text{ of } \text{I}\# x \rightarrow e \rrbracket_\theta^\Gamma(\overline{a}) \triangleq \text{case } C \llbracket e' \rrbracket_\theta^\Gamma(\varepsilon) \text{ of } \text{I}\#(x_{\text{IntR}}) \rightarrow C \llbracket e \rrbracket_{[x_{\text{IntR}}/x]}^{\Gamma, x:\text{Int}\#}(\overline{a})$$

Abstractions:

$$C \llbracket \lambda x:\sigma. e \rrbracket_\theta^\Gamma(a', \overline{a}) \triangleq C \llbracket e \rrbracket_{[a'/x]}^{\Gamma, x:\sigma}(\overline{a})$$

$$C \llbracket \lambda \chi. e \rrbracket_\theta^\Gamma(\overline{a}) \triangleq C \llbracket e \rrbracket_\theta^{\Gamma, \chi}(\overline{a})$$

$$C \llbracket \text{Clos}^\gamma e \rrbracket_\theta^\Gamma(\varepsilon) \triangleq \mathcal{P} \llbracket \text{Clos}^\gamma e \rrbracket_\theta^\Gamma$$

Passive expressions:

$$\mathcal{P} \llbracket c \rrbracket_\theta^\Gamma \triangleq c \quad \mathcal{P} \llbracket x \rrbracket_\theta^\Gamma \triangleq \theta(x) \quad \mathcal{P} \llbracket \text{I}\#^\gamma P \rrbracket_\theta^\Gamma \triangleq \text{I}\#(\mathcal{P} \llbracket P \rrbracket_\theta^\Gamma)$$

$$\mathcal{P} \llbracket P \phi \rrbracket_\theta^\Gamma \triangleq \mathcal{P} \llbracket P \rrbracket_\theta^\Gamma \quad \mathcal{P} \llbracket \lambda \chi. P \rrbracket_\theta^\Gamma \triangleq \mathcal{P} \llbracket P \rrbracket_\theta^\Gamma$$

$$\mathcal{P} \llbracket \text{Clos}^\gamma e \rrbracket_\theta^\Gamma \triangleq \text{Clos}^{\text{arity}(\eta)} \mathcal{E}_{\text{Call}[\text{arity}(\eta)]} \llbracket e \rrbracket_\theta^\Gamma \quad (\text{if } \Gamma \vdash e:\tau \text{ and } \Gamma \vdash \tau \rightsquigarrow^{\text{conv}} \eta)$$

$$\text{Calculating known representations and conventions: } \boxed{\Gamma \vdash \tau \rightsquigarrow^{\text{rep}} \pi} \text{ and } \boxed{\Gamma \vdash \tau \rightsquigarrow^{\text{conv}} \eta}$$

$$\frac{\Gamma \vdash \tau : \text{TYPE } \rho \vee \quad \rho = \pi}{\Gamma \vdash \tau \rightsquigarrow^{\text{rep}} \pi} \quad \frac{\Gamma \vdash \tau : \text{TYPE } \rho \vee \quad \vee = \eta}{\Gamma \vdash \tau \rightsquigarrow^{\text{conv}} \eta}$$

The levy of a known convention:

$$\text{lev}(\text{Eval}^\psi) = \psi \quad \text{lev}(\text{Call}[\overline{\pi}]) = \cup$$

Fig. 6. Compiling \mathcal{IL} to \mathcal{ML}

precisely one argument a is required. What if the user has written a partial application of error? Such partial applications are *always* η -expanded to be fully saturated, satisfying the requirement here. Compiling a variable x looks it up in the environment θ . However, note that there is different \mathcal{ML} code for *evaluating* x versus *calling* it, even when there are no arguments. In other words, a function variable x with the empty calling convention $\text{Call}[]$ compiles as $C[x]_{[y_{\text{PtrR}}/x]\theta}^{\Gamma}(\epsilon) = y_{\text{PtrR}}()$, which is a nullary function call, but a boxed integer variable x with the evaluation convention Eval^{\vee} compiles as $C[x]_{[y_{\text{PtrR}}/x]\theta}^{\Gamma}(\epsilon) = y_{\text{PtrR}}$ which merely looks up the pointer.

Closure applications are compiled straightforwardly, and an erasable arguments ϕ and binders $\lambda_{\chi}.e$ are simply dropped. Passive expressions of \mathcal{IL} can be compiled outright to passive expressions of \mathcal{ML} via $\mathcal{P}[P]_{\theta}^{\Gamma}$, which also performs type erasure. Compiling an application to a passive argument P depends on the nature of that passive argument:

- If $\mathcal{P}[P]_{\theta}^{\Gamma}$ is a variable or constant (after erasure type), then it can be passed directly.
- Otherwise, we name the argument with a **let** (respecting the A-normal form) and pass it by reference to the function. Notice that in this case, the compiled argument will always have the form $\text{I}\#(a)$ or $\text{Clos}^n P$, which means the **let**-binding will always be represented as a pointer into the heap.¹¹

In the variable case, we do not need to track the levity or representation of the argument, because any decisions around its convention have already been made, when the variable definition itself was compiled. Crucially, we did not have to look up any information in the typing environment to compile passive arguments; this is why no highlighted premises are needed in rule LAM-P-E.

In the case of a general application $e\ e'$ to an arbitrary argument that needs to be computed, corresponding to LAM-E, we always generate a **let** similar to the second case for $e\ P$. However, for LAM-E, we need to determine the representation, calling convention, *and* levity of the binding, which could truly be anything. This is why we require the additional restrictions in LAM-E that correspond to the highlighted side conditions of this case.

4.4 Correctness of compilation

Notice how the same polymorphism restrictions used in the typing rules also appear during compilation. Even though the defined compilation transformation is partial (that is, not every syntactically valid expression can be compiled), all well-typed \mathcal{IL} expressions with a known calling convention have a defined compilation to the lower-level, levity-monomorphic \mathcal{ML} . That is, we want to ensure that the compilation $\mathcal{E}_{\eta}[e]$ is well-defined for any well-typed closed expression $\vdash e : \tau : \text{TYPE } \rho\ \eta$, where the syntax of known calling conventions η is from Fig. 5.

In fact, we allow for a little more levity polymorphism during compilation: $\mathcal{E}_{\text{Eval}^{\vee}}[e]$, for an unknown levity g , is also allowed. That's because we generate the code that will be executed only when the expression is evaluated: in other words, when a computation is forced, there is no difference between eager or lazy. This added flexibility is essential for compiling levity polymorphic expressions appearing in strict contexts, such as in the discriminant of a **case** or argument of **App**. Although implicit, the C compilation function makes the same assumptions on the known convention and strictness of the expression it compiles. In contrast, the \mathcal{P} function does not assume the expression is being evaluated, because it is only 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 multi-arity functions) without erroneously introducing extra strictness.

¹¹It may be that in a more feature-rich language, this case could involve representations other than just a single pointer. However, the value returned from a passive expression is syntactically manifest. As such, the representation is always be apparent from the syntax of the expression, which means we do not require additional typing information here.

During compilation, we occasionally require the ability to calculate a known representation (π) or calling convention (η) for a sub-expression. This appears in Fig. 6 as the highlighted side conditions $\Gamma \vdash \tau \overset{rep}{\rightsquigarrow} \pi$ and $\Gamma \vdash \tau \overset{conv}{\rightsquigarrow} \eta$, respectively. In general, these calculations could fail if the representation of convention in the kind of τ are partially unknown—that is, contains free variables. But any closed representation has the form π and any closed convention is equivalent to a η . The places where this requirement appears corresponds exactly to the highlighted monomorphism restrictions in Fig. 2. For more details on the invariants for compiling open expressions, see Appendix B.

Theorem 1 (Closed Compilation). *If $\vdash e : \tau$ and $\vdash \tau : \text{TYPE } \rho \vee$ then $\mathcal{E}_v[[e]]$ is defined.*

This theorem just states when compilation translates a closed \mathcal{IL} expression to \mathcal{ML} code. We should also expect that compilation preserves the behavior of \mathcal{IL} expressions as well. In other words, if an expression is equal to some answer in \mathcal{IL} , then executing the compiled code should give the same answer in \mathcal{ML} . But note that we are not interested in evaluating functions directly—they must be 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. For more details on the semantic relationship between \mathcal{IL} and \mathcal{ML} , see Appendices A and B in the supplementary material.

Theorem 2 (Soundness and Completeness).

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

5 DYNAMIC ARITY RAISING

So far, both the intermediate language \mathcal{IL} and lower-level abstract machine \mathcal{ML} assume that the arity of all functions are statically known at compile-time. This is not the most ideal situation, especially since it has been shown that arity mismatches of “unknown” function calls can be efficiently accommodated in practical implementations [Marlow and Peyton Jones 2004]. This section thus proposes an optional extension to \mathcal{ML} , granting more flexibility at a small runtime cost when unpacking closures. Such additional flexibility requires only a little dynamic checking at run-time. Applying too few arguments creates a partial application, and applying too many is broken down into two (or more) separate calls. Both of these possibilities are already accounted for by Marlow and Peyton Jones [2004] with “unknown” function calls. In terms of \mathcal{ML} , we can extend its grammar with partial applications of the form $\text{Clos}^n f(\bar{a})$, where $n > 0$ is the number of remaining arguments expected before the function f can be called, and \bar{a} are the arguments applied so far. We can now express the extra rules for dynamic handling a run-time arity mismatch:

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

In practice, this lets us get away with *fusing* adjacent chains of Closures and Applications, so that the runtime system can dynamically choose a better calling convention when possible.

For example, consider the following two different closure wrappers around the `plus#` function:

$plus_1 : \text{Int}\# \rightsquigarrow \text{Int}\# \rightsquigarrow \text{Int}\#$ $plus_2 : \text{Int}\# \rightsquigarrow \text{Int}\# \rightsquigarrow \text{Int}\#$
 $plus_1 = \text{Clos}^1 \lambda x:\text{Int}\#. \text{Clos}^1 \lambda y:\text{Int}\#. \text{plus}\# x y$ $plus_2 = \text{Clos}^1 \lambda x:\text{Int}\#. \lambda y:\text{Int}\#. \text{plus}\# x y$

$plus_1$ is written to receive one argument at a time, whereas $plus_2$ expects two arguments at once, and the difference between the two calling conventions is expressed in *both* the expression and its

type. These arities are made more explicit when compiling to the lower level:

$$\begin{aligned}\mathcal{E}_{\text{Eval}^L}[\![plus_1]\!] &= \text{Clos}^1 \lambda(x_{\text{IntR}}). \text{Clos}^1 \lambda(y_{\text{IntR}}). \text{plus}\#(x_{\text{IntR}}, y_{\text{IntR}}) \\ \mathcal{E}_{\text{Eval}^L}[\![plus_2]\!] &= \text{Clos}^2 \lambda(x_{\text{IntR}}, y_{\text{IntR}}). \text{plus}\#(x_{\text{IntR}}, y_{\text{IntR}})\end{aligned}$$

If we support only the *Apply* rule at runtime, as done in Fig. 5, then this difference is important not only for performance, but *also* correctness: *Apply* fails if the arity of a function does not match the arity of an application exactly. However, with *PApp* and *OApp*, it is correct to use *plus₁* and *plus₂* interchangeably at runtime, albeit with different performance. Applying *plus₂* to one argument as in *App plus₂(1)* dynamically creates a partial application object $\text{Clos}^2(\lambda(x_{\text{IntR}}, y_{\text{IntR}}). \text{plus}\#(x, y))(1)$, whereas applying *plus₁* to two arguments hits the over-application case:

$$\langle plus_1 \mid \text{App}(1, 2); K \mid H \rangle \mapsto_{\text{OApp}} \langle (\lambda(x_{\text{IntR}}). \text{Clos}^1 \lambda(y_{\text{IntR}}). \text{plus}\#(x, y))(1) \mid \text{App}(2); K \mid H \rangle$$

In other words, it is safe to treat these two expressions interchangeably in our extended \mathcal{ML} .

But notice: if we compile only well-typed \mathcal{IL} expressions, then over application (*OApp*) and under application (*PApp*) will never happen! With the support that over- and under-application give for dynamically checking arity at runtime, the type system of \mathcal{IL} is overly conservative. Programs that are safe to execute in \mathcal{ML} would be rejected in \mathcal{IL} . That's because the actual arity can be checked at runtime, so an arity 2 closure is treated like it has arity 1 (by *PApp*) and vice versa (by *OApp*).

What is missing is a way to say that *plus₁* and *plus₂*, though they are different, are treated identically at runtime: under this extension, both would have their arity dynamically checked before any function is called. There is no harm in using values of one type in any context expecting the other. In more general terms, if we always dynamically check the arity of closures, we can safely equate the following two types without the need for wrapping and unwrapping closures:

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

One way is to capture both of these two operations as a type equality, or a type-safe coercion [Breitner et al. 2016]: the two types are represented identically at runtime as closure objects, with the only difference being the arity count stored inside—which will be checked dynamically.

Correspondingly, consider applying closures of different types to the same argument twice:

$$\begin{aligned}dup_1 : {}^L\{\text{Int}\# \rightsquigarrow {}^L\{\text{Int}\# \rightsquigarrow \text{Int}\#\}\} \rightsquigarrow \text{Int}\# \rightsquigarrow \text{Int}\# & \quad dup_2 : {}^L\{\text{Int}\# \rightsquigarrow \text{Int}\# \rightsquigarrow \text{Int}\#\} \rightsquigarrow \text{Int}\# \rightsquigarrow \text{Int}\# \\ dup_1 f x = \text{App} (\text{App} f x) x & \quad dup_2 f x = \text{App} f x x\end{aligned}$$

These compile to the following \mathcal{ML} functions:

$$\begin{aligned}\mathcal{E}_{\text{Call}[\text{PtrR}, \text{Int}\#]}[\![dup_1]\!] &= \lambda(f_{\text{PtrR}}, x_{\text{IntR}}). \text{App} (\text{App} f_{\text{PtrR}}(x_{\text{IntR}}))(x_{\text{IntR}}) \\ \mathcal{E}_{\text{Call}[\text{PtrR}, \text{Int}\#]}[\![dup_2]\!] &= \lambda(f_{\text{PtrR}}, x_{\text{IntR}}). \text{App} f_{\text{PtrR}}(x_{\text{IntR}}, x_{\text{IntR}})\end{aligned}$$

By checking arities at runtime, both *dup₁* and *dup₂* can be safely applied to either *plus₁* or *plus₂*. This can be necessary when the *plus* operations need to be stored in a list that requires all elements be of the same type. If it happens that *plus₂* is stored with *plus₁*'s type, and then later *dup₂* is applied to it with *dup₁*'s type, then we still get the benefit of passing two arguments at once at runtime, even if our type system cannot ensure it statically.

In the end, the details of dynamically-resolved, higher-arity, function calls can be fully captured by the type system. The arity of a function call is described *exactly* by the type and kind of the function, and the program must provide the right number of arguments and binders. However, as long as there is a top-most lifted abstraction, which corresponds to storing additional arity information at run-time, we can freely convert between functions of different arities. In other words, kinds fully capture calling conventions in programs, as the kind of a boxed closure does not mention any arity and thus must be subject to further checks.

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

Fig. 7. Compiling call-by-name and call-by-value System F to \mathcal{IL}

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. That is indeed possible, by simply having two different closure types (with their introduction and elimination forms). Then, an optimizing compiler, or an expert user, can select the appropriate form for the best performance.

6 COMPILATION TO \mathcal{IL} FROM A HIGHER LEVEL

\mathcal{IL} is, by design, a fairly low-level language that makes fine distinctions about representation, levity and so on. This allows it to act as a target for *both* eager *and* lazy languages. To make this claim concrete, we now give translations for call-by-name System F (Section 6.1) and call-by-value System F (Section 6.2) into \mathcal{IL} .

6.1 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 is capable of capturing all the types of the source language, namely $\star = \text{TYPE PtrR Eval}^L$. Each source-language type τ will be represented by an \mathcal{IL} type $\llbracket \tau \rrbracket$ with kind \star ; that is, a pointer to a lifted value, perhaps a thunk.

Figure 7 gives this type translation. To get the correct call-by-name semantics for numbers, we have to use the boxed integer type Int^L , which happily has the correct kind. However, even though the function type $\llbracket \tau \rrbracket \rightsquigarrow \llbracket \sigma \rrbracket$ has the correct semantics, it has the wrong kind $\text{TYPE PtrR Call}[\text{PtrR}]$. Therefore, the translation coerces the calling convention with a closure type. The only change to polymorphic type abstraction is to say that the type variable t ranges over the “uniform” kind \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 more descriptive kinds.

Call-by-need evaluation is often considered to be the more “practical” version of call-by-name, because it has better asymptotic complexity in many cases but still gives the same result for purely functional programs [Ariola and Felleisen 1997]. We make the same conflation here, where we interpret Eval^L to mean call-by-name in the high-level equational theory (Fig. 4), but actually

Call-by-name substitutable values:

$$V ::= e$$

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

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

$$(\beta_v) \quad (\lambda t.V) \tau = e[\tau/t]$$

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

Call-by-value substitutable values:

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

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

$$(\eta_v) \quad \lambda t.(V t) = V : \forall t.\sigma$$

Fig. 8. Equational theory of System F; call-by-value and call-by-name

implement call-by-need evaluation for Eval^L in the machine (Fig. 5). So to compile a pure call-by-need programming language, like Haskell, one can just apply the call-by-name compilation scheme directly, and get the correct performance characteristics.

A realistic compiler for a call-by-need source language should also interpret the lifted levity L as call-by-need in \mathcal{IL} itself. For example, the call-by-name (β_{\rightarrow}) rule in Figure 4 duplicates the argument S , which might duplicate an arbitrary amount of work. It is well known how to adjust the equational theory to enshrine call-by-need [Ariola and Felleisen 1997], and we do not elaborate here.

6.2 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. 7.

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 wrapped in a closure, as in the call-by-name case, but this time unlifted.

The translation of polymorphism is more complex, however, due to the standard semantics of the call-by-value System F. With call-by-value evaluation, the abstraction $\lambda t.\perp$ is a value, despite that \perp diverges, whereas in call-by-name they would be η -equivalent. We need to make sure that this abstraction is still a value even after the polymorphic λ 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 in expressions is 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 $\lambda t.\text{error } t n$ returns immediately in call-by-value System F, but the corresponding $\lambda t:\star.\text{error PtrR } U t$ 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 $\lambda().e$.

6.3 Correctness of Source-to- \mathcal{IL} Compilation

Compiling call-by-name and call-by-value System F into \mathcal{IL} are both correct, in terms of the preservation of types and the preservation of equalities between expressions from the source to the intermediate language. We assume the standard type system for System F, and review the calculus'

call-by-name and -value axiomatic semantics in Fig. 8. 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 3 (Type Preservation). *If $\Gamma \vdash e : \tau$ is derivable then so is $\llbracket \Gamma \rrbracket \vdash \llbracket e \rrbracket : \llbracket \tau \rrbracket$.*

More interesting is the translation of equalities from the source to the intermediate language, for which we use the axiomatic semantics given in Fig. 4. For more details of the proof, see Appendix C.

Theorem 4 (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 the call-by-name compilation, and likewise for call-by-value.*

7 USER-DEFINED LEVITY-POLYMORPHIC DATA TYPES

Thus far we have had a single built-in boxed data type, `Int`. It is easy to see how to generalize the idea to arbitrary user-defined algebraic data types [Graf 2020]. For example, rather than having `Int` built-in, we might define it like this

$$\begin{aligned} \text{data } \text{Int } (g:\text{Lev}) : \text{TYPE PtrR (Eval } g) \text{ where} \\ \text{I\#} : \forall (g:\text{Lev}). \text{Int\#} \leadsto \text{Int } g \end{aligned}$$

Now, the `Int` type is parameterized by a chosen levity, 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 strictly or lazily evaluated, and whose cells contain elements of any arity or evaluation strategy:

$$\begin{aligned} \text{data } \text{List } (g:\text{Lev}) (v:\text{Conv}) (a:\text{TYPE PtrR } v) : \text{TYPE PtrR (Eval } g) \text{ where} \\ \text{Nil} : \text{List } g \text{ } v \text{ } a \\ \text{Cons} : a \leadsto \text{List } g \text{ } v \text{ } a \leadsto \text{List } g \text{ } v \text{ } a \end{aligned}$$

Despite the 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 of any levity:

$$\begin{aligned} \text{sum} & : \forall (g:\text{Lev}). \text{List } g \text{ (Eval } g) (\text{Int } g) \leadsto \text{Int } g \\ \text{sum Nil} & = \text{I\# } 0 \\ \text{sum (Cons (I\# } x) \text{ xs)} & = \text{case sum xs of I\# } y \rightarrow \text{I\# (plus\# } x \text{ } y) \end{aligned}$$

This polymorphic definition is possible because the *sum* function is completely strict, no matter if it is given an evaluated list or an unevaluated thunk, the entire thing will be added together before a value is returned. Therefore, the same code can be used for either strict or lazy evaluation. Other functions, such as mapping over a list, could not be levity polymorphic in the same way, because the code is very different where a lazy map will delay the recursive call and first return a `Cons` cell, whereas a strict map will first recurse before returning anything.

Because this `List` type is also polymorphic in calling convention, it can be used to contain lists of functions of statically-known arity, as foreshadowed in Section 3.4. For example:

$$\text{Cons plus\# Nil} : \text{List (Call[IntR,IntR]) (Int\#} \leadsto \text{Int\#} \leadsto \text{Int\#})$$

8 RELATED WORK

The system presented in this paper is the culmination and consolidation of several independent lines of work on expressing performance issues directly in an intermediate language. The underlying

theme of each of these topics 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 lifted from unlifted types goes back to [Peyton Jones and Launchbury 1991], and has been used to great effect in GHC for nearly three decades. For most of that time the distinction has been static, but recent work has 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*—indeed, it recognizes no such distinction. Our contribution here is to separate the two completely, and show that one might want to be polymorphic in one but not the other.

One of the main requirements we have followed is to generate only *one* piece of code for every polymorphic definition. This requirement means that there are certain definitions that must be rejected, because the compilation depends on a choice of levity. An alternative approach by [Dunfield 2015] accepts more uses of levity polymorphism, but at the cost of generating *different* code for each choice—an exponential blowup of code size in practice—avoided by our approach.

8.2 Optimizing curried functions

Previous work has established methods for optimizing curried function calls dynamically at runtime, in order to avoid the overhead of naïvely calling $((f\ 1)\ 2)\ 3$ by passing one argument at a time, taking extra steps and allocating spurious closures in between each application. In practice, f will often expect all three arguments before doing any interesting work, so those calls should be fused when possible. This 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 have been able capture this dynamic type of optimization within the syntax and types of programs, as described in Section 5.

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 do those optimizations statically at compile time. Of course, this is easy to do when the compiler can statically find the definition of the called function [Marlow and Peyton Jones 2004]. However, 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 3.4), and polymorphism over calling conventions (Section 3.6). Another difference is that Downen et al. [2019] had two function arrows, $(\tau \rightarrow \sigma)$ and $(\tau \rightsquigarrow \sigma)$, whereas \mathcal{IL} has just one arrow $(\tau \rightsquigarrow \sigma)$, plus the closure type $^Y\{\tau\}$. The two are inter-convertible: $\tau \rightarrow \sigma \approx \{\tau \rightsquigarrow \sigma\}$ and $\{\tau\} \approx () \rightarrow \tau$. The approach here has a greater economy of concepts, and a nice correspondence with $\text{Int}\#$ and Int^Y . 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 levities. Indeed GHC goes further: rather than having a stratified zoo of different sorts of things (types,

kinds, levities, representations, calling conventions) as we do in Fig. 1, they are all *types* [Weirich et al. 2013] kept separate by their *kinds*. This is a fantastic simplification, and immediately allows polymorphism over all these conceptually-different things. This does mean, however, that in GHC it is hard *not* to have polymorphism! Returning to Section 3.4, it would be hard to prevent a forall-quantified type variable from being instantiated by an arrow type, because you would have to say something like “this quantified variable can have any kind *other than* *Call*”. So GHC’s infrastructure strongly encourages fully-fledged polymorphism, just as we present in this paper.

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.

One of the connections that previous work on unboxed types and extensional functions has noticed is that *lifting*—in the sense of denotational semantics—corresponds to a *mismatch* between machine primitives and the semantics of a programming language. The unlifted versions of types, like integers, can be implemented more directly—and therefore more efficiently—in a machine. But preventing lifting is different depending on the type: unlifted integers need to be call-by-value whereas unlifted function chains need to be call-by-name. The first way to reconcile this tension 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 that foundational line of work.

The same correspondence between types and evaluation also corresponds to *polarity* developed for focused proof search [Andreoli 1992; Laurent 2002]. Polarity and focusing has been found to correspond to functional programming [Zeilberger 2008, 2009], as well as semantics and computation [Munch-Maccagnoni 2009, 2013]. More recently, these mixed evaluation strategy languages have been extended with more 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 unknown function type constructor to a “down shift” ($\{\tau \rightsquigarrow \sigma\} = \Downarrow(\tau \rightsquigarrow \sigma)$). Note that for the sake of usability, \mathcal{IL} performs implicit polarity conversions of types based on their context. For example, closing over a non-function type like 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*, *calling convention*, 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 directly capture performance in the intermediate language, we would also like to be able to characterize the cost of computation in its semantics as well.

ACKNOWLEDGMENTS

This material is based upon work supported by the National Science Foundation under Grant No. nnnnnnn and Grant No. mmmmmmm. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author and do not necessarily reflect the views of the National Science Foundation.

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>
- Arvind and Kattamuri Ekanadham. 1988. Future Scientific Programming on Parallel Machines. *J. Parallel and Distrib. Comput.* 5, 5 (1988), 460–493.
- 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, 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.

- 1275 Guillaume Munch-Maccagnoni. 2013. *Syntax and Models of a non-Associative Composition of Programs and Proofs*. Ph.D.
1276 Dissertation. Université Paris Diderot.
- 1277 Simon L. Peyton Jones. 1992. Implementing Lazy Functional Languages on Stock Hardware: The Spineless Tagless G-machine.
1278 *Journal of Functional Programming* 2, 2 (1992), 127–202.
- 1279 Simon L. Peyton Jones and John Launchbury. 1991. Unboxed Values As First Class Citizens in a Non-Strict Functional
1280 Language. In *Proceedings of the 5th ACM Conference on Functional Programming Languages and Computer Architecture*.
Springer-Verlag, London, UK, UK, 636–666.
- 1281 Amr Sabry and Matthias Felleisen. 1993. Reasoning About Programs in Continuation-Passing Style. *Lisp and Symbolic*
1282 *Computation* 6, 3-4 (Nov. 1993), 289–360.
- 1283 Alex Theriault. 2019. GHC Proposal 209: levity-polymorphic Lift. [https://github.com/ghc-proposals/ghc-proposals/blob/
1284 master/proposals/0209-levity-polymorphic-lift.rst](https://github.com/ghc-proposals/ghc-proposals/blob/master/proposals/0209-levity-polymorphic-lift.rst)
- 1285 Philip Wadler, Walid Taha, and David MacQueen. 1998. How to add laziness to a strict language without even being odd. In
1286 *Proceedings of the Standard ML Workshop*.
- 1287 Stephanie Weirich, Justin Hsu, and Richard A. Eisenberg. 2013. System FC with Explicit Kind Equality. In *International*
1288 *Conference on Functional Programming (ICFP '13)*. ACM.
- 1289 Noam Zeilberger. 2008. On the Unity of Duality. *Annals of Pure and Applied Logic* 153, 1 (2008), 660–96.
- 1290 Noam Zeilberger. 2009. *The Logical Basis of Evaluation Order and Pattern-Matching*. Ph.D. Dissertation. Carnegie Mellon
1291 University.
- 1292
- 1293
- 1294
- 1295
- 1296
- 1297
- 1298
- 1299
- 1300
- 1301
- 1302
- 1303
- 1304
- 1305
- 1306
- 1307
- 1308
- 1309
- 1310
- 1311
- 1312
- 1313
- 1314
- 1315
- 1316
- 1317
- 1318
- 1319
- 1320
- 1321
- 1322
- 1323

A OPERATIONAL SEMANTICS FOR \mathcal{IL}

A.1 Notation

We use \mapsto to denote the single step relation of an operational semantics, $\mapsto^?$ to denote its reflexive closure, \mapsto^+ to denote its transitive closure, and \mapsto^* to denote its reflexive-transitive closure. This same notation is used for other reduction arrows.

A.2 Substitution-based (call-by-name) operational semantics

$$A \in \text{Answer} ::= V \mid E[\text{error } \rho \ \gamma \ \tau \ V]$$

$$V \in \text{Value} ::= x \ \bar{\phi} \mid c \ \bar{\phi} \mid \text{I}\#^Y V \mid \text{Clos}^Y e \mid \lambda x:\tau.e \mid \lambda \chi.V$$

$$E \in \text{EvalCxt} ::= \square \mid e \ E \mid E \ e \mid \text{App } E \mid \text{I}\#^Y E \mid \text{case } E \text{ of } \text{I}\# \ x \rightarrow e \mid E \ \phi \mid \lambda \chi.E$$

Note that the definition of *value* only differs from *passive* in that we don't allow for β_V redexes. In other words, every passive expression reduces (by the following operational semantics) to some value: for all P there is a V such that $P \mapsto_{\beta_V}^* V$. A substitutable value can be bound to a variable, and is determined in part by the type of the expression as given by the following rules:

$$\frac{}{\Gamma \vdash V \ \mathbf{val}} \quad \frac{\Gamma \vdash e : \tau \quad \Gamma \vdash \tau : \text{TYPE } \rho \ \text{Call}[\alpha]}{\Gamma \vdash e \ \mathbf{val}} \quad \frac{\Gamma \vdash e : \tau \quad \Gamma \vdash \tau : \text{TYPE } \rho \ \text{Eval}^L}{\Gamma \vdash e \ \mathbf{val}}$$

The decomposition of an expression into an evaluation context surrounding an expression is also defined in part by typing information, written with the judgement $\Gamma \vdash E @ e : \tau \xRightarrow{\Gamma'} \sigma$, by the following inference rules:

$$\begin{array}{c} \frac{\Gamma \vdash e : \tau}{\Gamma \vdash \square @ e : \tau \xRightarrow{\Gamma'} \tau} \quad \frac{\Gamma \vdash E @ e : \tau \xRightarrow{\Gamma'} \text{Int}\#}{\Gamma \vdash \text{I}\#^Y E @ e : \tau \xRightarrow{\Gamma'} \text{Int}^Y} \quad \frac{\Gamma \vdash E @ e : \tau \xRightarrow{\Gamma'} \text{Int}^Y \quad \Gamma, x : \text{Int}\# \vdash e' : \sigma}{\Gamma \vdash \text{case } E \text{ of } \text{I}\# \ x \rightarrow e' @ e : \tau \xRightarrow{\Gamma'} \sigma} \\[10pt] \frac{\Gamma \vdash E @ e : \tau \xRightarrow{\Gamma'} \forall\{\sigma\}}{\Gamma \vdash \text{App } E @ e : \tau \xRightarrow{\Gamma'} \sigma} \quad \frac{\Gamma \vdash E @ e : \tau \xRightarrow{\Gamma'} \forall \chi. \sigma}{\Gamma \vdash E \ \phi @ e : \tau \xRightarrow{\Gamma'} \sigma[\phi/\chi]} \quad \frac{\Gamma, \chi \vdash E @ e : \tau \xRightarrow{\Gamma'} \sigma}{\Gamma \vdash \lambda \chi. E @ e : \tau \xRightarrow{\Gamma', \chi} \forall \chi. \sigma} \\[10pt] \frac{\Gamma \vdash E @ e : \tau \xRightarrow{\Gamma'} \tau' \quad \Gamma \vdash \tau' : \text{TYPE } \rho \ \text{Eval}^U \quad \Gamma \vdash e : \tau' \rightsquigarrow \sigma \quad \vdash \rho \ \mathbf{rep}}{\Gamma \vdash e' E @ e : \tau \xRightarrow{\Gamma'} \sigma} \\[10pt] \frac{\Gamma \vdash E @ e : \tau \xRightarrow{\Gamma'} \tau' \rightsquigarrow \sigma \quad \Gamma \vdash e' : \tau' \quad \Gamma \vdash e' \ \mathbf{val} \quad \Gamma \vdash \tau' \ \mathbf{mono-rep} \quad \Gamma \vdash \tau' \ \mathbf{mono-conv}}{\Gamma \vdash E e' @ e : \tau \xRightarrow{\Gamma'} \sigma} \end{array}$$

The (typed) operational stepping relation is written $\Gamma \vdash e \mapsto e' : \sigma$, and always assumes as a precondition that $\Gamma \vdash e : \sigma$ holds. The primary reduction steps are by the following rules (where we omit the typing precondition on the left-hand side):

$$\begin{array}{ll} (\beta_{\lambda}) & \Gamma \vdash (\lambda x:\tau.e) \ e' \mapsto e[e'/x] : \sigma \quad (\text{if } \Gamma \vdash e' \ \mathbf{val}) \\ (\beta_V) & \Gamma \vdash (\lambda \chi.V) \ \phi \mapsto V[\phi/\chi] : \sigma[\phi/\chi] \\ (\beta_{\{\}}) & \Gamma \vdash \text{App } (\text{Clos}^Y e) \mapsto e : \sigma \\ (\beta_{\text{Int}}) & \Gamma \vdash \text{case } \text{I}\#^Y V \text{ of } \text{I}\# \ x \rightarrow e \mapsto e[V/x] : \sigma \end{array}$$

Additionally, the operational steps are closed under evaluation contexts:

$$\frac{\Gamma \vdash E @ e : \tau \xRightarrow{\Gamma'} \sigma \quad \Gamma, \Gamma' \vdash e \mapsto e' : \tau}{\Gamma \vdash E[e] \mapsto E[e'] : \sigma} \text{ compat}$$

Lemma 1 (Typed Decomposition). *If $\Gamma \vdash E @ e : \tau \xRightarrow{\Gamma'} \sigma$ then $\Gamma, \Gamma' \vdash e : \tau$ and $\Gamma \vdash E[e] : \sigma$.*

PROOF. By induction on the derivation of the decomposition $\Gamma \vdash E @ e : \tau \xRightarrow{\Gamma'} \sigma$. \square

Lemma 2 (Unique Decomposition). *For every \mathcal{IL} expression $\Gamma \vdash e : \sigma$, either:*

- (1) $e \in \text{Value}$, or
- (2) there is a unique $\Gamma \vdash E @ e' : \tau \xRightarrow{\Gamma'} \sigma$ such that either
 - (a) e' is a variable or constant, or
 - (b) $\Gamma \vdash e' \mapsto e'' : \tau$ directly (i.e., not by compat).

PROOF. By induction on the typing derivation of $\Gamma \vdash e : \sigma$. \square

Corollary 1 (Determinism). *If $\Gamma \vdash e \mapsto e_1 : \sigma$ and $\Gamma \vdash e \mapsto e_2 : \sigma$ then $e_1 =_\alpha e_2$.*

PROOF. Follows directly from Lemma 2 and the fact that the operational rules don't overlap. \square

Lemma 3 (Stability under Substitution).

- (1) If $\Gamma \vdash e'' : \sigma$ and $\Gamma \vdash e'' \text{ val}$ then $\Gamma, x : \sigma \vdash e \mapsto e' : \tau$ implies $\Gamma \vdash e[e''/x] \mapsto e'[e''/x] : \sigma$.
- (2) If $\Gamma \vdash [\phi/\chi] \text{ poly}$ then $\Gamma, \chi \vdash e \mapsto e' : \tau$ if and only if $\Gamma[\phi/\chi] \vdash e[\phi/\chi] \mapsto e'[\phi/\chi] : \tau[\phi/\chi]$.

Similarly, the decomposition of evaluation contexts E is stable under substitution.

PROOF. By induction on the derivation of the reduction, or the syntax of the context. \square

Lemma 4 (Typed Substitution). (1) If $\Gamma, x : \sigma \vdash e : \tau$ and $\Gamma \vdash e' : \sigma$ then $\Gamma \vdash e[e'/x] : \tau$.

- (2) If $\Gamma, t : \kappa \vdash e : \tau$ and $\Gamma \vdash \sigma : \kappa$ then $\Gamma \vdash e[\sigma/t] : \tau[\sigma/t]$.
- (3) If $\Gamma, g \vdash e : \tau$ and $\Gamma \vdash \gamma \text{ lev}$ then $\Gamma[\gamma/g] \vdash e[\gamma/g] : \tau[\gamma/g]$.
- (4) If $\Gamma, r \vdash e : \tau$ and $\Gamma \vdash \rho \text{ rep}$ then $\Gamma[\rho/r] \vdash e[\rho/r] : \tau[\rho/r]$.
- (5) If $\Gamma, v \vdash e : \tau$ and $\Gamma \vdash v \text{ conv}$ then $\Gamma[v/v] \vdash e[v/v] : \tau[v/v]$.

PROOF. By induction on the given typing derivation of e . \square

Lemma 5 (Progress). *If $\vdash e : \tau$ then either:*

- (1) $e \in \text{Answer}$, or
- (2) $\vdash e \mapsto e' : \tau$.

PROOF. Follows directly from Lemma 2 and compatibility. \square

Lemma 6 (Preservation). *If $\Gamma \vdash e \mapsto e' : \tau$ then $\Gamma \vdash e' : \tau$.*

PROOF. Follows by induction on the decomposition of expressions into evaluation contexts and redexes and cases on the operational rules, using the fact that well-typed substitution preserves typing (Lemma 4). \square

A.3 Environment-based (call-by-need) operational semantics

The primary difference between the call-by-need operational semantics—using environments to handle variables—versus the call-by-name operational semantics—using substitution instead—is *sharing*: when a lazy expression is bound in call-by-need, it is computed *at most once*. In order to make this distinction more syntactically apparent, here we officially at **let**-expressions to the language, with the typing rules given by the encoding

$$\mathbf{let} x:\tau = e \mathbf{in} e' \triangleq (\lambda x:\tau. e') e$$

The refined definition of answers, evaluation contexts, *etc.* is given as follows:

$$\begin{aligned} V \in \mathbf{Value} &::= c \bar{\phi} \mid R \mid \lambda \chi. V \\ R \in \mathbf{Reference} &::= \mathbf{I}^\gamma a \mid \mathbf{Clos}^\gamma e \mid \lambda x:\tau. e \mid \lambda \chi. R \\ a \in \mathbf{Arg} &::= c \bar{\phi} \mid x \bar{\phi} \mid \lambda \chi. a \\ A \in \mathbf{Answer} &::= H[V] \mid E[\mathbf{error} \ \rho \ \gamma \ \tau \ a] \\ E \in \mathbf{EvalCxt} &::= \square \mid F[E] \mid B[E] \\ H \in \mathbf{ClosCxt} &::= \square \mid B[H] \\ F \in \mathbf{FrameCxt} &::= \square a \mid \mathbf{App} \ \square \mid \mathbf{case} \ \square \mathbf{of} \ \mathbf{I}^\# x \rightarrow e \mid \mathbf{let} x:\tau = \square \mathbf{in} e \mid \square \phi \\ B \in \mathbf{BindCxt} &::= \mathbf{let} x:\tau = e \mathbf{in} \square \mid \lambda \chi. \square \end{aligned}$$

Notice that to support call-by-need evaluation, we further distinguish *references* to values, which may be copied at will, versus *bindable* expressions, which are merely allocated. Intuitively, referenced values may be copied, whereas bound expressions must be fully evaluated first before they can be copied.

$$\begin{array}{c} \frac{}{\Gamma \vdash R \mathbf{ref}} \quad \frac{\Gamma \vdash e : \tau \quad \Gamma \vdash \tau : \mathbf{TYPE} \ \mathbf{PtrR} \ \mathbf{Call}[\alpha]}{\Gamma \vdash e \mathbf{ref}} \\ \frac{\Gamma \vdash e \mathbf{ref}}{\Gamma \vdash e \mathbf{bind}} \quad \frac{\Gamma \vdash e : \tau \quad \Gamma \vdash \tau : \mathbf{TYPE} \ \mathbf{PtrR} \ \mathbf{Eval}^L}{\Gamma \vdash e \mathbf{bind}} \end{array}$$

The type-driven decomposition of expressions into evaluation contexts and their sub-expression is the same for all syntactic forms except for applications and **let**-expressions (which are new for call-by-need). Decomposition for these cases are defined as:

$$\begin{array}{c} \frac{\Gamma \vdash E @ e : \tau \xRightarrow{\Gamma'} \tau' \rightsquigarrow \sigma \quad \Gamma \vdash a : \tau'}{\Gamma \vdash E a @ e : \tau \xRightarrow{\Gamma'} \sigma} \\ \frac{\Gamma, x : \tau' \vdash E @ e : \tau \xRightarrow{\Gamma'} \sigma \quad \Gamma \vdash e' : \tau' \quad \Gamma \vdash e' \mathbf{bind} \quad \Gamma \vdash \tau' \mathbf{mono-rep}}{\Gamma \vdash \mathbf{let} x:\tau' = e' \mathbf{in} E @ e : \tau \xRightarrow{\Gamma', x:\tau'} \sigma} \\ \frac{\Gamma \vdash E @ e : \tau \xRightarrow{\Gamma'} \tau' \quad \Gamma, x : \tau' \vdash e' : \sigma \quad \Gamma \vdash \tau' : \mathbf{TYPE} \ \rho \ \mathbf{Eval}^U \quad \Gamma \vdash \tau' \mathbf{mono-rep}}{\Gamma \vdash \mathbf{let} x:\tau' = E \mathbf{in} e' @ e : \tau \xRightarrow{\Gamma'} \sigma} \\ \frac{\Gamma, x : \tau' \vdash E_1 @ x : \tau' \xRightarrow{\Gamma_1} \sigma \quad \Gamma \vdash E_2 @ e : \tau \xRightarrow{\Gamma_2} \tau' \quad \Gamma \vdash \tau' : \mathbf{TYPE} \ \mathbf{PtrR} \ \mathbf{Eval}^L}{\Gamma \vdash \mathbf{let} x:\tau' = E_2 \mathbf{in} E_1[x] @ e : \tau \xRightarrow{\Gamma_2} \sigma} \end{array}$$

As with decomposing an evaluation context, we also need to confirm that answers are properly decomposed into the appropriate context surrounding either a value or an error:

$$\frac{\Gamma \vdash H @ V : \tau \xRightarrow{\Gamma'} \sigma}{\Gamma \vdash H[V] \text{ ans}} \quad \frac{\Gamma \vdash E @ \text{error } \rho \gamma \tau a : \tau \xRightarrow{\Gamma'} \sigma}{\Gamma \vdash E[\text{error } \rho \gamma \tau a] \text{ ans}}$$

The reduction rules for call-by-need are also updated from the call-by-name semantics. Of note, β_{\sim} is much more restricted to only substituting arguments which (after type erasure) are variables or constants. The analogous substitution is repeated for renaming in a **let**-expression.

$$\begin{aligned} (\beta_{\sim}) \quad & \Gamma \vdash (\lambda x:\tau. e) a \mapsto e[a/x] : \sigma \\ (\beta_{\forall}) \quad & \Gamma \vdash (\lambda \chi. A) \phi \mapsto A[\phi/\chi] : \sigma[\phi/\chi] \\ (\beta_{\{\}}) \quad & \Gamma \vdash \text{App } (\text{Clos}^Y e) \mapsto e : \sigma \\ (\beta_{\text{Int}}) \quad & \Gamma \vdash \text{case } I\#^Y a \text{ of } I\# x \rightarrow e \mapsto e[a/x] : \sigma \\ (\text{rename}) \quad & \Gamma \vdash \text{let } x:\tau = a \text{ in } e \mapsto e[a/x] : \sigma \\ (\text{name}_{I\#}) \quad & \Gamma \vdash I\#^Y e \mapsto \text{let } x:\text{Int}\# = e \text{ in } I\#^Y x : \text{Int}^Y \quad (\text{if } e \notin \text{Arg}) \\ (\text{name}) \quad & \Gamma \vdash e' e \mapsto \text{let } x:\tau = e \text{ in } e' x : \sigma \quad (\text{if } \Gamma \vdash e : \tau \text{ and } e \notin \text{Arg}) \end{aligned}$$

Because the grammar of evaluation context is simplified to be primarily based on **lets**—rather than chains of curried function application—we convert more complex applications to an alternative **let** form, as done by the *name* and *name_{I#}* rules. These **let**-expressions are interpreted by variable lookup (which inlines the definition of a **let** only when its needed in an evaluation context, and only when the definition is a copyable reference). The presence of delayed **let** bindings also necessitates administrative *commuting* conversions, which push bindings out of the way to bring frames of an evaluation context closer to the root of an answer.

$$\begin{aligned} & \frac{\Gamma \vdash E @ x : \tau \xRightarrow{\Gamma'} \sigma \quad x \notin \text{dom}(\Gamma) \quad \Gamma \vdash e \text{ ref}}{\Gamma \vdash \text{let } x:\tau = e \text{ in } E[x] \mapsto \text{let } x:\tau = e \text{ in } E[e] : \sigma} \text{ look} \\ & \frac{\Gamma \vdash F @ \text{let } x:\tau' = e \text{ in } A : \tau \xRightarrow{\Gamma'} \sigma \quad \Gamma \vdash \text{let } x:\tau' = e \text{ in } A \text{ ans}}{\Gamma \vdash F[\text{let } x:\tau' = e \text{ in } A] \mapsto \text{let } x:\tau' = e \text{ in } F[A] : \sigma} \text{ comm} \\ & \frac{\Gamma \vdash E @ e : \tau \xRightarrow{\Gamma'} \sigma \quad \Gamma, \Gamma' \vdash e \mapsto e' : \tau}{\Gamma \vdash E[e] \mapsto E[e'] : \sigma} \text{ compat} \end{aligned}$$

Lemma 7 (Typed Decomposition). *If $\Gamma \vdash E @ e : \tau \xRightarrow{\Gamma'} \sigma$ then $\Gamma, \Gamma' \vdash e : \tau$ and $\Gamma \vdash E[e] : \sigma$.*

PROOF. By induction on the derivation of the decomposition $\Gamma \vdash E @ e : \tau \xRightarrow{\Gamma'} \sigma$. \square

Lemma 8 (Unique Decomposition). *For every \mathcal{IL} expression $\Gamma \vdash e : \sigma$, either:*

- (1) $e \in \text{Answer}$, or
- (2) there is a unique $\Gamma \vdash E @ e' : \tau \xRightarrow{\Gamma'} \sigma$ such that either
 - (a) e' is a variable or constant, or
 - (b) $\Gamma \vdash e' \mapsto e'' : \tau$ directly (i.e., not by *compat*).

PROOF. By induction on the typing derivation of $\Gamma \vdash e : \sigma$. \square

Corollary 2 (Determinism). *If $\Gamma \vdash e \mapsto e_1 : \sigma$ and $\Gamma \vdash e \mapsto e_2 : \sigma$ then $e_1 =_{\alpha} e_2$.*

PROOF. Follows directly from Lemma 8 and the fact that the operational rules don't overlap. \square

Lemma 9 (Progress). *If $\Gamma \vdash e : \tau$ then either:*

- (1) $e \in \text{Answer}$, or
- (2) $\Gamma \vdash e \mapsto e' : \tau$.

PROOF. Follows directly from Lemma 8 and compatibility. \square

Lemma 10 (Preservation). *If $\Gamma \vdash e \mapsto e' : \tau$ then $\Gamma \vdash e' : \tau$.*

PROOF. Follows by induction on the decomposition of expressions into evaluation contexts and redexes and cases on the operational rules. \square

A.4 Bisimulation between call-by-name and call-by-need

Definition 1 (Unwinding Simulation). The base simulation relation between the call-by-need and call-by-name operational semantics of \mathcal{IL} , written $\Gamma \vdash_{\mathcal{U}} e \sim e' : \tau$ where e may contain let-expressions but e' cannot, is defined inductively by the following rules for unwinding lets:

$$\frac{\Gamma \vdash_{\mathcal{U}} e_e \sim e_s : \tau \quad \Gamma, \overline{y_i} : \tau \vdash_{\mathcal{U}} e'_e \sim e'_s : \sigma \quad \Gamma \vdash e_s \mathbf{val} \quad \Gamma \vdash e_{si} \mapsto^* e_s : \tau}{\Gamma \vdash_{\mathcal{U}} \mathbf{let} x : \tau = e_e \mathbf{in} e'_e[\overline{x/y_i}] \sim e'_s[\overline{e_{si}/y_i}] : \sigma} \text{copy}$$

$$\frac{\Gamma \vdash_{\mathcal{U}} e_e \sim e_s : \tau \quad \Gamma, x : \tau \vdash_{\mathcal{U}} e'_e \sim e'_s : \sigma}{\Gamma \vdash_{\mathcal{U}} \mathbf{let} x : \tau = e_e \mathbf{in} e'_e \sim (\lambda x : \tau. e'_e) e_s : \sigma} \text{share}$$

plus rules for compatibility with all other syntactic forms. The full simulation is then defined by inlining some unneeded names in the call-by-name semantics:

$$\frac{\Gamma \vdash_{\mathcal{U}} e \sim e'' : \tau \quad \Gamma \vdash e'' \mapsto_{\text{unname}, \text{unname}_{\text{I\#}}}^* e' : \tau}{\Gamma \vdash_{\mathcal{N}} e \sim e' : \tau}$$

$$\begin{aligned} (\text{unname}) \quad & \Gamma \vdash (\lambda x : \tau. e' x) e \mapsto e' e : \sigma \quad (\text{if } x \notin \text{FV}(e')) \\ (\text{unname}_{\text{I\#}}) \quad & \Gamma \vdash (\lambda x : \text{Int\#}. \text{I\#}^Y x) e \mapsto \text{I\#}^Y e : \text{Int}^Y \end{aligned}$$

Lemma 11 (Context Unwinding). (1) *If $\Gamma \vdash E_e @ e_e : \tau \xRightarrow{\Gamma', \Gamma''} \sigma$ in the call-by-need \mathcal{IL} , $\Gamma \vdash_{\mathcal{U}} E_e[e_e] \sim e''_s : \sigma$, and $\Gamma \not\vdash e_e \mathbf{val}$, then $\Gamma \vdash e''_s \mapsto^* E_s[e_s[\theta]] : \tau$ for some evaluation context E_s and substitution $\theta = [\overline{e/x}]$ in the call-by-name \mathcal{IL} such that $\Gamma \vdash E_s @ e_s[\theta] : \tau \xRightarrow{\Gamma'} \sigma$ and $\Gamma, \Gamma', \Gamma'' \vdash_{\mathcal{U}} e_e \sim e_s : \tau$. Furthermore, if $\Gamma, \Gamma', \Gamma'' \vdash e_e \mapsto^* e'_e : \tau$ and $\Gamma, \Gamma', \Gamma'' \vdash e_s \mapsto^* e'_s : \tau$ such that $\Gamma, \Gamma', \Gamma'' \vdash e'_e \sim e'_s : \tau$, then $\Gamma E_e[e'_e] \sim E_s[e'_s[\theta]] : \tau$.*

(2) *If $\Gamma \vdash H_e @ e_e : \tau \xRightarrow{\Gamma'} \sigma$ and $\Gamma \vdash e_e \mathbf{val}$ in the call-by-need \mathcal{IL} and $\Gamma \vdash_{\mathcal{U}} H_e[e_e] \sim e'_s : \sigma$, then $\Gamma \vdash e'_s \mapsto^* e_s[\theta] : \tau$ for some substitution $\theta = [\overline{e/x}]$ in the call-by-name \mathcal{IL} such that $\Gamma, \Gamma' \vdash_{\mathcal{U}} e_e \sim e_s : \tau$.*

PROOF. By simultaneous induction on the evaluation context decomposition and the unwinding simulation (note that the part (2) follows analogously to part (1) for bound lets):

- The cases for all non-let evaluation contexts follow from the inductive hypothesis.
- $\Gamma \vdash \mathbf{let} x : \tau' = e'_e \mathbf{in} E_e @ e_e : \tau \xRightarrow{\Gamma, \Gamma'', x : \tau'} \sigma$ because $\Gamma, x : \tau' \vdash E_e @ e_e : \tau \xRightarrow{\Gamma', \Gamma''} \sigma$ and $\Gamma \vdash e'_e \mathbf{bind}$. This case depends on the rule of the simulation:
 - *share*: We have

$$\frac{\Gamma, x : \tau' \vdash_{\mathcal{U}} E_e[e_e] \sim e_s : \sigma \quad \Gamma \vdash_{\mathcal{U}} e'_e \sim e'_s : \tau'}{\Gamma \vdash_{\mathcal{U}} \mathbf{let} x : \tau' = e'_e \mathbf{in} E_e[e_e] \sim (\lambda x : \tau'. e_s) e'_s : \sigma}$$

By the inductive hypothesis, we get that

$$\Gamma \vdash e_s \mapsto^* E_{s1}[e_{s1}[\theta]] : \sigma \quad \Gamma, x : \tau' \vdash E_s @ e_{s1}[\theta] : \tau \xRightarrow{\Gamma'} \sigma \quad \Gamma, x : \tau', \Gamma', \Gamma'' \vdash_{\mathcal{U}} e_e \sim e_{s1} : \tau$$

In other words, we reduce to

$$(\lambda x : \tau'. e_s) e'_s \mapsto_{\beta_{\sim}} e_s[e'_s/x] \mapsto^* E_s[e'_s/x][e_{s1}[\theta, e'_s/x]]$$

which gives our new evaluation context and substitution, since evaluation contexts are stable under substitution (Lemma 3). Any further reduction of e_e and e_{s1} preserve the relation by the *copy* rule.

– *copy*: We have $E_{e0}[e_{e0}][x/y_i] = E_e[e_e]$ such that

$$\frac{\Gamma, \overline{y_i : \tau'} \vdash_{\mathcal{U}} E_{e0}[e_{e0}] \sim e_s : \sigma \quad \Gamma \vdash_{\mathcal{U}} e'_e \sim e'_s : \tau' \quad \Gamma \vdash e'_{si} \mapsto e'_s : \tau'}{\Gamma \vdash_{\mathcal{U}} \text{let } x : \tau' = e'_e \text{ in } E_{e0}[e_{e0}][x/y_i] \sim e_s[e'_{si}/y_i] : \sigma}$$

By the inductive hypothesis, we get that

$$\Gamma, \overline{y_i : \tau'} \vdash e_s \mapsto^* E_{s0}[e_{s0}[\theta]] : \sigma \quad \Gamma, \overline{y_i : \tau'} \vdash E_{s0} @ e_{s0}[\theta] : \tau \xRightarrow{\Gamma'} \sigma \quad \Gamma, \overline{y_i : \tau'} \vdash e_{e0} \sim e_{s0} : \sigma$$

In other words, by stability of call-by-name reduction under substitution (Lemma 3), we reduce to $e_s[e'_{si}/y_i] \mapsto^* E_{s0}[e_{s0}[\theta, e'_{si}/y_i]]$. Any further reduction of e_e and e_{s0} preserve the relation by the *copy* rule.

- $\Gamma \vdash \text{let } x : \tau' = E_e \text{ in } e'_e @ e_e : \tau \xRightarrow{\Gamma', \Gamma''} \sigma$ because $\Gamma \vdash E_e @ e_e : \tau \xRightarrow{\Gamma', \Gamma''} \tau'$ and $\Gamma \vdash \tau' : \text{TYPE } \rho \text{ Eval}^{\text{U}}$. Only the *share* rule applies, so we have that

$$\frac{\Gamma, x : \tau' \vdash_{\mathcal{U}} e'_e \sim e'_s : \sigma \quad \Gamma \vdash_{\mathcal{U}} E_e[e_e] \sim e_s : \tau'}{\Gamma \vdash_{\mathcal{U}} \text{let } x : \tau' = E_e[e_e] \text{ in } e'_e \sim (\lambda x : \tau'. e'_s) e_s : \sigma}$$

By the inductive hypothesis, we get that

$$\Gamma \vdash e_s \mapsto^* E_{s0}[e_{s0}[\theta]] : \tau' \quad \Gamma \vdash E_{s0} @ e_{s0}[\theta] : \tau \xRightarrow{\Gamma'} \sigma \quad \Gamma, \Gamma', \Gamma'' \vdash_{\mathcal{U}} e_e \sim e_{s0} : \tau$$

In other words, we reduce to $(\lambda x : \tau'. e'_s) e_s \mapsto^* (\lambda x : \tau'. e'_s) E_{s0}[e_{s0}[\theta]]$. Any further reduction of e_e and e_{s0} preserve the relation by the *share* rule.

- $\Gamma \vdash \text{let } x : \tau' = E_{e2} \text{ in } E_{e1}[x] @ e_e : \tau \xRightarrow{\Gamma_2, \Gamma''} \sigma$ because $\Gamma \vdash E_{e2} @ e_e : \tau \xRightarrow{\Gamma_2, \Gamma''} \tau'$ and $\Gamma \vdash E_{1e} @ x : \tau' \xRightarrow{\Gamma_2, \Gamma'', x : \tau'} \sigma$ and $\Gamma \vdash \tau' : \text{TYPE } \rho \text{ Eval}[\text{L}]$. This case depends on the rule of the simulation:

– *share*: We have

$$\frac{\Gamma, x : \tau' \vdash_{\mathcal{U}} E_{e1}[x] \sim e_s : \sigma \quad \Gamma \vdash_{\mathcal{U}} E_{e2}[e_e] \sim e'_s : \tau'}{\Gamma \vdash_{\mathcal{U}} \text{let } x : \tau' = E_{e2}[e_e] \text{ in } E_{e1}[x] \sim (\lambda x : \tau'. e'_s) e_s : \sigma}$$

By the inductive hypothesis, we get that

$$\begin{array}{ll} \Gamma, x : \tau' \vdash e_s \mapsto^* E_{s1}[e_{s1}[\theta]] : \sigma & \Gamma \vdash e'_s \mapsto^* E_{s2}[e''_s[\theta]] : \tau' \\ \Gamma, x : \tau' \vdash E_{s1} @ e_{s1}[\theta] : \tau' \xRightarrow{\Gamma_1} \sigma & \Gamma \vdash E_{s2} @ e''_s[\theta] : \tau \xRightarrow{\Gamma_2} \tau' \\ \Gamma, x : \tau' \vdash_{\mathcal{U}} x \sim e_{s1} : \tau' & \Gamma, \Gamma_2 \vdash_{\mathcal{U}} e_e \sim e''_s : \tau \end{array}$$

It follows that e_{s1} must be x and $x \notin \text{dom}(\theta)$, so because reduction is stable under substitution (Lemma 3) we reduce to

$$(\lambda x : \tau'. e_s) e'_s \mapsto_{\beta_{\sim}} e_s[e'_s/x] \mapsto^* E_{s1}[e'_s] \mapsto^* E_{s1}[E_{s2}[e''_s[\theta]]]$$

Any further reduction of e_e and e''_s preserve the relation by the *copy* rule, notably expanding the other instances of e'_s .

– *copy*: We have $E'_{e_1}[\overline{x/y_i}] = E_{e_1}$ and (without loss of generality) $x \notin \text{FV}(E_{e_1})$ such that

$$\frac{\Gamma, x : \tau', \overline{y_i : \tau'} \vdash_{\mathcal{U}} E_{e_1}[x] \sim e_s : \sigma \quad \Gamma \vdash_{\mathcal{U}} E_{e_2}[e_e] \sim e'_s : \tau' \quad \Gamma \vdash e'_{s_i} \mapsto^* e'_s : \tau'}{\Gamma \vdash_{\mathcal{U}} \text{let } x : \tau' = E_{e_2}[e_e] \text{ in } E_{e_1}[x] \sim e_s[e'_{s_0}/x, \overline{e'_{s_i}/y_i}] : \sigma}$$

By the inductive hypothesis, we get that

$$\begin{aligned} \Gamma, x : \tau', \overline{y_i : \tau'} \vdash e_s \mapsto^* E_{s_1}[e_{s_1}[\theta]] : \sigma & \quad \Gamma \vdash e'_s \mapsto^* E_{s_2}[e'_s[\theta]] : \tau' \\ \Gamma, x : \tau', \overline{y_i : \tau'} \vdash E_{s_1} @ e_{s_1}[\theta] : \tau' \xRightarrow{\Gamma_1} \sigma & \quad \Gamma \vdash E_{s_2} @ e'_s[\theta] : \tau \xRightarrow{\Gamma_2} \tau' \\ \Gamma, x : \tau', \overline{y_i : \tau'} \vdash_{\mathcal{U}} x \sim e_{s_1} : \tau' & \quad \Gamma, \Gamma_2 \vdash_{\mathcal{U}} e_e \sim e'_s : \tau \end{aligned}$$

It follows that e_{s_1} must be x and $x \notin \text{dom}(\theta)$, so we reduce to

$$e_s[e_{s_0}/x, \overline{e'_{s_i}/y_i}] \mapsto^* E_{s_1}[e_{s_0}] \mapsto^* E_{s_1}[e'_s] \mapsto^* E_{s_1}[E_{s_2}[e'_s[\theta]]]$$

Any further reduction of e_e and e'_s preserve the relation by the *copy* rule, further expanding the other instances of e'_s .

□

Lemma 12 (Answer Preservation). *If $\Gamma \vdash_{\mathcal{U}} A \sim e : \tau$ then $e \in \text{Answer}$, and if $\Gamma \vdash_{\mathcal{U}} e \sim A : \tau$ then $e \in \text{Answer}$. Likewise, if $\Gamma \vdash_{\mathcal{N}} A \sim e : \tau$ then $\Gamma \vdash e \mapsto^* A' : \tau$, and if $\Gamma \vdash_{\mathcal{N}} e \sim A : \tau$ then $e \in \text{Answer}$.*

PROOF. The first part about \mathcal{U} follows immediately from Lemma 11 and the fact that answers are closed under substitution. The second part about \mathcal{U} follows by induction on the simulation relation. The remaining statement about \mathcal{N} follows by the reductions of the simulation, which are all instances of β_{\sim} .

Lemma 13 (Forward Simulation). *If $\Gamma \vdash_{\mathcal{N}} e_1 \sim e_2 : \tau$ and $\Gamma \vdash e_1 \mapsto e'_1 : \tau$ by the call-by-need semantics then $\Gamma \vdash e_2 \mapsto^* e'_2 : \tau$ by the call-by-name semantics.*

PROOF. First consider the cases where an operational step is applied directly to e_1 (i.e., not *compat*). These cases are as follows:

- β_{\sim} , β_V , $\beta_{\{\}}$, and β_{Int} in call-by-need follow the rule of the same name in call-by-name.
- *rename* is either immediate (due to an application of *copy*) or follows by β_{\sim} (due to an application of *share*).
- *name*_{I#} and *name* follow from the *share* unwinding rules followed by a step of *unname*_{I#} or *unname*, respectively.
- *look* and *comm* follow immediately when the *let* is simulated by *copy*, and follows by β_{\sim} when simulated by *share* or *name*.

compat, follows by Lemma 11 and the fact that reductions are stable under substitution (Lemma 3). In more detail, suppose that

$$\frac{\Gamma \vdash E_e @ e_e : \tau \xRightarrow{\Gamma'} \sigma \quad \Gamma, \Gamma' \vdash e_e \mapsto e'_e : \tau}{\Gamma \vdash E_e[e_e] \mapsto E_e[e'_e] : \sigma}$$

Then we know that, given any $\Gamma \vdash_{\mathcal{N}} E_e[e_e] \sim e_s : \sigma$, we have $\Gamma \vdash e_s \mapsto^* E_{s_1}[e_{s_1}[\theta]] \mapsto^* E_{s_1}[e'_{s_1}[\theta]] : \sigma$ such that $\Gamma \vdash_{\mathcal{N}} E_e[e'_e] \sim E_{s_1}[e'_{s_1}[\theta]]$. □

Lemma 14 (Context Rewinding). *If $\Gamma \vdash E_s @ e_s : \tau \xRightarrow{\Gamma'} \sigma$ in the call-by-name \mathcal{IL} and $\Gamma \vdash_{\mathcal{U}} e'_e \sim E_s[e_s] : \sigma$, then*

PROOF. By induction on the unwinding relation first, then on the decomposition of evaluation contexts.

- The cases for compatibility of unwinding with non-**let** syntax all follow by the inductive hypothesis.
- *share*: Where

$$\frac{\Gamma \vdash_{\mathcal{U}} e_e \sim e_s : \tau' \quad \Gamma, x : \tau' \vdash_{\mathcal{U}} e'_e \sim e'_s : \sigma}{\Gamma \vdash_{\mathcal{U}} \text{let } x:\tau' = e'_e \text{ in } e_e \sim (\lambda x:\tau' e_s) e'_s : \sigma}$$

The only possible decompositions of $(\lambda x:\tau' e_s) e'_s$ into an evaluation context is the empty context (for which the result follows immediately), or else $(\lambda x:\tau' e_s) E_s$ when $\Gamma \vdash \tau : \text{TYPE } \rho \text{ Eval}^U$ (for which the result follows from the inductive hypothesis).

- *copy*: Where

$$\frac{\Gamma \vdash_{\mathcal{U}} e'_e \sim e'_s : \tau' \quad \Gamma, \overline{y_i : \tau'} \vdash_{\mathcal{U}} e_e \sim e_s : \sigma \quad \Gamma \vdash e'_s \text{ val} \quad \Gamma \vdash e'_{s_i} \mapsto^* e'_s : \tau}{\Gamma \vdash_{\mathcal{U}} \text{let } x:\tau' = e'_e \text{ in } e_e \sim e_s[\overline{x/y_i}] \sim e_s[\overline{e'_{s_i}/y_i}] : \sigma} \text{ copy}$$

There are two possibilities for the decomposition of $e_s[\overline{e'_{s_i}/y_i}]$ into an evaluation context surrounding a sub-expression (let $\theta = [e'_{s_0}/y_0, \theta']$ and $\theta' = [e'_{s_1}/y_1, \dots]$):

- Commutation $e_s[\theta] = E_{s_1}[\theta][e_{s_1}[\theta]] = E_{s_1}[e_{s_1}][\theta]$ where the substitution does not interact with decomposition. In this case, we also have that $\Gamma, \overline{y_i : \tau'} \vdash E_{s_1} @ e_{s_1} : \tau \xRightarrow{\Gamma'} \sigma$ before substitution, and so the result follows from the inductive hypothesis.
- Interference $e_s[\theta] = E_{s_1}[\theta][E_{s_2}[\theta][e_{s_2}[\theta]]] = E_{s_1}[y_0][E_{s_2}[e_{s_2}]/y_0, \theta']$ where the substitution lands in the middle of an evaluation context (without loss of generality, assume this is the only free appearance of y_0) and that substituted expression is further decomposed. In this case, the evaluation context is composed of $E_s = E_{s_1}[E_{s_2}]$ which includes part of an expression in the substitution θ . In other words, we have

$$e_s[\theta] = E_{s_1}[\theta'] [e'_{s_0}] = E_{s_1}[\theta'] [E_{s_2}[e_{s_2}]] \mapsto^* E_{s_1}[\theta'] [e'_s]$$

where we know already that $\Gamma \vdash_{\mathcal{U}} e'_e \sim e'_s : \tau'$.

□

Lemma 15 (Backward Simulation). *If $\Gamma \vdash_{\mathcal{N}} e_1 \sim e_2 : \tau$ and $\Gamma \vdash e_2 \mapsto e'_2 : \tau$ by the call-by-name semantics then $\Gamma \vdash e_1 \mapsto^* e'_1 : \tau$ by the call-by-need semantics such that $\Gamma_{\mathcal{N}} \vdash e'_1 \sim e'_2 : \tau$.*

PROOF. Note that any unnamings expansions are either β_{\sim} redexes, or preserve the target of the evaluation context (in the case of a strict Eval^U application), but these must be done first. Assuming there are additional reductions afterward, we can proceed by induction on the derivation the number of reduction steps in $\Gamma \vdash e_2 \mapsto e'_2 : \tau$ first and the underlying unwinding $\Gamma \vdash_{\mathcal{U}} e_1 \sim e_2 : \tau$ first then. In the cases for compatibility of unwinding, we either reduce a sub-expression, which follows by induction, or the top-level expression itself. The latter cases may be:

- β_{\vee} , $\beta_{\{\}}$, and β_{Int} all follow by the inductive hypothesis by applying the call-by-need rule of the same name, using Lemma 12 to ensure the body of the λ -abstraction is still an answer in the case of β_{\vee} . Note that it is possible that a redex of the form $F[e]$ (where e **ans**) might have a **let** inserted in between the frame context and the answer as $F[B[e]]$. In this case, an additional *comm* reduction is required.
- β_{\sim} : note that in call-by-need we have

$$\Gamma \vdash (\lambda x:\tau. e_e) e'_e \mapsto_{\text{name}} \text{let } x:\tau. = e'_e \text{ in } (\lambda x:\tau. e_e) x \mapsto_{\beta_{\sim}} \text{let } x:\tau. = e'_e \text{ in } e_e : \sigma$$

such that the *copy* rule relates this to the call-by-name β_{\sim} reduct like so

$$\frac{\Gamma \vdash_{\mathcal{U}} e'_e \sim e'_s : \tau \quad \Gamma, x:\tau \vdash_{\mathcal{U}} e_e \sim e_s : \sigma}{\Gamma \vdash_{\mathcal{U}} \text{let } x:\tau. = e'_e \text{ in } e_e \sim e_s[e'_s/x] : \sigma} \text{ copy}$$

This reduction is still possible if an additional **let**-binding is inserted in the middle of the β_V redex, pushing the application to the argument a inside via *comm*.

Finally, we have the cases for unwinding a **let**-expression, which are as follows:

- *share*: Where

$$\frac{\Gamma \vdash_{\mathcal{U}} e_e \sim e_s : \tau' \quad \Gamma, x : \tau' \vdash_{\mathcal{U}} e'_e \sim e'_s : \sigma}{\Gamma \vdash_{\mathcal{U}} \text{let } x : \tau' = e'_e \text{ in } e_e \sim (\lambda x : \tau' e_s) e'_s : \sigma} \text{share}$$

The only possible decompositions of $(\lambda x : \tau' e_s) e'_s$ into an evaluation context are

- the empty context, for which the only possible reduction is β_{\sim} , whose reduct is related to the original expression by *copy* instead, or
- else $(\lambda x : \tau' e_s) e_s$ when $\Gamma \vdash \tau : \text{TYPE } \rho \text{ Eval}^U$, for which the result follows from the inductive hypothesis.

- *copy*: Where

$$\frac{\Gamma \vdash_{\mathcal{U}} e'_e \sim e'_s : \tau' \quad \Gamma, \overline{y_i : \tau'} \vdash_{\mathcal{U}} e_e \sim e_s : \sigma \quad \Gamma \vdash e'_s \text{ val} \quad \Gamma \vdash e'_{si} \mapsto^* e'_s : \tau}{\Gamma \vdash_{\mathcal{U}} \text{let } x : \tau' = e'_e \text{ in } e_e[x/\overline{y_i}] \sim e_s[\overline{e'_{si}/\overline{y_i}}] : \sigma} \text{copy}$$

There are two possibilities for the decomposition of $e_s[\overline{e'_{si}/\overline{y_i}}]$ into an evaluation context surrounding a sub-expression (let $\theta = [e'_{s0}/y_0, \theta']$ and $\theta' = [e'_{s1}/y_1, \dots]$):

- Commutation $e_s[\theta] = E_{s1}[\theta][e_{s1}[\theta]] = E_{s1}[e_{s1}][\theta]$ where the substitution does not interact with decomposition. In this case, we also have that $\Gamma, \overline{y_i : \tau'} \vdash E_{s1} @ e_{s1} : \tau \xRightarrow{\Gamma'} \sigma$ before substitution, and so the result follows from the inductive hypothesis.
- Interference $e_s[\theta] = E_{s1}[\theta][E_{s2}[\theta][e_{s2}[\theta]]] = E_{s1}[y_0][E_{s2}[e_{s2}]/y_0, \theta']$ where the substitution lands in the middle of an evaluation context (without loss of generality, assume this is the only free appearance of y_0) and that substituted expression is further decomposed. In this case, the evaluation context is composed of $E_s = E_{s1}[E_{s2}]$ which includes part of an expression in the substitution θ . In other words, we have

$$e_s[\theta] = E_{s1}[\theta'][e'_{s0}] = E_{s1}[\theta'][E_{s2}[e_{s2}]] \mapsto^* E_{s1}[\theta'][e'_s]$$

where we know already that $\Gamma \vdash_{\mathcal{U}} e'_e \sim e'_s : \tau'$. Therefore, the reduction in call-by-name is either catching up to *copy*, or if there are any steps remaining afterward, the result follows by the inductive hypothesis. \square

A.5 Correspondence to the equational theory

Here we show the correspondence between the equational theory (given in Section 3.10) and the call-by-name operational semantics for \mathcal{IL} (given in Appendix A.2). Note that, crucially, the correspondence will hold specifically for expressions of answer types τ_{ans} :

$$\tau_{ans}, \sigma_{ans} ::= \text{Int}\# \mid \text{Int}^V$$

This restriction prevents the η_{\sim} and $\eta_{\{\}}$ rules from exposing any underlying computation that needs to be evaluated, and thus, they are unnecessary for evaluating the final answer.

Theorem 5. *If $\Gamma \vdash e \mapsto^* e' : \tau$ then $\Gamma \vdash e = e' : \tau$.*

PROOF. Each call-by-name operational step is an instance of an equational axiom of \mathcal{IL} . \square

For the other direction, we can show via confluence and standarizaton. First, let reduction relation

$$\Gamma \vdash e \rightarrow e' : \tau$$

be defined as the generalization of the call-by-name operational semantics in Appendix A.2 so that compatibility (*compat*) applies to *any* context, as well as the generalization of the β_V rule to the left-to-right reading of β_V rule in Fig. 4. In other words, β_V applies to any polymorphic instantiation $(\lambda\chi.e)\phi$, not just ones where the body e is an answer. Note that the reflexive-transitive-symmetric closure of this reduction theory the same as the equational theory presented in Fig. 4, even though the β_{\sim} rule only applies to *values* of kind Eval^U rather than *passive expressions*. That's because all passive expressions evaluate to values, anyway.

Lemma 16 (Passive Value). *If $\Gamma \vdash P : \tau$ then $\Gamma \vdash P \mapsto V : \tau$.*

Theorem 6 (Confluence). *If $\Gamma \vdash e \rightarrow^* e_1 : \tau$ and $\Gamma \vdash e \rightarrow^* e_2 : \tau$ then $\Gamma \vdash e_1 \rightarrow^* e' : \tau$ and $\Gamma \vdash e_2 \rightarrow^* e' : \tau$ for some e' .*

PROOF. The only critical pairs in the reduction rules are between matching β and η rules, which join together immediately (in zero steps). As such, the reduction theory forms a combinatorial, orthogonal rewrite system. \square

Corollary 3 (Church-Rosser). *$\Gamma \vdash e = e' : \tau$ if and only if $\Gamma \vdash e \rightarrow^* e'' \leftarrow^* e' : \tau$.*

PROOF. Follows from Theorem 6. \square

Definition 2 (Internal Reduction). An *internal reduction*, written $\Gamma \vdash e \rightarrow e' : \tau$, is any reduction $\Gamma \vdash e \rightarrow e' : \tau$ such that $\Gamma \vdash e \not\mapsto e' : \tau$.

Internal reductions (that is, ones that differ from the next operational step) are useful because they are never needed to convert expressions into answers.

Lemma 17 (Answer Stability).

- (1) *If $\Gamma \vdash e \rightarrow e' : \tau$, then $e \in \text{Answer}$ if e' is, and $\Gamma \vdash e \text{ val}$ if e' is.*
- (2) *If $\Gamma \vdash e \rightarrow e' : \tau_{\text{ans}}$, then $e' \in \text{Answer}$ if e is, and $\Gamma \vdash e' \text{ val}$ if e is.*

PROOF. By induction on the syntax of answers and derivation of **val**, then cases on the possible internal reductions. Note that the restriction on the type in the second part prevents an η_{\sim} or $\eta_{\{\}}$ reduction from exposing a non-answer expression. \square

Definition 3 (Parallel Reduction). The *parallel reduction relation*, written $\Gamma \vdash e \Rightarrow e' : \tau$, is given by the restriction on $\Gamma \vdash e \rightarrow^* e' : \tau$ to only non-overlapping redexes that are all present originally in e . The *internal parallel reduction relation*, written $\Gamma \vdash e \Rightarrow e' : \tau$, is the restriction of parallel reduction to only internal reductions.

Parallel reduction is interesting because, unlike an ordinary single reduction step, it commutes with substitution in one parallel step:

Lemma 18 (Parallel Substitution).

If $\Gamma, x : \tau \vdash e : \sigma$ and $\Gamma \vdash e' \Rightarrow e'' : \tau$ then $\Gamma \vdash e[e'/x] \Rightarrow e[e''/x] : \sigma$.

PROOF. By induction on the derivation of $\Gamma, x : \tau \vdash e : \sigma$. \square

Lemma 19 (Internal Decomposition). *If $\Gamma \vdash e \Rightarrow E_2[e_2] : \sigma_{\text{ans}}$ and $\Gamma \vdash E_2 @ e_2 : \tau \xRightarrow{\Gamma'} \sigma_{\text{ans}}$ then there is an E_1 and $\theta = \overline{[\phi/\chi]}$ and τ' such that*

- $e =_{\alpha} E_1[e_1]$ and $\Gamma \vdash E_1 @ e_1 : \tau' \xRightarrow{\Gamma', \bar{\chi}} \sigma_{\text{ans}}$,
- $e_2 =_{\alpha} e_2[\theta]$ and $\Gamma, \Gamma', \chi \vdash e_1 \Rightarrow e_2 : \tau$, and
- $\Gamma \vdash E[e_3] \Rightarrow E'[e_3[\theta]] : \sigma_{\text{ans}}$ for any $\Gamma, \Gamma', \bar{\chi} \vdash e_3 : \tau$.

PROOF. By induction on the syntax of evaluation contexts, and inversion on the possible internal reductions. The substitutions θ are possible due to internal applications of β_v to an abstraction $\lambda\chi.e$ with e not an answer. Note that the restriction on the answer type required to prevent an internal η_{\sim} or $\eta_{\{\}}$ reduction from exposing a deeper evaluation context inside the body of a λ -abstraction or closure. As such, each application of these two η rules must either occur outside the eye of the evaluation context, or inside a β_{\sim} or $\beta_{\{\}}$ redex. In the latter case, such η reductions mimic β operational steps, therefore making them non-internal and ruling them out. \square

Lemma 20 (Standard Preponement). *If $\Gamma \vdash e \Rightarrow e_1 \mapsto^* e' : \tau_{ans}$ then $\Gamma \vdash e \mapsto^* e_2 \Rightarrow e' : \tau_{ans}$.*

PROOF. First, we show the case for a single step: If $\Gamma \vdash e \Rightarrow e_1 \mapsto e' : \tau$ then $\Gamma \vdash e \mapsto^* e_2 \Rightarrow e' : \tau$ for some e_2 . The cases for applying an operational step directly to an expression are as follow from Lemma 18. For example, the β_{\sim} step is as follows: assume that

$$\Gamma \vdash (\lambda x:\tau.e_1) e_2 \Rightarrow (\lambda x:\tau.e'_1) e'_2 \mapsto_{\beta_{\sim}} e'_1[e'_2/x] : \sigma$$

because $\Gamma, x : \tau \vdash e_1 \Rightarrow e'_1 : \sigma$ and $\Gamma \vdash e_2 \Rightarrow e'_2 : \tau$ and $\Gamma \vdash e'_2 \mathbf{val}$. Note that $\Gamma \vdash e_2 \mathbf{val}$ due to Lemma 17, so we can do the β_{\sim} step first, like so:

$$\Gamma \vdash (\lambda x:\tau.e_1) e_2 \mapsto_{\beta_{\sim}} e_1[e_2/x] \Rightarrow e'_1[e'_2/x]$$

The other cases are similar. Compatibility of an operational steps inside an evaluation context follows from induction on the derivation and Lemma 19. Notably, if we have

$$\frac{\Gamma \vdash E_2 @ e_2 : \sigma \xRightarrow{\Gamma'} \tau_{ans} \quad \Gamma, \Gamma' \vdash e_2 \mapsto e'_2 : \sigma}{\Gamma \vdash e \Rightarrow E_2[e_2] \mapsto E_2[e'_2] : \tau_{ans}}$$

Lemmas 3 and 19 ensure there is an E_1, e_1 , and $\theta = [\overline{\phi/\chi}]$ such that the beginning $e =_{\alpha} E_1[e_1]$, the reduced $e_2 =_{\alpha} e_3[\theta]$, and:

$$\Gamma, \Gamma', \bar{\chi} \vdash e_1 \Rightarrow e_3 \mapsto e'_3 : \tau_{ans} \quad \Gamma, \Gamma \vdash e_1[\theta] \Rightarrow e_2 \mapsto e'_3[\theta] : \tau_{ans}$$

meaning that $e'_2 =_{\alpha} e'_3[\theta]$ by Corollary 1. Therefore, from the inductive hypothesis and compatibility:

$$\Gamma, \Gamma \vdash e_1 \mapsto^* e_4 \Rightarrow e'_3 : \tau_{ans} \quad \Gamma \vdash E_1[e_1] \mapsto^* E_1[e_4] \Rightarrow E_2[e'_3[\theta]] =_{\alpha} E_2[e'_2] : \tau_{ans}$$

Finally, the case for multiple steps of the operational semantics follows from the single-step case by induction on the transitive closure of the stepping relation. \square

Lemma 21 (Internal Postponement). *If $\Gamma \vdash e \rightarrow^* e_1 \mapsto^* e' : \tau_{ans}$ then $\Gamma \vdash e \mapsto^* e_2 \rightarrow^* e' : \tau_{ans}$.*

PROOF. Note that \rightarrow^* and \Rightarrow^* are the same relation, so the theorem is equivalent to: If $\Gamma \vdash e \Rightarrow^* e_1 \mapsto^* e' : \tau_{ans}$ then $\Gamma \vdash e \mapsto^* e_2 \Rightarrow^* e' : \tau_{ans}$. This follows from Lemma 20 by induction on the reflexive-transitive closure of \Rightarrow^* . \square

Lemma 22 (Standard Order). *If $\Gamma \vdash e \rightarrow^* e' : \tau_{ans}$ then $\Gamma \vdash e \mapsto^* e'' \rightarrow^* e' : \tau_{ans}$*

PROOF. Every reduction sequence corresponds to alternations between operational steps and internal reductions:

$$e_1 \rightarrow^* e_n \text{ iff } e_1 \rightarrow^* e_1 \mapsto^* \dots \rightarrow^* e_{n-1} \mapsto^* e_n$$

Therefore, by induction on the number of these alternations, we can reorder all the operational steps to come first with Lemma 21. \square

Theorem 7 (Standardization). *If $\Gamma \vdash e \rightarrow^* A : \tau_{ans}$ then there is an $\Gamma \vdash A' \rightarrow^* A : \tau_{ans}$ such that $\Gamma \vdash e \mapsto^* A' : \tau_{ans}$.*

PROOF. From Lemma 22, we know that $\Gamma \vdash e \mapsto^* e' \rightarrow^* A : \tau$ for some e' , and from Lemma 17 e' must be an answer. \square

Corollary 4. *If $\Gamma \vdash e = A : \tau_{ans}$ then $\Gamma \vdash e \mapsto^* A' : \tau_{ans}$ such that $\Gamma \vdash A' = A : \tau_{ans}$.*

PROOF. From Corollary 3 we know that $\Gamma \vdash e \rightarrow^* A' \leftarrow^* A : \tau_{ans}$, and from Theorem 7 we know that $\Gamma \vdash e \mapsto^* A' \rightarrow^* A' \leftarrow^* A : \tau_{ans}$. \square

B CORRECTNESS OF \mathcal{IL} -TO- \mathcal{ML} COMPILATION

In this section, assume the use of black holes to mask forced thinks:

$$\begin{aligned} (\text{Force}) \quad & \langle x_{\text{PtrR}} \mid K \mid [x := \text{memo } e]H \rangle \mapsto \langle e \mid \text{set } x; K \mid [x := \bullet]H \rangle \\ (\text{Memo}) \quad & \langle R \mid \text{set } x; K \mid [x := \bullet]H \rangle \mapsto \langle R \mid K \mid [x := R]H \rangle \end{aligned}$$

B.1 Well-typed Expressions Can be Compiled

While we might only be interested in applying the compiler to whole, closed programs, we still need to be able to handle fragments of that program during compilation. In general, we must consider compiling open expressions of different types which may have free variables in them. To do so, we need to calculate the primitive representation of function arguments a , which is defined as:

$$x_\pi \overset{\text{rep}}{\rightsquigarrow} \pi \qquad n \overset{\text{rep}}{\rightsquigarrow} \text{IntR} \qquad \text{error} \overset{\text{rep}}{\rightsquigarrow} \text{PtrR}$$

These free variables are tracked in both a typing context Γ (used for type checking sub-expressions during compilation) as well as a renaming environment θ (used for replacing \mathcal{IL} variables with \mathcal{ML} ones). We need to check these two forms of environments correspond to one another, and also that each \mathcal{ML} variable has a known representation matching the type in Γ , as follows:

$$\Gamma \vdash \theta \text{ mono-rep} \text{ iff for any } x:\tau \in \Gamma, \Gamma \vdash \tau \overset{\text{rep}}{\rightsquigarrow} \pi \text{ and } \theta(x) \overset{\text{rep}}{\rightsquigarrow} \pi$$

Note that the main compilation of expressions, $C\llbracket e \rrbracket_\theta^\Gamma(\bar{a})$ generates code which evaluates e . In other words, this operation is always *strict* in e . As such, it does not matter if e is lifted or unlifted: evaluation of e is being forced either way. In general, compilation needs to know about the convention of e —if it is just being evaluated or called with a list of parameters—in order to generate the appropriate \mathcal{ML} code. However, it does *not* need to know its levity, since the code that is generated will only be run when the result is needed anyway. To state this assumption formally, we need a more relaxed notion of calculating the calling convention of a type even if the levity is unknown—in other words, a levity polymorphic version of $\tau \overset{\text{conv}}{\rightsquigarrow} v$ —as defined by the following rules for $\Gamma \vdash \tau \overset{\text{conv-lp}}{\rightsquigarrow} \eta$:

$$\frac{\Gamma \vdash \tau \overset{\text{conv}}{\rightsquigarrow} \eta}{\Gamma \vdash \tau \overset{\text{conv-lp}}{\rightsquigarrow} \eta} \qquad \frac{\Gamma \vdash \tau : \text{TYPE } \rho \text{ Eval}^l}{\Gamma \vdash \tau \overset{\text{conv-lp}}{\rightsquigarrow} \text{Eval}^l}$$

With the notion of a levity-polymorphic, but otherwise statically known, calling convention, we can state the invariants for when static compilation is defined:

Lemma 23 (Open Compilation). $\mathcal{E}_v\llbracket e \rrbracket_\theta^\Gamma$ is defined if:

- (1) $\Gamma \vdash e : \tau$ is derivable for some type τ ,
- (2) $\Gamma \vdash \tau \overset{\text{conv-lp}}{\rightsquigarrow} v$ for some v , and
- (3) $\Gamma \vdash \theta \text{ mono-rep}$.

PROOF. By induction on the given typing derivation of $\Gamma \vdash e : \tau$. Note the following invariants to ensure that $C[[e]]_{\theta}^{\Gamma}(\bar{a})$ is defined:

- (1) $\Gamma \vdash e : \tau$ is derivable for some type τ ,
- (2) $\Gamma \vdash \tau \xrightarrow{\text{conv-lp}} \nu$ for some ν , such that $|\bar{a}| = |\text{arity}(\nu)|$, and $a_i \xrightarrow{\text{rep}} \pi_i$ for each $\pi_i \in \text{arity}(\nu)$, and
- (3) $\Gamma \vdash \theta$ **mono-rep**.

Furthermore, the only invariants required for $\mathcal{P}[[P]]_{\theta}^{\Gamma}$ to be defined are that (1) $\Gamma \vdash e : \tau$ is derivable for some type τ , (2) $\Gamma \vdash \theta$ **mono-rep**. The side conditions in rules requiring $\Gamma \vdash \tau \xrightarrow{\text{rep}} \pi$ and $\Gamma \vdash \tau \xrightarrow{\text{conv}} \eta$ are ensured by the following facts about the monomorphism restrictions:

- (1) If $\vdash \rho$ **rep**, then $\rho = \pi$ for some π .
- (2) If $\vdash \gamma$ **lev**, then $\gamma = \psi$ for some ψ .
- (3) If $\vdash \nu$ **conv**, then either $\nu = \eta$ for some η .

□

B.2 An Expression-based Operational Semantics for \mathcal{ML}

$A \in \text{Answer} ::= H[V] \mid E[\text{error}(n)]$

$E \in \text{EvalCxt} ::= \square \mid F[E] \mid B[E]$

$H \in \text{ClosCxt} ::= \square \mid B[H]$

$F \in \text{FrameCxt} ::= \text{App } \square(\bar{a}) \mid \text{case } \square \text{ of } I\#(x_{\text{InTR}}) \rightarrow e \mid \text{let } x_{\pi}^{\text{U}} = \square \text{ in } e \mid \text{let } x_{\text{PtrR}}^{\text{L}} = \square \text{ in } E[x_{\text{PtrR}}]$

$B \in \text{BindCxt} ::= \text{let } x_{\text{PtrR}}^{\text{U}} = R \text{ in } \square \mid \text{let } x_{\text{PtrR}}^{\text{L}} = e \text{ in } \square$

Main reductions:

(Call)	$(\lambda(\bar{x}_{\pi}).e)(\bar{a}) \mapsto e[\bar{a}/\bar{x}_{\pi}]$	
(Apply)	$\text{App } (\text{Clos}^n P)(\bar{a}) \mapsto P(\bar{a})$	(if $ \bar{a} = n$)
(Move)	$\text{let } x_{\pi}^{\text{U}} = c \text{ in } e \mapsto e[c/x_{\pi}]$	
(Unbox)	$\text{case } I\#(a) \text{ of } I\#(x_{\text{InTR}}) \rightarrow e \mapsto e[a/x_{\text{InTR}}]$	

Pointer lookup and memoization (assume E does not bind x_{PtrR}):

(Fun)	$\text{let } x_{\text{PtrR}}^{\text{U}} = R \text{ in } E[x_{\text{PtrR}}(\bar{a})] \mapsto \text{let } x_{\text{PtrR}}^{\text{U}} = R \text{ in } E[R(\bar{a})]$
(Look)	$\text{let } x_{\text{PtrR}}^{\text{U}} = R \text{ in } E[x_{\text{PtrR}}] \mapsto \text{let } x_{\text{PtrR}}^{\text{U}} = R \text{ in } E[R]$
(Memo)	$\text{let } x_{\text{PtrR}}^{\text{L}} = R \text{ in } E[x_{\text{PtrR}}] \mapsto \text{let } x_{\text{PtrR}}^{\text{U}} = R \text{ in } E[R]$

Percolating frame contexts out of binding contexts (assume $x \notin \text{FV}(F)$):

(LAlloc)	$F[\text{let } x_{\text{PtrR}}^{\text{L}} = e \text{ in } A] \mapsto \text{let } x_{\text{PtrR}}^{\text{L}} = e \text{ in } F[A]$
(SAlloc)	$F[\text{let } x_{\text{PtrR}}^{\text{U}} = R \text{ in } A] \mapsto \text{let } x_{\text{PtrR}}^{\text{U}} = R \text{ in } F[A]$

B.3 Bisimulation between \mathcal{ML} 's operational semantics and abstract machine

Translating evaluation contexts to stacks and heaps:

$\mathcal{K}[\square] \triangleq \varepsilon$	$\mathcal{H}[\square] \triangleq \varepsilon$
$\mathcal{K}[F[E]] \triangleq \mathcal{K}[E] \circ \mathcal{K}[F]$	$\mathcal{H}[F[E]] \triangleq \mathcal{H}[E] \circ \mathcal{H}[F]$
$\mathcal{K}[B[E]] \triangleq \mathcal{K}[E]$	$\mathcal{H}[B[E]] \triangleq \mathcal{H}[B] \circ \mathcal{H}[E]$

$$\begin{aligned}
& \mathcal{K}[\llbracket \text{App } \square(\bar{a}) \rrbracket] \triangleq \text{App}(\bar{a}); \varepsilon & \mathcal{H}[\llbracket \text{App } \square(\bar{a}) \rrbracket] &\triangleq \varepsilon \\
& \mathcal{K}[\llbracket \text{case } \square \text{ of } \text{I}\#(x_{\text{InTR}}) \rightarrow e \rrbracket] \triangleq \text{case } \text{I}\#(x_{\text{InTR}}) \rightarrow e; \varepsilon & \mathcal{H}[\llbracket \text{case } \square \text{ of } \text{I}\#(x_{\text{InTR}}) \rightarrow e \rrbracket] &\triangleq \varepsilon \\
& \mathcal{K}[\llbracket \text{let } x_{\pi}^{\text{U}} = \square \text{ in } e \rrbracket] \triangleq \text{let } x_{\pi} \text{ in } e; \varepsilon & \mathcal{K}[\llbracket \text{let } x_{\pi}^{\text{U}} = \square \text{ in } e \rrbracket] &\triangleq \varepsilon \\
& \mathcal{K}[\llbracket \text{let } x_{\text{PtrR}}^{\text{L}} = \square \text{ in } E[x_{\text{PtrR}}] \rrbracket] \triangleq \text{set } x; \mathcal{K}[\llbracket E \rrbracket] & \mathcal{H}[\llbracket \text{let } x_{\text{PtrR}}^{\text{L}} = \square \text{ in } E[x_{\text{PtrR}}] \rrbracket] &\triangleq \mathcal{H}[\llbracket E \rrbracket][x := \bullet] \\
& \mathcal{H}[\llbracket \text{let } x_{\text{PtrR}}^{\text{U}} = R \text{ in } \square \rrbracket] \triangleq [x := R] \\
& \mathcal{H}[\llbracket \text{let } x_{\text{PtrR}}^{\text{L}} = e \text{ in } \square \rrbracket] \triangleq [x := \text{memo } e]
\end{aligned}$$

Definition 4 (Refocusing). The *refocusing* steps of the abstract machine are *PshApp*, *PshCase*, *PshLet*, *LAlloc*, *SAlloc*, *Force*, and *Error*. We write $m \mapsto_F m'$ for a transition by one of the refocusing steps, and $m \mapsto_R m'$ for a transition by a non-refocusing reduction step.

Lemma 24. $\langle E[e] \mid K \mid H \rangle \mapsto_F^* \langle e \mid \mathcal{K}[\llbracket E \rrbracket] \circ K \mid \mathcal{H}[\llbracket E \rrbracket] \circ H \rangle$

PROOF. By induction on the syntax of E . The most interesting case is for an evaluation context of the form $\text{let } x_{\text{PtrR}}^{\text{L}} = E_2 \text{ in } E_1[x_{\text{PtrR}}]$, which proceeds as follows:

$$\begin{aligned}
\langle \text{let } x_{\text{PtrR}}^{\text{L}} = E_2[e] \text{ in } E_1[x_{\text{PtrR}}] \mid K \mid H \rangle &\mapsto_{\text{Alloc}} \langle E_1[x_{\text{PtrR}}] \mid K \mid [x := \text{memo } E_2[e]]H \rangle \\
&\mapsto_{\text{IH}} \langle x_{\text{PtrR}} \mid \mathcal{K}[\llbracket E_1 \rrbracket] \circ K \mid \mathcal{H}[\llbracket E_1 \rrbracket] \circ [x := \text{memo } E_2[e]]H \rangle \\
&\mapsto_{\text{Force}} \langle E_2[e] \mid \mathcal{K}[\llbracket E_1 \rrbracket] \circ K \mid \mathcal{H}[\llbracket E_1 \rrbracket] \circ [x := \bullet]H \rangle \\
&\mapsto_{\text{IH}} \langle e \mid \mathcal{K}[\llbracket E_2 \rrbracket] \circ \mathcal{K}[\llbracket E_1 \rrbracket] \circ K \mid \mathcal{H}[\llbracket E_2 \rrbracket] \circ \mathcal{H}[\llbracket E_1 \rrbracket] \circ [x := \bullet]H \rangle
\end{aligned}$$

And note that the corresponding stack is $\mathcal{K}[\llbracket \text{let } x_{\text{PtrR}}^{\text{L}} = E_2 \text{ in } E_1[x_{\text{PtrR}}] \rrbracket] = \mathcal{K}[\llbracket E_2 \rrbracket] \circ \mathcal{K}[\llbracket E_1 \rrbracket]$ and heap is $\mathcal{H}[\llbracket \text{let } x_{\text{PtrR}}^{\text{L}} = E_2 \text{ in } E_1[x_{\text{PtrR}}] \rrbracket] = \mathcal{H}[\llbracket E_2 \rrbracket] \circ \mathcal{H}[\llbracket E_1 \rrbracket] \circ [x := \bullet]$ \square

Corollary 5 (Answer Preservation). If $A \sim m$ then $m \mapsto_F^* A'$. If $e \sim A$, then $e \in \text{Answer}$.

PROOF. By definition of answer expressions and machine states in \mathcal{ML} , from Lemma 24, taking a final *Error* step in the erroneous case $E[\text{error}(n)]$. \square

Definition 5 (Machine Simulation). The simulation relation between \mathcal{ML} expressions and machine states is

$$e \sim m \text{ iff } \langle e \mid \varepsilon \mid \varepsilon \rangle \mapsto_F^* m$$

Lemma 25 (Forward Simulation). (1) If $e \mapsto e'$ by *LAlloc*, or *SAlloc* and $e \sim m$ then $e' \sim m$.

(2) If $e \mapsto e'$ by any other step and $e \sim m$ then $m \mapsto_F^* m' \mapsto_R m''$ and $e' \sim m''$.

Therefore, if $e \mapsto^* e'$ and $e \sim m$ then $m \mapsto^* m'$ and $e' \sim m'$ for some m' .

PROOF. First, consider the cases where $e \mapsto e'$ by applying a single step of the reduction rules directly. Each of the rules in the operational semantics corresponds to the machine rule of the same name.

- *LAlloc* is, due Lemma 24 and the fact that the heap is unordered:

$$\begin{aligned}
\langle F[\llbracket \text{let } x_{\text{PtrR}}^{\text{L}} = e \text{ in } A \rrbracket] \mid \varepsilon \mid \varepsilon \rangle &\mapsto_F^* \langle \text{let } x_{\text{PtrR}}^{\text{L}} = e \text{ in } A \mid \mathcal{K}[\llbracket F \rrbracket] \mid \mathcal{H}[\llbracket F \rrbracket] \rangle \\
&\mapsto_{\text{LAlloc}} \langle A \mid \mathcal{K}[\llbracket F \rrbracket] \mid [x := \text{memo } e] \mathcal{H}[\llbracket F \rrbracket] \rangle \\
&\leftarrow_F^* \langle F[A] \mid \varepsilon \mid [x := \text{memo } e] \rangle \\
&\leftarrow_F \langle \text{let } x_{\text{PtrR}}^{\text{L}} = e \text{ in } F[A] \mid \varepsilon \mid \varepsilon \rangle
\end{aligned}$$

- *SAlloc* is similar to *LAlloc*.
- *Call* is $\langle (\lambda(\bar{x}_{\pi}).e)(\bar{a}) \mid \varepsilon \mid \varepsilon \rangle \mapsto_{\text{Call}} \langle e[\bar{a}/\bar{x}_{\pi}] \mid \varepsilon \mid \varepsilon \rangle$.
- *Move* is similar to *Call*.

- *Apply* is $\langle \text{App} (\text{Clos}^n P)(\bar{a}) \mid \varepsilon \mid \varepsilon \rangle \mapsto_{\text{PshApp}} \langle \text{Clos}^n P \mid \text{App}(\bar{a}); \varepsilon \mid \varepsilon \rangle \mapsto_{\text{Apply}} \langle P(\bar{a}) \mid \varepsilon \mid \varepsilon \rangle$.
- *Unbox* is similar to *Apply*
- *Look* is

$$\begin{aligned}
\langle \text{let } x_{\text{PtrR}}^U = R \text{ in } E[x_{\text{PtrR}}] \mid K \mid H \rangle &\mapsto_F^* \langle E[x_{\text{PtrR}}] \mid K \mid [x := R]H \rangle \\
&\mapsto_F^* \langle x_{\text{PtrR}} \mid \mathcal{K}[\![E]\!] \circ K \mid \mathcal{K}[\![E]\!] \circ [x := R]H \rangle \\
&\mapsto_{\text{Look}} \langle R \mid \mathcal{K}[\![E]\!] \circ K \mid \mathcal{K}[\![E]\!] \circ [x := R]H \rangle \\
&\leftarrow_F^* \langle E[R] \mid K \mid [x := R]H \rangle \\
&\leftarrow_F^* \langle \text{let } x_{\text{PtrR}}^U = R \text{ in } E[R] \mid K \mid H \rangle
\end{aligned}$$

- *Fun* is similar.
- *Memo* is

$$\begin{aligned}
\langle \text{let } x_{\text{PtrR}}^L = R \text{ in } E[x_{\text{PtrR}}] \mid K \mid H \rangle &\mapsto_F \langle E[x_{\text{PtrR}}] \mid K \mid [x := \text{memo } R]H \rangle \\
&\mapsto_F^* \langle x_{\text{PtrR}} \mid \mathcal{K}[\![E]\!] \circ K \mid \mathcal{H}[\![E]\!] \circ [x := \text{memo } R]H \rangle \\
&\mapsto_F \langle R \mid \text{set } x; \mathcal{K}[\![E]\!] \circ K \mid \mathcal{H}[\![E]\!] \circ [x := \bullet]H \rangle \\
&\mapsto_{\text{Memo}} \langle R \mid \mathcal{K}[\![E]\!] \circ K \mid \mathcal{H}[\![E]\!] \circ [x := R]H \rangle \\
&\leftarrow_{\text{Look}} \langle x_{\text{PtrR}} \mid \mathcal{K}[\![E]\!] \circ K \mid \mathcal{H}[\![E]\!] \circ [x := R]H \rangle \\
&\leftarrow_F^* \langle E[x_{\text{PtrR}}] \mid K \mid [x := R]H \rangle \\
&\leftarrow_F^* \langle \text{let } x_{\text{PtrR}}^U = R \text{ in } E[x_{\text{PtrR}}] \mid K \mid H \rangle
\end{aligned}$$

Now, assume that $e \mapsto e'$ by applying one of the reduction rules as in the above steps, by Lemma 24 and due to the fact that refocusing steps are deterministic, reduction within an evaluation context $E[e] \mapsto E[e']$ proceeds like so:

$$\begin{aligned}
\langle E[e] \mid K \mid H \rangle &\mapsto_F^* m \mapsto_F^* \langle e \mid \mathcal{K}[\![E]\!] \circ K \mid \mathcal{H}[\![E]\!] \circ H \rangle \\
&\mapsto^* m' \leftarrow_F^* \langle e' \mid \mathcal{K}[\![E]\!] \circ K \mid \mathcal{H}[\![E]\!] \circ H \rangle_F^* \leftarrow \langle E[e'] \mid K \mid H \rangle
\end{aligned}$$

so that $E[e'] \sim \langle e' \mid \mathcal{K}[\![E]\!] \circ K \mid \mathcal{H}[\![E]\!] \circ H \rangle$. Finally, multiple reductions $e \mapsto^* e'$ follows from transitivity by induction on the number of steps. \square

Lemma 26 (Backward Simulation). *If $m \mapsto^* m'$ and $e \sim m$ then $e' \sim m'$ and $e \mapsto^* e'$ for some e' .*

PROOF. Note that if $m \mapsto_F^* m'$ then $e \sim m'$ directly. Otherwise, if $m \mapsto_R m'$ then $e \mapsto_{\text{LAlloc, SAlloc}}^* e'' \mapsto e'$ and by forward simulation Lemma 25 and determinism of the abstract machine, we know that $e' \sim m'$. \square

B.4 Bisimulation between \mathcal{IL} and \mathcal{ML}

Here, we relate the environment-based call-by-need operational semantics for \mathcal{IL} with the operational semantics for \mathcal{ML} , so in the following, assume those definitions given in Appendix A.3. The simulation relationship between the two languages is given by compilation. Note that the compilation of **let**-expressions can be derived from the encoding of **lets** as applied λ -abstractions.

Definition 6 (Compilation Simulation). The simulation relation between open, typed \mathcal{IL} expressions $\Gamma \vdash e : \tau : \text{TYPE } \rho \vee$ and \mathcal{ML} expressions is defined as

$$\Gamma \vdash_C e \sim e' : \tau \text{ iff } e' \mapsto_{\text{Unlift}}^* C[\![e]\!]_{\theta}^{\Gamma}(\bar{a})$$

for some θ and \bar{a} such that $\Gamma \vdash \theta$ **mono-rep** and $\bar{a} \stackrel{\text{rep}}{\rightsquigarrow} \text{arity}(\nu)$, where the *Unlift* reduction is:

$$(\text{Unlift}) \quad \text{let } x_{\text{PtrR}}^L = R \text{ in } e \mapsto \text{let } x_{\text{PtrR}}^U = R \text{ in } e$$

For the purpose of this simulation, we treat an *Apply* followed by a *Call* as a single step in \mathcal{ML} .

Lemma 27 (Passive Stability). $\mathcal{E}_v \llbracket P \rrbracket_\theta^\Gamma \triangleq \mathcal{P} \llbracket P \rrbracket_\theta^\Gamma$ and $C \llbracket P' \rrbracket_\theta^\Gamma(\varepsilon) \triangleq \mathcal{P} \llbracket P' \rrbracket_\theta^\Gamma$ when defined and when P' is not a variable.

PROOF. By induction on the syntax of P . \square

Lemma 28 (Answer Preservation). If $\Gamma \vdash A$ **ans** and $\Gamma \vdash_C A \sim e : \tau$ then $e \in \text{Answer}$. If $\Gamma \vdash_C e \sim A : \tau$, then $e \in \text{Answer}$.

PROOF. The first part follows by induction on the syntax of \mathcal{IL} answers. Of note, in the base cases are where A are references $\text{I}\#^y \text{arg}$ or $\text{Clos}^y e$ (which have a convention of Eval^y and cannot be applied to any arguments) or a constant $c \bar{\phi}$ (which could have any calling convention but cannot reduce further, even when applied). Dually, the only \mathcal{IL} expressions that can compile to \mathcal{ML} answers are \mathcal{IL} answers. \square

Translation of \mathcal{IL} evaluation contexts (surrounding an expression of type τ) to \mathcal{ML} evaluation contexts is written as $C_\tau \llbracket E \rrbracket_\theta^\Gamma(\bar{a})$. The definition for most cases, where strictness is obvious from the syntax, are defined as follows:

$$C_\tau \llbracket \square \rrbracket_\theta^\Gamma(\bar{a}) \triangleq \square \bar{a}$$

$$C_\tau \llbracket E a' \rrbracket_\theta^\Gamma(\bar{a}) \triangleq C_\tau \llbracket E \rrbracket_\theta^\Gamma(\mathcal{P} \llbracket a' \rrbracket_\theta^\Gamma, \bar{a})$$

$$C_\tau \llbracket \text{case } E \text{ of } \text{I}\# x \rightarrow e \rrbracket_\theta^\Gamma(\bar{a}) \triangleq \text{case } C_\tau \llbracket E \rrbracket_\theta^\Gamma(\varepsilon) \text{ of } \text{I}\# x \rightarrow C \llbracket e \rrbracket_\theta^\Gamma(\bar{a})$$

$$C_\tau \llbracket \text{App } E \rrbracket_\theta^\Gamma(\bar{a}) \triangleq \text{App } (C_\tau \llbracket E \rrbracket_\theta^\Gamma(\varepsilon))(\bar{a})$$

$$C_\tau \llbracket E \phi \rrbracket_\theta^\Gamma(\bar{a}) \triangleq C_\tau \llbracket E \rrbracket_\theta^\Gamma(\bar{a})$$

$$C_\tau \llbracket \lambda \chi. E \rrbracket_\theta^\Gamma(\bar{a}) \triangleq C_\tau \llbracket E \rrbracket_\theta^\Gamma(\bar{a})$$

The most complicated cases are for **let**-expressions, which rely on the types of the variable binding to determine the difference between strictness and laziness. These cases are defined (top-to-bottom) as follows:

$$C_\tau \llbracket \text{let } x:\sigma = R \text{ in } E \rrbracket_\theta^\Gamma(\bar{a}) \triangleq \text{let } x_{\text{PtrR}}^U = \mathcal{P} \llbracket R \rrbracket_\theta^\Gamma \text{ in } C_\tau \llbracket E \rrbracket_{[x_{\text{PtrR}}/x]\theta}^{\Gamma, x:\sigma}(\bar{a})$$

$$C_\tau \llbracket \text{let } x:\sigma = e \text{ in } E \rrbracket_\theta^\Gamma(\bar{a}) \triangleq \text{let } x_{\text{PtrR}}^{\text{lev}(\eta)} = \mathcal{E}_\eta \llbracket e \rrbracket_\theta^\Gamma \text{ in } C_\tau \llbracket E \rrbracket_{[x_{\text{PtrR}}/x]\theta}^{\Gamma, x:\sigma}(\bar{a})$$

$$(\text{if } \Gamma \vdash e \text{ bind and } \Gamma \vdash \sigma \xrightarrow{\text{conv}} \eta)$$

$$C_\tau \llbracket \text{let } x:\sigma = E \text{ in } e \rrbracket_\theta^\Gamma(\bar{a}) \triangleq \text{let } x_\pi^U = C_\tau \llbracket E \rrbracket_\theta^\Gamma(\varepsilon) \text{ in } C \llbracket e \rrbracket_{[x_\pi/x]\theta}^{\Gamma, x:\sigma}(\bar{a})$$

$$(\text{if } \Gamma \vdash \sigma \xrightarrow{\text{conv}} \text{Eval}^U \text{ and } \Gamma \vdash \sigma \xrightarrow{\text{rep}} \pi)$$

$$C_\tau \llbracket \text{let } x:\sigma = E_2 \text{ in } E_1[x] \rrbracket_\theta^\Gamma(\bar{a}) \triangleq \text{let } x_{\text{PtrR}}^L = C_\tau \llbracket E \rrbracket_\theta^\Gamma(\varepsilon) \text{ in } C_\sigma \llbracket E_1 \rrbracket_{[x_{\text{PtrR}}/x]\theta}^{\Gamma, x:\sigma}(\bar{a})[x_{\text{PtrR}}]$$

$$(\text{if } \Gamma \vdash \sigma \xrightarrow{\text{conv}} \text{Eval}^L)$$

Note that \mathcal{IL} closing, frame, and binding contexts are all special cases of evaluation contexts, and are defined as above, and all translate to their corresponding special case in \mathcal{ML} .

Lemma 29 (Context Compilation). If $\Gamma \vdash E @ e : \tau \xrightarrow{\Gamma'} \sigma$ then

$$C \llbracket E[e] \rrbracket_\theta^\Gamma(\bar{a}) \triangleq C_\tau \llbracket E \rrbracket_\theta^\Gamma(\bar{a})[C \llbracket e \rrbracket_\theta^\Gamma(\varepsilon)]$$

PROOF. By induction on the derivation of $\Gamma \vdash E @ e : \tau \xrightarrow{\Gamma'} \sigma$. \square

Lemma 30 (Instantiation). (1) For any $\Gamma, x : \tau \vdash e : \sigma$ and $\Gamma \vdash a' : \tau$, renaming an argument for a variable commutes with compilation: $C\llbracket e[a'/x] \rrbracket_{\theta}^{\Gamma}(\bar{a}) \triangleq C\llbracket e \rrbracket_{[\mathcal{P}\llbracket a' \rrbracket_{\theta}^{\Gamma}/x]_{\theta}}^{\Gamma, x:\tau}(\bar{a})$.

(2) For any $\Gamma, \chi \vdash e : \sigma$ and $\Gamma \vdash [\phi/\chi]$ **poly**, instantiation of a type variable commutes with compilation: $C\llbracket e[\phi/\chi] \rrbracket_{\theta}^{\Gamma}(\bar{a}) \triangleq C\llbracket e \rrbracket_{\theta}^{\Gamma, \chi}(\bar{a})$.

PROOF. By induction on the syntax of e . \square

Lemma 31 (Forward Simulation). For any \mathcal{IL} expression e_i and \mathcal{ML} expression e_m ,

- (1) If $\Gamma \vdash_C e_i \sim e_m : \tau$ and $\Gamma \vdash e_i \mapsto e'_i : \tau$ by a β_{\sim}, β_V , rename, name, or $\text{name}_{I\#}$ step, then $\Gamma \vdash e'_i \sim e_m : \tau$.
- (2) If $\Gamma \vdash_C e_i \sim e_m : \tau$ and $e_i \mapsto e'_i$ by a $\beta_{\{\}}, \beta_{\text{Int}}$, look, or comm step, then $e_m \mapsto e'_m$ and $\Gamma \vdash e'_i \sim e_m : \tau$ for some e'_m .

Therefore, if $\Gamma \vdash_C e_i \sim e_m : \tau$ and $\Gamma \vdash e_i \mapsto^* e'_i : \tau$ then $e_m \mapsto^* e'_m$ and $\Gamma \vdash e'_i \sim e'_m : \tau$ for some e'_m .

PROOF. First consider the cases where $\Gamma \vdash e_i \mapsto e'_i : \tau$ by applying a single step of the reduction rules directly (i.e., not *compat*). The first case of \mathcal{IL} reductions, which are erased by compilation are:

- β_{\sim} $C\llbracket (\lambda x:\tau. e) a' \rrbracket_{\theta}^{\Gamma}(\bar{a}) \triangleq C\llbracket \lambda x:\tau. e \rrbracket_{\theta}^{\Gamma}(\mathcal{P}\llbracket a' \rrbracket_{\theta}^{\Gamma}, \bar{a}) \triangleq C\llbracket e \rrbracket_{[\mathcal{P}\llbracket a' \rrbracket_{\theta}^{\Gamma}/x]_{\theta}}^{\Gamma, x:\tau}(\bar{a}) \triangleq C\llbracket e[a'/x] \rrbracket_{\theta}^{\Gamma}(\bar{a})$ which follows from Lemma 30.
- *rename* is the same as β_{\sim} , via the encoding of **let** as applied λ -abstractions.
- β_V is $C\llbracket (\lambda \chi. e) \phi \rrbracket_{\theta}^{\Gamma}(\bar{a}) \triangleq C\llbracket \lambda \chi. e \rrbracket_{\theta}^{\Gamma}(\bar{a}) \triangleq C\llbracket e \rrbracket_{\theta}^{\Gamma, \chi}(\bar{a}) \triangleq C\llbracket e \rrbracket_{\theta}^{\Gamma}(\bar{a})$ which follows from Lemma 30.
- $\text{name}_{I\#}$ is $C\llbracket I\#^Y e \rrbracket_{\theta}^{\Gamma}(\varepsilon) \triangleq \text{let } x_{\text{IntR}}^U = C\llbracket e \rrbracket_{\theta}^{\Gamma}(\varepsilon) \text{ in } I\#(x_{\text{IntR}}) \triangleq C\llbracket \text{let } x:\text{Int}\# = e \text{ in } I\#^Y x \rrbracket_{\theta}^{\Gamma}(\varepsilon)$.
- *name* is similar to $\text{name}_{I\#}$, and depends on the argument of application and its kind.

The second case of \mathcal{IL} reductions that are mirrored by compilation are:

- $\beta_{\{\}}$ is

$$\begin{aligned} C\llbracket \text{App } (\text{Clos}^Y e) \rrbracket_{\theta}^{\Gamma}(\bar{a}) &\triangleq \text{App } (\text{Clos}^n \mathcal{E}_{\eta}\llbracket e \rrbracket_{\theta}^{\Gamma})(\bar{a}) \\ &\mapsto_{\text{Apply}} \mathcal{E}_{\eta}\llbracket e \rrbracket_{\theta}^{\Gamma}(\bar{a}) \\ &\mapsto_{\text{Call}}^? C\llbracket e \rrbracket_{\theta}^{\Gamma}(\bar{a}) \end{aligned}$$

assuming $\Gamma \vdash e : \tau$ and $\Gamma \vdash \tau \xrightarrow{\text{conv}} \eta$, which forces $\bar{\pi} = \text{arity}(\eta)$ and $n = |\bar{\pi}| = |\bar{a}|$ and $\bar{a} \xrightarrow{\text{rep}} \bar{\pi}$.

- β_{Int} follows from Lemma 30

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

- *comm* is, for some v and ψ

$$\begin{aligned} C\llbracket F[\text{let } x:\tau = e \text{ in } A] \rrbracket_{\theta}^{\Gamma}(\bar{a}) &\leftarrow_{\text{Unlift}}^* C_{\sigma}\llbracket F \rrbracket_{\theta}^{\Gamma}(\bar{a})[\text{let } x_{\text{PtrR}}^{\psi} = \mathcal{E}_v\llbracket e \rrbracket_{\theta}^{\Gamma} \text{ in } C\llbracket A \rrbracket_{\theta}^{\Gamma}(\varepsilon)] \\ &\mapsto_{\text{Alloc}} \text{let } x_{\text{PtrR}}^{\psi} = \mathcal{E}_v\llbracket e \rrbracket_{\theta}^{\Gamma} \text{ in } C_{\sigma}\llbracket F \rrbracket_{\theta}^{\Gamma}(\bar{a})[C\llbracket A \rrbracket_{\theta}^{\Gamma}(\varepsilon)] \\ &\triangleq C\llbracket \text{let } x:\tau = e \text{ in } F[A] \rrbracket_{\theta}^{\Gamma}(\bar{a}) \end{aligned}$$

where $\text{Alloc} = \text{LAlloc}$ if $\psi = \text{L}$ and $\text{Alloc} = \text{SAlloc}$ if $\psi = \text{U}$.

- *look* is, by Lemma 27, for some v and ψ

$$\begin{aligned}
C[\llbracket \text{let } x:\tau = e \text{ in } E[x] \rrbracket_\theta^\Gamma(\bar{a})] &\leftarrow_{\text{Unlift}}^* \text{let } x_{\text{PtrR}}^\psi = \mathcal{E}_v[\llbracket e \rrbracket_\theta^\Gamma] \text{ in } C_\sigma[\llbracket E \rrbracket_\theta^\Gamma(\bar{a})][x_{\text{PtrR}}] \\
&\mapsto_{\text{Inline}} \text{let } x_{\text{PtrR}}^u = \mathcal{E}_v[\llbracket e \rrbracket_\theta^\Gamma] \text{ in } C_\sigma[\llbracket E \rrbracket_\theta^\Gamma(\bar{a})][\mathcal{E}_v[\llbracket e \rrbracket_\theta^\Gamma]] \\
&\triangleq C[\llbracket \text{let } x:\tau = e \text{ in } E[e] \rrbracket_\theta^\Gamma(\bar{a})]
\end{aligned}$$

where *Inline* = *Look* or *Inline* = *Memo* if $v = \text{Eval}^Y$ and *Inline* = *Fun* if $v = \text{Call}[\bar{\pi}]$.

Now, assume that $\Gamma \vdash e \mapsto e' : \tau$ by applying one of the reduction rules as in the above steps. Reduction within an evaluation context $E[e] \mapsto E[e']$ proceeds like so:

$$\begin{aligned}
C[\llbracket E[e] \rrbracket_\theta^\Gamma(\bar{a})] &\triangleq C_\tau[\llbracket E \rrbracket_\theta^\Gamma(\bar{a})][C[\llbracket e \rrbracket_\theta^\Gamma(\epsilon)]] \leftarrow_{\text{Unlift}}^* E''[e''] \\
&\mapsto C_\tau[\llbracket E \rrbracket_\theta^\Gamma(\bar{a})][C[\llbracket e' \rrbracket_\theta^\Gamma(\epsilon)]] \\
&\mapsto_{\text{Unlift}}^* C[\llbracket E[e'] \rrbracket_\theta^\Gamma(\bar{a})]
\end{aligned}$$

where additional applications of the *Unlift* step may be needed when a bound lazy expression (of kind Eval^L) is finally reduced to a reference R . \square

Lemma 32 (Backward Simulation). *If $\Gamma \vdash_C e_i \sim e_m : \tau$ and $e_m \mapsto^* e'_m$ then $\Gamma \vdash e_i \mapsto^* e'_i : \tau$ and $\Gamma \vdash_C e'_i \sim e'_m$ for some e'_i .*

PROOF. Note that if $\Gamma \vdash_C e_i \sim e_m : \tau$ then $e_m \not\mapsto$ by the *Call* or *Move* steps. Otherwise, given $e_m \mapsto e'_m$ by any other step, it is the case that $\Gamma \vdash e_i \mapsto^* e''_i : \tau$ by $\beta_{\rightarrow}, \beta_V, \text{rename}, \text{name}$, and $\text{name}_{I\#}$ steps, followed by $\Gamma \vdash e''_i \mapsto e'_i : \tau$ by one of the other steps. Therefore, by forward simulation Lemma 25 and determinism of the \mathcal{ML} operational semantics, we know that $\Gamma \vdash_C e'_i \sim e'_m$. \square

B.5 Correctness of closed compilation

The correspondence between the equational theory of \mathcal{IL} and the abstract machine of \mathcal{ML} is based on four parts:

- (1) Standardization and confluence relating the equational theory of \mathcal{IL} to the call-by-name operational semantics of \mathcal{IL} (defined as a sub-relation of the equational theory).
- (2) A bisimulation between the call-by-name and call-by-need operational semantics of \mathcal{IL} based on *unwinding* (i.e., copying) *let*-bindings.
- (3) A bisimulation between the call-by-need operational semantics of \mathcal{IL} and the operational semantics of \mathcal{ML} based on the compilation given in Fig. 6.
- (4) A bisimulation between the operational semantics and abstract machine of \mathcal{ML} .

Each of these four parts are bi-directional, and bisimulations compose with one another. Therefore, we can trace high-level equalities in \mathcal{IL} all the way down to the \mathcal{ML} machine, and back. The only restriction imposed is the types of answers allowed: a necessary restriction to respect the full η -extensionality we're after.

Theorem 8 (Soundness and Completeness).

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

PROOF. We show the case for $\text{Int}\#$ as Int^Y is analogous.

First, from the assumption that $\vdash e = n : \text{Int}\#$, we know:

- $\vdash e \mapsto^* n : \text{Int}\#$ in call-by-name \mathcal{IL} from Corollary 4,
- $\vdash e \mapsto^* H[n] : \text{Int}\#$ in call-by-need \mathcal{IL} from Lemmas 12 and 15 and the fact that n is closed,
- $\mathcal{E}_{\text{Eval}^U}[\llbracket e \rrbracket] \mapsto^* H'[n]$ in \mathcal{ML} from Lemmas 28 and 31, and

Types σ corresponding to the source call-by-name System F

$$\sigma ::= a \mid \text{Int}^\gamma \mid \gamma\{\tau\} \mid \forall t:\kappa.\sigma \quad \tau ::= \sigma \rightsquigarrow \sigma' \quad \kappa ::= \text{TYPE PtrR Eval}^\gamma \quad \gamma ::= \text{L}$$

Types σ corresponding to the source call-by-value System F

$$\tau ::= a \mid \text{Int}^\gamma \mid \gamma\{\tau\} \quad \tau ::= \sigma \rightsquigarrow \sigma' \mid \forall t:\kappa.\sigma \quad \kappa ::= \text{TYPE PtrR Eval}^\gamma \quad \gamma ::= \text{U}$$

Decompilation of source types $\llbracket \sigma \rrbracket^{-1}$ and target-only types $\boxed{\tau}$

$$\llbracket a \rrbracket^{-1} \triangleq a \quad \llbracket \text{Int}^\gamma \rrbracket^{-1} \triangleq \text{Int} \quad \llbracket \gamma\{\tau\} \rrbracket^{-1} \triangleq \llbracket \tau \rrbracket^{-1} \quad \llbracket \forall t:\kappa.\sigma \rrbracket^{-1} \triangleq \forall t.\llbracket \sigma \rrbracket^{-1} \quad \llbracket \sigma \rightsquigarrow \sigma' \rrbracket^{-1}$$

Fig. 9. Decompiling \mathcal{IL} types back to System F.

Expressions corresponding to call-by-name -value System F (invariant: $x : \sigma$)

$$e ::= x \mid \text{I}\#^\gamma n \mid e \, e' \mid \lambda x:\sigma.e \mid e \, \sigma \mid \lambda t:\kappa.e \mid \text{App } e \mid \text{Clos}^\gamma S$$

Decompilation of serious expressions $\llbracket e \rrbracket^{-1}$ (invariant: $e : \tau$ or $e : \sigma$)

$$\begin{aligned} \llbracket x \rrbracket^{-1} &\triangleq x & \llbracket \text{I}\#^\gamma n \rrbracket^{-1} &\triangleq n \\ \llbracket e \, e' \rrbracket^{-1} &\triangleq \llbracket e \rrbracket^{-1} \llbracket e' \rrbracket^{-1} & \llbracket \lambda x:\sigma.e \rrbracket^{-1} &\triangleq \lambda x:\sigma.\llbracket e \rrbracket^{-1} \\ \llbracket e \, \sigma \rrbracket^{-1} &\triangleq \llbracket e \rrbracket^{-1} \llbracket \sigma \rrbracket^{-1} & \llbracket \lambda t:\kappa.e \rrbracket^{-1} &\triangleq \lambda t.\llbracket e \rrbracket^{-1} \\ \llbracket \text{App } e \rrbracket^{-1} &\triangleq \llbracket e \rrbracket^{-1} & \llbracket \text{Clos}^\gamma V \rrbracket^{-1} &\triangleq \eta \llbracket V \rrbracket^{-1} \end{aligned}$$

η -expanded decompilation of values $\eta \llbracket e \rrbracket^{-1}$ (invariant: $e : \tau$ and e is a value)

$$\begin{aligned} \eta \llbracket \lambda x:\sigma.e \rrbracket^{-1} &\triangleq \lambda x:\llbracket \sigma \rrbracket^{-1}.\llbracket e \rrbracket^{-1} & \eta \llbracket \text{App } e \rrbracket^{-1} &\triangleq \lambda x:\llbracket \sigma \rrbracket^{-1}.\llbracket e \rrbracket^{-1} x \quad (\text{if } \text{App } e : \sigma \rightsquigarrow \sigma') \\ \eta \llbracket \lambda t:\kappa.e \rrbracket^{-1} &\triangleq \lambda t.\llbracket e \rrbracket^{-1} & \eta \llbracket \text{App } e \rrbracket^{-1} &\triangleq \lambda t.\llbracket e \rrbracket^{-1} t \quad (\text{if } \text{App } e : \forall t:\kappa.\sigma) \end{aligned}$$

Fig. 10. Decompiling \mathcal{IL} expressions back to System F.

- $\langle \mathcal{E}_{\text{Eval}}^U \llbracket e \rrbracket \mid \varepsilon \mid \varepsilon \rangle \mapsto^* \langle n \mid \varepsilon \mid H \rangle$ from Lemma 25 and Corollary 5.

Second, from the assumption that $\langle \mathcal{E}_{\text{Eval}}^U \llbracket e \rrbracket \mid \varepsilon \mid \varepsilon \rangle \mapsto^* \langle n \mid \varepsilon \mid H \rangle$, we know:

- $\mathcal{E}_{\text{Eval}}^U \llbracket e \rrbracket \mapsto^* H[n]$ in \mathcal{ML} from Lemma 26 and Corollary 5,
- $\vdash e \mapsto^* H'[n] : \text{Int}\#$ in call-by-need \mathcal{IL} from Lemmas 28 and 32,
- $\vdash e \mapsto^* n : \text{Int}\#$ in call-by-name \mathcal{IL} from Lemmas 12 and 13 and the fact that n is closed, and
- $\vdash e = n : \text{Int}\#$ from Theorem 5.

□

C CORRECTNESS OF SYSTEM F-TO- \mathcal{IL} COMPILE

Decompilation of types and expressions from \mathcal{IL} back to System F is shown in Figs. 9 and 10 for call-by-name and call-by-value, respectively, which is defined over the sublanguage of \mathcal{IL} that is reachable from compiling System F. Note that this sublanguage is closed under reduction.

We now show that both call-by-name and call-by-value compilation a decompilation form an equational correspondence between System F and \mathcal{IL} .

Lemma 33 (η -Expansion). *For both the call-by-value and call-by-name \mathcal{IL} sublanguages:*

- (1) if $\Gamma \vdash e : \sigma \rightsquigarrow \sigma'$, then $\eta \llbracket e \rrbracket^{-1} =_{\beta_{\text{U}}} \lambda x:\llbracket \sigma \rrbracket^{-1}.\llbracket e \rrbracket^{-1} x$, and

(2) if $\Gamma \vdash e : \forall t:\kappa.\sigma$, then $\eta\llbracket e \rrbracket^{-1} =_{\beta_V} \lambda t:\kappa.\llbracket e \rrbracket^{-1} t$.

PROOF. By cases on the form of e . First, note that either type of $\eta\llbracket \text{App } e \rrbracket^{-1}$ is definitionally equal to the right-hand side. If instead e is a λ -abstraction (of either type), then $\eta\llbracket e \rrbracket^{-1}$ β -expands to the right-hand side. \square

Lemma 34 (Value Preservation). *For both call-by-name and call-by-value:*

- (1) Given any value $\Gamma \vdash V : \tau$ in System F, $\llbracket V \rrbracket$ is passable in \mathcal{IL} .
- (2) Given any $\Gamma \vdash e$ **pass** in the \mathcal{IL} sublanguage, $\llbracket e \rrbracket^{-1}$ is a value in System F.
- (3) Given any $\Gamma \vdash e$ **pass** in the \mathcal{IL} sublanguage, $\eta\llbracket e \rrbracket^{-1}$ is a value in System F.

PROOF. Call-by-name value preservation is immediate, because every expression is a value in call-by-name System F, and every System F expression compiles to one with a type of kind TYPE PtrR Eval^L so it is a value in \mathcal{IL} .

Call-by-value value preservation follows by cases on the forms of values in the source System F and target \mathcal{IL} sublanguage. Note that the η -expansion of $\eta\llbracket \text{App } e \rrbracket^{-1}$ ensures that this is always a value for types like $\sigma \rightsquigarrow \sigma'$, from which value preservation follows for $\llbracket e \rrbracket^{-1}$ for expressions of type σ which all have kind TYPE PtrR Eval^U . \square

Lemma 35 (Forward Inverse). *In call-by-value or call-by-name System F:*

- (1) $\llbracket \llbracket \tau \rrbracket \rrbracket^{-1} \triangleq \tau$, and
- (2) if $\Gamma \vdash e : \tau$ then $\llbracket \llbracket e \rrbracket \rrbracket^{-1} \triangleq e$.

PROOF. By induction on the syntax of τ and e , each case following directly from the inductive hypothesis. \square

Lemma 36 (Backward Inverse). *In the \mathcal{IL} sublanguage, for both call-by-value and call-by-name (de)compilation:*

- (1) $\llbracket \llbracket \tau \rrbracket^{-1} \rrbracket \triangleq \tau$ and $\llbracket \llbracket \sigma \rrbracket^{-1} \rrbracket \triangleq \forall \{\sigma\}$,
- (2) if $\Gamma \vdash e : \sigma$ then $\llbracket \llbracket e \rrbracket^{-1} \rrbracket =_{\beta_{\{\}} \eta_{\rightsquigarrow} \eta_V} e$,
- (3) if $\Gamma \vdash e : \tau$ then $\text{App } \llbracket \llbracket e \rrbracket^{-1} \rrbracket =_{\beta_{\{\}} \eta_{\rightsquigarrow} \eta_V} e$, and
- (4) if $\Gamma \vdash e : \tau$ and $\Gamma \vdash e$ **pass** then $\llbracket \eta\llbracket e \rrbracket^{-1} \rrbracket =_{\beta_{\{\}} \eta_{\rightsquigarrow} \eta_V} \text{Clos}^V e$.

PROOF. The first part follows directly on the syntax of τ and σ , and the remaining parts follow simultaneously by induction on the structure of the syntax of e . The cases for x and $\text{I}\#^V n$ are immediate, and the remaining cases are as follows:

- (2) $\llbracket \llbracket e e' \rrbracket^{-1} \rrbracket \triangleq \text{App } \llbracket \llbracket e \rrbracket^{-1} \rrbracket \llbracket \llbracket e' \rrbracket^{-1} \rrbracket =_{(3)} e \llbracket \llbracket e' \rrbracket^{-1} \rrbracket =_{(2)} e e'$
- (2) $\llbracket \llbracket e \sigma \rrbracket^{-1} \rrbracket =_{(2,3)} e \sigma$ is analogous to the previous case for call-by-value, and for call-by-name is: $\llbracket \llbracket e \sigma \rrbracket^{-1} \rrbracket \triangleq \llbracket \llbracket e \rrbracket^{-1} \rrbracket \llbracket \llbracket \sigma \rrbracket^{-1} \rrbracket =_{(2)} e \llbracket \llbracket \sigma \rrbracket^{-1} \rrbracket =_{(1)} e \sigma$
- (2) $\llbracket \llbracket \lambda t:\kappa.e \rrbracket^{-1} \rrbracket \triangleq \lambda t:\kappa.\llbracket \llbracket e \rrbracket^{-1} \rrbracket =_{(2)} \lambda t:\kappa.e$ Note this case only applies in call-by-name (de)compilation.
- (2) $\llbracket \llbracket \text{Clos}^V S \rrbracket^{-1} \rrbracket \triangleq \llbracket \eta\llbracket S \rrbracket^{-1} \rrbracket =_{(4)} \text{Clos}^V S$
- (3) $\text{App } \llbracket \llbracket \lambda x:\sigma.e \rrbracket^{-1} \rrbracket \triangleq \text{App } \llbracket \lambda x:\llbracket \sigma \rrbracket^{-1}.\llbracket e \rrbracket^{-1} \rrbracket \triangleq \text{App } (\text{Clos}^V \lambda x:\llbracket \llbracket \sigma \rrbracket^{-1} \rrbracket.\llbracket \llbracket e \rrbracket^{-1} \rrbracket)$
 $=_{(1,2)} \text{App } (\text{Clos}^V \lambda x:\sigma.e) =_{\beta_{\{\}} } \lambda x:\sigma.e$
- (3) $\text{App } \llbracket \llbracket \lambda t:\kappa.e \rrbracket^{-1} \rrbracket =_{(1,2)\beta_{\{\}} } \lambda t:\kappa.e$ follows analogously to the previous case. Note this case only applies in call-by-value (de)compilation.
- (3) $\text{App } \llbracket \llbracket \text{App } e \rrbracket^{-1} \rrbracket \triangleq \text{App } \llbracket \llbracket e \rrbracket^{-1} \rrbracket =_{(2)} \text{App } e$
- (4) $\llbracket \eta\llbracket \lambda x:\sigma.e \rrbracket^{-1} \rrbracket \triangleq \llbracket \lambda x:\llbracket \sigma \rrbracket^{-1}.\llbracket e \rrbracket^{-1} \rrbracket \triangleq \text{Clos}^V \lambda x:\llbracket \llbracket \sigma \rrbracket^{-1} \rrbracket.\llbracket \llbracket e \rrbracket^{-1} \rrbracket =_{(1,2)} \text{Clos}^V \lambda x:\sigma.e$

- (4) $\llbracket \eta[\lambda t:\kappa.e]^{-1} \rrbracket =_{(1,2)} \text{Clos}^Y \lambda t:\kappa.e$ follows analogously to the previous case. Note this case only applies in call-by-value (de)compilation.
- (4) Given $\Gamma \vdash e : {}^Y\{\sigma \rightsquigarrow \sigma'\}$:

$$\begin{aligned} \llbracket \eta[\text{App } e]^{-1} \rrbracket &\triangleq \llbracket \lambda x:\llbracket \sigma \rrbracket^{-1}. \llbracket e \rrbracket^{-1} x \rrbracket \triangleq \text{Clos}^Y (\lambda x:\llbracket \sigma \rrbracket^{-1}. \text{App } \llbracket e \rrbracket^{-1} x) \\ &=_{(1,2)} \text{Clos}^Y (\lambda x:\sigma. \text{App } e x) =_{\eta_{\rightsquigarrow}} \text{Clos}^Y (\text{App } e) \end{aligned}$$

- Given $\Gamma \vdash e : {}^Y\{\forall t:\kappa.\sigma\}$, $\llbracket \eta[\text{App } e]^{-1} \rrbracket =_{(1,2)\eta_V} \text{App } e$ analogously to the previous case. \square

Lemma 37 (Forward Soundness). *For any $\Gamma \vdash e_i : \tau$ in System F, if $\Gamma \vdash e_1 = e_2 : \tau$ in either call-by-value or call-by-name then $\llbracket \Gamma \rrbracket \vdash \llbracket e_1 \rrbracket = \llbracket e_2 \rrbracket : \llbracket \tau \rrbracket$ with regards to the corresponding compilation.*

PROOF. Note that, since the compilation translation is compositional, substitution commutes with compilation (i.e., $\llbracket e \rrbracket[\llbracket V \rrbracket/x] \triangleq \llbracket e[V/x] \rrbracket$ and similarly for type substitution). The equational axioms of call-by-name System F are sound w.r.t. compilation as follows:

- $(\beta_{\{\}})$

$$\begin{aligned} \llbracket (\lambda x:\tau.e) V \rrbracket &\triangleq \text{App } (\text{Clos}^Y \lambda x:\llbracket \tau \rrbracket. \llbracket e \rrbracket) \llbracket V \rrbracket \\ &=_{\beta_{\{\}}} (\lambda x:\llbracket \tau \rrbracket. \llbracket e \rrbracket) \llbracket V \rrbracket =_{\beta_{\rightsquigarrow}} \llbracket e \rrbracket[\llbracket V \rrbracket/x] \triangleq \llbracket e[V/x] \rrbracket \end{aligned}$$

- (β_V)

$$\llbracket (\lambda t.e) \tau \rrbracket \triangleq (\lambda t:\kappa. \llbracket e \rrbracket) \llbracket \tau \rrbracket =_{\beta_V} \llbracket e \rrbracket[\llbracket \tau \rrbracket/t] \triangleq \llbracket e[\tau/t] \rrbracket$$

- $(\eta_{\{\}})$

$$\llbracket \lambda x:\tau.(V x) \rrbracket \triangleq \text{Clos}^Y \lambda x:\llbracket \tau \rrbracket. (\text{App } \llbracket V \rrbracket x) =_{\eta_{\rightsquigarrow}} \text{Clos}^Y (\text{App } \llbracket V \rrbracket) =_{\eta_{\{\}}} \llbracket V \rrbracket$$

- (η_V)

$$\llbracket \lambda t.V t \rrbracket \triangleq \lambda t. \llbracket V \rrbracket t =_{\eta_V} \llbracket V \rrbracket$$

- $(name)$

$$\begin{aligned} \llbracket (\lambda x:\tau.e x) e' \rrbracket &\triangleq \text{App } (\text{Clos}^Y \lambda x:\llbracket \tau \rrbracket. \text{App } \llbracket e \rrbracket x) \llbracket e' \rrbracket \\ &=_{\beta_{\{\}}} (\lambda x:\llbracket \tau \rrbracket. \text{App } \llbracket e \rrbracket x) \llbracket e' \rrbracket =_{\eta_{\rightsquigarrow}} \text{App } \llbracket e \rrbracket \llbracket e' \rrbracket \triangleq \llbracket e e' \rrbracket \end{aligned}$$

The equational axioms of call-by-value System F are similarly sound, where the only difference is that the compilation of the β_V and η_V more closely resemble $\beta_{\{\}}$ and $\eta_{\{\}}$. \square

Lemma 38 (Backward Soundness). *In either the call-by-value or call-by-name sublanguage of \mathcal{IL} :*

- (1) *for any $\Gamma \vdash e_i : \sigma$, if $\Gamma \vdash e_1 = e_2 : \sigma$ then $\llbracket \Gamma \rrbracket^{-1} \vdash \llbracket e_1 \rrbracket^{-1} = \llbracket e_2 \rrbracket^{-1} : \llbracket \sigma \rrbracket^{-1}$, and*
- (2) *for any $\Gamma \vdash e_i : \tau$, if $\Gamma \vdash e_1 = e_2 : \tau$ then $\llbracket \Gamma \rrbracket^{-1} \vdash \eta[\llbracket e_1 \rrbracket^{-1}] = \eta[\llbracket e_2 \rrbracket^{-1}] : \llbracket \tau \rrbracket^{-1}$,*

in the corresponding System F.

PROOF. Note that, since the compilation translation is compositional, substitution commutes with decompilation (i.e., $\llbracket e \rrbracket^{-1}[\llbracket e' \rrbracket^{-1}/x] \triangleq \llbracket e[e'/x] \rrbracket^{-1}$ and similarly for type substitution). The equational axioms of the call-by-name \mathcal{IL} sublanguage sound w.r.t. decompilation as follows:

- $(\beta_{\rightsquigarrow})$

$$\llbracket (\lambda x:\tau.e) e' \rrbracket^{-1} \triangleq (\lambda x:\tau. \llbracket e \rrbracket^{-1}) \llbracket e' \rrbracket^{-1} =_{\beta_{\{\}}} \llbracket e \rrbracket^{-1}[\llbracket e' \rrbracket^{-1}/x] = \llbracket e[e'/x] \rrbracket^{-1}$$

- (β_V)

$$\llbracket (\lambda t:\kappa.e) \tau \rrbracket^{-1} \triangleq (\lambda t:\kappa. \llbracket e \rrbracket^{-1}) \llbracket \tau \rrbracket^{-1} =_{\beta_V} \llbracket e \rrbracket^{-1}[\llbracket \tau \rrbracket^{-1}/t] = \llbracket e[\tau/t] \rrbracket^{-1}$$

- ($\beta_{\{\}} \})$ Given that $\Gamma \vdash e : \sigma \rightsquigarrow \sigma'$, we have

$$\eta[\llbracket \text{App } (\text{Clos}^\vee e) \rrbracket]^{-1} \triangleq \lambda x : \llbracket \sigma \rrbracket^{-1}. \llbracket \text{Clos}^\vee e \rrbracket^{-1} x \triangleq \lambda x : \llbracket \sigma \rrbracket^{-1}. \eta[\llbracket e \rrbracket]^{-1} x =_{\eta_{\{\}}} \eta[\llbracket e \rrbracket]^{-1}$$

The case in call-by-value where $\Gamma \vdash e : \forall t : \kappa. \sigma$ is analogous.

- (η_{\rightsquigarrow})

$$\eta[\llbracket \lambda x : \sigma. (e \ x) \rrbracket]^{-1} \triangleq \lambda x : \llbracket \sigma \rrbracket^{-1}. (\llbracket e \rrbracket^{-1} x) =_{\beta_{\{\}}} \eta[\llbracket e \rrbracket]^{-1}$$

- (η_{\vee}) In call-by-name, we have

$$\llbracket \lambda t : \kappa. (e \ t) \rrbracket^{-1} \triangleq \lambda t. (\llbracket e \rrbracket^{-1} t) =_{\eta_{\vee}} \llbracket e \rrbracket^{-1}$$

whereas in call-by-value, we have an analogous equality to the previous case.

- ($\eta_{\{\}} \})$ Given that $\Gamma \vdash V : \vee \{ \sigma \rightsquigarrow \sigma' \}$

$$\llbracket \text{Clos}^\vee (\text{App } e) \rrbracket^{-1} \triangleq \eta[\llbracket \text{App } e \rrbracket]^{-1} \triangleq \lambda x : \llbracket \sigma \rrbracket^{-1}. \llbracket e \rrbracket^{-1} x =_{\eta_{\{\}}} \llbracket V \rrbracket^{-1}$$

The case in call-by-value where $\Gamma \vdash e : \vee \{ \forall t : \kappa. \sigma \}$ is analogous.

Congruence of equality follows by induction on the syntax of expressions, where the only cases that do not follow immediately from the inductive hypothesis are those which introduce $\llbracket e \rrbracket$ when $\Gamma \vdash e : \tau$. This can occur with function application and App e expansion of type $\sigma \rightsquigarrow \sigma'$, which simplify to a known case as follows:

$$\llbracket e \ e' \rrbracket^{-1} \triangleq \llbracket e \rrbracket^{-1} \llbracket e' \rrbracket^{-1} =_{\text{name}} (\lambda x : \llbracket \sigma \rrbracket^{-1}. \llbracket e \rrbracket^{-1} x) \llbracket e' \rrbracket^{-1} = \eta[\llbracket e \rrbracket]^{-1} \llbracket e' \rrbracket^{-1}$$

$$\eta[\llbracket \text{App } e \rrbracket]^{-1} \triangleq \lambda x : \llbracket \sigma \rrbracket^{-1}. \llbracket e \rrbracket^{-1} x =_{\beta_{\{\}}} \eta[\llbracket e \rrbracket]^{-1}$$

An analogous derivation is also required for polymorphic application and App e expansion in call-by-value, which both follow from the β_{\vee} axiom. \square

Corollary 6 (Equational Correspondence). *There is an equational correspondence between both call-by-name and call-by-value System F and the corresponding sublanguage of \mathcal{IL} . Namely, the following properties hold:*

- (1) If $\Gamma \vdash e_1 = e_2 : \tau$ in System F , then $\llbracket \Gamma \rrbracket \vdash \llbracket e_1 \rrbracket = \llbracket e_2 \rrbracket : \llbracket \tau \rrbracket$ in \mathcal{IL} .
- (2) If $\Gamma \vdash e_1 = e_2 : \sigma$ in \mathcal{IL} , then $\llbracket \Gamma \rrbracket^{-1} \vdash \llbracket e_1 \rrbracket^{-1} = \llbracket e_2 \rrbracket^{-1} : \llbracket \sigma \rrbracket^{-1}$ in System F .
- (3) For all $\Gamma \vdash e : \tau$ in System F , $\Gamma \vdash \llbracket \llbracket e \rrbracket \rrbracket^{-1} = e : \tau$.
- (4) For all $\Gamma \vdash e : \sigma$ in \mathcal{IL} , $\Gamma \vdash \llbracket \llbracket e \rrbracket^{-1} \rrbracket = e : \sigma$.

Corollary 7 (Soundness and Completeness). *For any $\Gamma \vdash e_i : \tau$ in either call-by-name or call-by-value System F , $\Gamma \vdash e_1 = e_2 : \tau$ if and only if $\llbracket \Gamma \rrbracket \vdash \llbracket e_1 \rrbracket = \llbracket e_2 \rrbracket : \llbracket \tau \rrbracket$.*