# Closure-Free Functional Programming in a Two-Level Type Theory

ANONYMOUS AUTHOR(S)

There are many abstraction tools in modern functional programming which heavily rely on general-purpose compiler optimization to achieve adequate performance. For example, monadic binding is a higher-order function which yields runtime closures in the absence of sufficient compile-time inlining and beta-reductions, thereby significantly degrading performance. In current systems such as the Glasgow Haskell Compiler, there is no strong guarantee that general-purpose optimization can eliminate abstraction overheads, and users only have indirect and fragile control over code generation through inlining directives and compiler options. In this paper we propose using a two-stage language to simultaneously get strong code generation guarantees and strong abstraction features. The object language is a simply typed first-order language where all function calls are statically known, and which can be compiled without runtime closures. The compile-time language is a dependent type theory. The two are integrated in a two-level type theory. We develop some abstraction tools in this setting. First, we develop monads and monad transformers. Second, we develop fusion for push and pull streams. Most of our results are also adapted to a proof-of-concept library in typed Template Haskell.

## 1 INTRODUCTION

Modern functional programming supports many convenient abstractions. These often come with significant runtime overheads. Sometimes the overheads are acceptable, but in other cases compiler optimization is crucial. Monads in Haskell is an example for the latter. Even the Reader monad, which is one of the simplest in terms of implementation, yields large overheads when compiled without optimizations. Consider the following snippet.

$$f :: \mathsf{Int} \to \mathsf{Reader\ Bool\ Int}$$
$$f\ x = \textbf{do}\ \{b \leftarrow \mathsf{ask}; \textbf{if}\ x\ \textbf{then return}\ (x + 10)\ \textbf{else return}\ (x + 20)\}$$

With optimizations enabled, GHC compiles this roughly to the following:

$$f :: \mathsf{Int} \to \mathsf{Bool} \to \mathsf{Int}$$
$$f = \lambda\ x\ b.\ \textbf{if}\ b\ \textbf{then}\ x + 10\ \textbf{else}\ x + 20$$

---

Without optimizations we roughly get:

$$f = \lambda\, x.\, (\gg\!=)\ \text{MonadReaderDict ask}\ (\lambda\, b.\ \textbf{if}\ b$$
$$\textbf{then return}\ \text{MonadReaderDict}\ (x + 10)$$
$$\textbf{else return}\ \text{MonadReaderDict}\ (x + 20))$$

Here, MonadReaderDict is a runtime dictionary, containing the methods of the Monad instance for Reader, and $(\gg\!=)$ MonadReaderDict is a field projection. Here, a runtime closure will be created for the $\lambda\, b. \ldots$ function, and $(\gg\!=)$, ask and **return** will create additional dynamic closures.

The difference between optimized and unoptimized code is already large here, and it gets even larger when we consider monad transformers or code that is polymorphic over monads. In Haskell, such code is pervasive, even in fairly basic programs which do not use fancy abstractions. The mapM function in the Haskell Prelude, which maps over a list in some monad, is a third-order and second-rank polymorphic function in disguise, because its monad dictionary argument contains the polymorphic second-order $(\gg\!=)$ method.

$$\text{mapM} :: \text{Monad}\ m \Rightarrow (a \to m\, b) \to [a] \to m\, [b]$$

Compiling mapM efficiently relies on inlining the instance dictionary, then inlining the methods contained there, and also inlining the functions that the higher-order binding method is applied to.

REORG intro, trim, contrib, motivate monadtrans & stream fusion

GHC's optimization efforts are respectable, and it has gotten quite good over its long history of development. However, there is no strong guarantee that certain optimizations will happen. Control over optimizations remains tricky, fragile and non-compositional. INLINE and REWRITE pragmas can be used to control code generation, but without any strong guarantee, and their sophisticated usage requires knowledge of GHC internals. For example, correctly specifying the *ordering* of certain rule applications is often needed. We have to also care about function arities. Infamously, the function composition operator is defined as $(.)\, f\, g = \lambda\, x \to f\, (g\, x)$ in the base libraries, instead of as $(.)\, f\, g\, x = f\, (g\, x)$, to get better inlining behavior — as explained in a source comment right next to the definition [?].

Although there are numerous tricks and idioms that are used in high-performance Haskell programming, for reliable performance it is necessary that programmers periodically *review* GHC's optimized code output. Needless to say, this is rather time-consuming and poorly scalable.

## 1.1 Staged Compilation

In this paper we use staged compilation to address issues of robustness. The idea is to shift as much as possible work from general-purpose optimization to metaprograms.

Metaprograms can be deterministic, transparent, and can be run efficiently, using fast interpreters or machine code compilation. In contrast, general-purpose optimizers are slower to run, less transparent and less reliable. Also, metaprogramming allows library authors to exploit domain-specific optimizations, while it is not realistic for general-purpose optimizers to know about all domains.

On the other hand, metaprogramming requires some additional care and input from programmers. Historically, there have been problems with ergonomics as well:

- Code generation might fail *too late* in the pipeline, producing incomprehensible errors; this is often caused by not having enough static guarantees about the well-formedness of code output.
- Tooling for the object language (debugging, profiling, IDE support) often does not work for metaprogramming. This is more likely if the object and meta layers are wildly different.

- Metaprogramming may introduce heavy noise and boilerplate, obscuring the logic of programs, or impose restrictions on how code can be structured. For instance, Template Haskell mandates that metaprograms used in splices in some module are defined in a different module.

The idea of **two-level type theory** (2LTT) is to use a highly expressive dependent type theory for compile-time computation, but generate code in possibly different, simpler object language. In this paper, we use 2LTT to sweeten the deal of staged compilation, aiming for a combination of strong guarantees, good ergonomics, high level of abstraction and easy-to-optimize code output.

We develop a particular two-level type theory for this purpose, which we call **CFTT**, short for "closure-free type theory". This consists of

- A simply-typed object theory with first-order functions, general recursion and finitary algebraic data types. This language is easy to optimize and compile downstream in the pipeline, but it lacks many convenience features.
- A standard Martin-Löf type theory for the compile-time language. This allows us to recover many features by metaprogramming.

In particular, since the object language is first-order, we guarantee that all programs in CFTT can be ultimately compiled without any dynamic closures, using only calls and jumps to statically known code locations. Why focus on closures? They are the foundation to almost all abstraction tools in functional programming:

- Higher-order functions in essentially all functional languages are implemented with closures.
- Type classes in Haskell use dictionary-passing, which relies on closures for function methods.
- Functors and first-class modules in OCaml and other ML-s rely on closures.

Hence, doing functional programming without closures is a clear demonstration that we can get rid of abstraction overheads.

It turns out that surprisingly little practical programming relies essentially on closures. Most of the time, programmers use higher-order functions for *abstraction*, such as when mapping over lists, where it is expected that the mapping function will be inlined. In other cases, the use of closures can be eliminated by small-scale defunctionalization. For example, difference lists [?] are often implemented with closures, but they can be also implemented as binary trees, with the same programming interface.

Essential closures: CPS? Threaded interpreters?

Whole-program defunctionalization is also possible, and it is notably used by the MLton compiler [?]. While this can be practical and effective, it can be also expensive and require whole-program processing. Also, conceptually speaking, it does not eliminate closures but instead makes them more transparent to optimizations. In this paper we do not use defunctionalization.

We note though that our setup is compatible with closures as well, and it can support two separate type formers for closure-based and non-closure-based ("static") functions. Having both of these would be desirable in a practical system. In the current work we focus on the closure-free case because it is much less known and developed, and it is quite interesting to see how far we can go with it.

## 1.2 Contributions

CONTRIBUTIONS

## 2 BASICS OF CFTT

In the following we give an overview of CFTT features. Here we focus on examples and informal explanations. We defer the more formal details to Sections ?? and ??. We first review the meta-level language, then the object-level one, and finally the staging operations which bridge between the two.

### 2.1 The Meta Level

**MetaTy** is the universe of types in the compile-time language. We will often use the term "metatype" to refer to inhabitants of MetaTy, and use "metaprogram" for inhabitants of metatypes. MetaTy supports dependent functions, products and indexed inductive types [?].

Formally, MetaTy is additionally indexed by universes levels, and we have $\text{MetaTy}_i : \text{MetaTy}_{i+1}$. However, this adds a bit of noise, and it is not very relevant to the current paper, so we shall omit levels. Note that universe levels have nothing to do with staged compilation; they are about sizing for the purpose of logical consistency.

> If we won't use Sigma notation, get rid of it

We use Agda-like syntax for functions and implicit arguments. A basic example:

$$id : \{A : \text{MetaTy}\} \rightarrow A \rightarrow A$$
$$id = \lambda x. x$$

Here, the type argument is implicit, and it gets inferred when we use the function. For example, $id$ True is elaborated to $id$ {Bool} True, where the braces mark an explicit application for the implicit argument. Similarly, we can introduce implicit lambdas explicitly:

$$id = \lambda \{A : \text{MetaTy}\}(x : A). x$$

We write $(a : A) \times B$ for $\Sigma$-types in MetaTy. For the field projections, for $t : (a : A) \times B$, we have fst $t : A$ and snd $t : B[a \mapsto \text{fst}\, t]$, and pairing is written simply as $(t, u)$. Additionally, we use syntactic sugar for field projections: we can write a type as $(field_1 : A) \times (field_2 : B) \times (field_3 : C)$, and write $t.field_2$ for a named field projection. The type itself denotes a right-nested $\Sigma$-type. We can also write $(t, u, v)$ for a right-nested iterated pairing. In Haskell style, we write () for the unit type, and also for its inhabitant.

Inductive types can be introduced using a Haskell-like ADT notation, or with a GADT-style one:

$$\textbf{data}\ \text{Bool}_M : \text{MetaTy} = \text{True}_M \mid \text{False}_M$$

$$\textbf{data}\ \text{Bool}_M : \text{MetaTy}\ \textbf{where}$$
$$\text{True}_M : \text{Bool}_M$$
$$\text{False}_M : \text{Bool}_M$$

Note that we added an $_M$ subscript to the type; when analogous types can be defined both on the meta and object levels, we will sometimes use this subscript to disambiguate the meta-level version. We will also use Haskell-like newtype notation, such as in **newtype** Wrap $A = $ Wrap $\{unWrap : A\}$.

Importantly, we can only produce meta-level values by induction on meta-level data. In general, all construction and elimination rules for type formers in MetaTy stay within MetaTy.

### 2.2 The Object Level

**Ty** is the universe of types in the object language. It is itself a metatype, so so we have Ty : MetaTy. Similarly as in the case of MetaTy, all construction and elimination rules of the object language stay within Ty. However, we further split Ty to two sub-universes.

First, **ValTy** : MetaTy is the universe of *value types*. ValTy supports parameterized algebraic data types, where parameters can have arbitrary types, but all constructor field types must be in ValTy.

Since ValTy is a sub-universe of Ty, we have that when $A$ : ValTy then also $A$ : Ty. Formally, this is specified as an explicit embedding operation but informally it is nicer to have an implicit subtyping.

Second, **CompTy** : MetaTy is the universe of *computation types*. This is also a sub-universe of Ty with implicit coercions. For now, we only specify that CompTy contains functions whose domains are value types:

$$- \to - : \mathsf{ValTy} \to \mathsf{Ty} \to \mathsf{CompTy}$$

For instance, if Bool : ValTy is defined as an object-level ADT, then Bool $\to$ Bool : CompTy, hence also Bool $\to$ Bool : Ty. However, (Bool $\to$ Bool) $\to$ Bool is ill-formed, since the domain is not a value type. Let us look at an example for an object-level program, where we already have natural numbers declared as **data** Nat := Zero | Suc Nat:

$$add : \mathsf{Nat} \to \mathsf{Nat} \to \mathsf{Nat}$$
$$add := \textbf{letrec } go\, n\, m := \textbf{case } n \textbf{ of}$$
$$\mathsf{Zero} \to m;$$
$$\mathsf{Suc}\, n \to \mathsf{Suc}\, (go\, n\, m);$$
$$go$$

Every recursive definition must be introduced with **letrec** . The general syntax is **letrec** $x : A := t; u$, where the $A$ type annotation can be omitted. **letrec** can be only used to define computations, not values (hence, only functions can be recursive so far).

Object-level definitions use := as notation, instead of the = that is used for meta-level ones. Later on, we may use some implicit CFTT notations, but we will always disambiguate the stages of **let**-s in this way. Non-recursive **let** is also allowed, and can be used to shadow binders:

$$f : \mathsf{Nat} \to \mathsf{Nat}$$
$$f\, x := \textbf{let } x := x + 10; \textbf{let } x := x + 20; x * 10$$

We also allow **newtype** definitions, both in ValTy and CompTy. These are assumed to be erased at runtime. In the Haskell adaptation they are important for aiding type class resolution, and we think that the explicit wrapping makes many definitions more comprehensible in CFTT as well.

Values are call-by-value at runtime; they are computed eagerly in function applications and **let**-s. **let**-definitions can be used to define inhabitants of any type, and the type of the **let** body can be also arbitrary. Additionally, the right hand sides of **case** branches can also have arbitrary types. So the following is well-formed:

$$f : \mathsf{Bool} \to \mathsf{Nat} \to \mathsf{Nat}$$
$$f\, b := \textbf{case } b \textbf{ of } \mathsf{True} \to (\lambda\, x.\, x + 10); \mathsf{False} \to (\lambda\, x.\, x * 10)$$

What about **definitional equality** at the object level? This is a distinct notion from the runtime semantics; object programs are embedded in CFTT, which is dependently typed, and we need to decide definitional equalities during type checking. The setup is simple: we have no $\beta$ or $\eta$ rules for object programs at all, nor any rule for **let**-unfolding. The reasons are the following. First, since the object language has general recursion, most $\beta$-rules are only valid in the operational semantics up to some restrictions, and $\beta$-conversion is not decidable to begin with. Second, in metaprogramming we care about the size and efficiency of generated code, and these properties are not stable under $\beta$ and $\eta$. Hence, the sensible choice is to do metaprogramming up to strict syntactic equality of object programs.

Let us discuss the object language. First, notice that there is no polymorphism or any kind of type dependency. Although we can define lists as **data** List $A$ : ValTy := Nil | Cons $A$ (List $A$),

the parameterization is just a shorthand; all concrete instantiations of the type are distinct. This monomorphism makes it easy to use different memory layouts for different types. For example, types which look like products may be unboxed. We could also make a distinction between boxed and unboxed sum types. We do not explore this in detail, we just note that monomorphic types make it easy to control memory layouts, and we believe that this is an important part of performance optimization.

Second, there are no higher-order functions, and functions also cannot be stored in data structures. Hence, locally defined functions can never escape their scope, and all function calls are to functions defined in the current scope. This makes it possible to run object programs without using dynamic closures. This latter point might not be completely straightforward; what about the previous $f$ function which has $\lambda$-expressions under a **case**, should that require closures?

We say that it should not. We make this formal formal in Section **??**; here we only give an intuitive explanation. In short, we choose a call-by-name runtime semantics for functions, which means that the only way we can compute with a function is by applying it to all arguments and extracting the resulting value. Hence, the only way to compute with $f$ is to apply it to *two* arguments, so $f$ is operationally equivalent to $\lambda\,b\,x.$ **case** $b$ **of** True $\rightarrow x + 10$; False $\rightarrow x * 10$. In Section **??** we show that all object programs can be transformed to a *saturated* form, where the arity of every function call matches the number of topmost $\lambda$-binders in the definition of the called function. Then, the local functions can be compiled either to top-level functions by lambda-lifting, or if they are only tail-called, left in place as join points [**?**].

Why not just make the object language less liberal, e.g. by disallowing $\lambda$ under **case** or **let**, thereby making call saturation easier or more obvious? There is a trade-off between making the object language more restricted, and thus easier to compile, and making *metaprogramming* more convenient. We will see that the ability to insert **let**-s without restriction is very convenient in code generation, and likewise the ability to have arbitrary object expressions in **case** bodies. In this paper we go with the most liberal object syntax, at the cost of needing more downstream processing. The call-by-name nature of computation types requires a bit of a change of thinking from programmers, but we believe that it is well worth to have for the metaprogramming convenience.

On the other side of the spectrum, one might imagine generating object code in low-level A-normal form. This is more laborious, but it could be also interesting, because it forces us to invent abstractions for manipulating ANF in the metalanguage. In a similar vein, Allais recently proposed metaprogramming quantum circuits in a two-level type theory [**?**]. We leave such setups with low-level object languages to future investigation.

Finally, one might compare our object language to call-by-push-value (CBPV). Indeed, we took inspiration from CBPV, and there are similarities, but also differences. In both systems there is a value-computation distinction, and values are call-by-value, and computations are call-by-name. However, our object language allows variable binding at arbitrary types, while CBPV only supports it at value types. In CBPV, a let-definition for a function is only possible by first packing it up as a closure value (or: "thunk"), which clearly does not work for us. Also, CBPV makes a judgment-level structural distinction between values and computations, while we use type universes for that. Generally speaking, type-based restrictions are easier to work with in dependent type theories than structural restrictions.

## 2.3 Staging

With what we have seen so far, there is no interaction between the meta and object levels. We make such interaction possible with *staging operations*.

- For $A$ : Ty, we have $\Uparrow A$ : MetaTy, pronounced as "lift $A$". This is the type of metaprograms that produce $A$-typed object programs.
- For $A$ : Ty and $t$ : $A$, we have $\langle t \rangle$ : $\Uparrow A$, pronounced "quote $t$". This is the metaprogram which immediately returns $t$.
- For $t$ : $\Uparrow A$, we have $\sim t$ : $A$, pronounced "splice $t$". This inserts the result of a metaprogram into an object term. *Notation:* splicing binds stronger than function application, so $f \sim x$ is parsed as $f (\sim x)$.
- We have $\langle \sim t \rangle \equiv t$ and $\sim \langle t \rangle \equiv t$ as definitional equalities.

A CFTT program is a mixture of object-level and meta-level top-level definitions and declarations. **Staging** means running all metaprograms in splices and inserting their output into object code, keeping all object-level top entries and discarding all meta-level ones. Thus, staging takes a CFTT program as input, and produces output which is purely in the object-level fragment, with no metatypes and metaprograms remaining. The output is guaranteed to be well-typed. This staging is formalized in detail in [?]. In Section ?? we describe the modifications to ibid. that are used in this paper.

Let us look at some basic staging examples. Recall the meta-level identity function; it can be used at the object-level too, by applying it to quoted terms:

$$\textbf{let } n : \mathsf{Nat} := \sim(id \langle 10 + 10 \rangle); \ldots$$

Here, $id$ is used at type $\Uparrow\mathsf{Nat}$. During staging, the expression in the splice is evaluated, so we get $\sim\langle 10 + 10 \rangle$, which is definitionally the same as $10 + 10$, which is our staging output here. Boolean short-circuiting is another basic use-case:

$$and : \Uparrow\mathsf{Bool} \to \Uparrow\mathsf{Bool} \to \Uparrow\mathsf{Bool}$$
$$and \, x \, y = \langle \textbf{case } \sim x \textbf{ of } \mathsf{True} \to \sim y; \mathsf{False} \to \mathsf{False} \rangle$$

Since the $y$ expression is inlined under a **case** branch at every use site, it is only computed at runtime when $x$ evaluates to True. In many situations, staging can be used instead of laziness to implement short-circuiting, and with generally better runtime performance, avoiding the overhead of thunking. Consider the *map* function now:

$$map : \{A \, B : \mathsf{ValTy}\} \to (\Uparrow A \to \Uparrow B) \to \Uparrow(\mathsf{List} \, A) \to \Uparrow(\mathsf{List} \, B)$$
$$map \, f \, as = \langle \textbf{letrec } go \, as := \textbf{case } as \textbf{ of}$$
$$\mathsf{Nil} \quad \to \mathsf{Nil};$$
$$\mathsf{Cons} \, a \, as \to \mathsf{Cons} \sim(f \langle a \rangle) \, (go \, as);$$
$$go \sim as \rangle$$

For example, this can be used as **let** $f \, as$ : List Nat $\to$ List Nat $:= \sim(map \, (\lambda \, x. \langle \sim x + 10 \rangle) \langle as \rangle)$. This is staged to a recursive definition where the mapping function is inlined into the Cons case as Cons $a \, as \to$ Cons $(a + 10) \, (go \, as)$. Note that *map* has to abstract over value types, since lists can only contain values, not functions. Also, the mapping function has type $\Uparrow A \to \Uparrow B$, instead of $\Uparrow(A \to B)$. The former type is often preferable to the latter in staging; the former is a metafunction with useful computational content, while the latter is merely a black box that computes object code. If we have $f : \Uparrow(A \to B)$, and $f$ is staged to $\langle \lambda \, x. t \rangle$, then $\sim f \, u$ is staged to an undesirable "administrative" $\beta$-redex $(\lambda \, x. t) \, u$.

## 3  MONADS & MONAD TRANSFORMERS

In this section we build a library for monads and monad transformers. We believe that this is a good demonstration of CFTT-s abilities, since monads are ubiquitous in practical Haskell programming, and they also introduce a great amount of abstraction and that should be optimized away.

### 3.1  Binding-Time Improvements

We start with some preparatory work before getting to monads. We saw that $\Uparrow A \to \Uparrow B$ is usually preferable to $\Uparrow(A \to B)$. The two types are actually equivalent up to the runtime semantics of object programs, and we can convert back and forth in CFTT:

$$up : \Uparrow(A \to B) \to \Uparrow A \to \Uparrow B \qquad down : (\Uparrow A \to \Uparrow B) \to \Uparrow(A \to B)$$
$$up\ f\ a = \langle {\sim} f\ {\sim}a \rangle \qquad\qquad down\ f = \langle \lambda\, a.\, {\sim}(f\ \langle a \rangle) \rangle$$

We can not show internally, using propositional equality, that these functions are inverses, since we do not have $\beta\eta$-rules for object functions; but we will not need this proof in the rest of the paper.

In the staged compilation and partial evaluation literature, the term **binding time improvement** is used to refer to such conversions, where the "improved" version supports more compile-time computation. A general strategy for generating efficient "fused" programs, is to try to work as much as possible with improved representations, and only convert back to object code at points where runtime control dependencies are unavoidable.

We use a Haskell-style type class for binding time improvements:

$$\textbf{class}\ \mathsf{Improve}\ (A : \mathsf{Ty})\ (B : \mathsf{MetaTy})\ \textbf{where}$$
$$up\quad : \Uparrow A \to B$$
$$down : B \to \Uparrow A$$
$$\textbf{instance}\ \mathsf{Improve}\ (A \to B)\ (\Uparrow A \to \Uparrow B)\ \textbf{where}\ \dots$$

We will use type classes in an informal way, without precisely specifying how they work. However, in the Haskell adaptation of this paper we do make crucial use of type classes. There are some necessary differences between the CFTT and Haskell versions though, which we will summarize in Section **??**. Let us look at improvement for product types now:

$$\textbf{instance}\ \mathsf{Improve}\ (A,\ B)\ (\Uparrow A,\ \Uparrow B)\ \textbf{where}$$
$$up\ x\qquad = (\langle \mathsf{fst}\ {\sim}x \rangle,\ \langle \mathsf{snd}\ {\sim}x \rangle)$$
$$down\ (a,\ b) = \langle ({\sim}a,\ {\sim}b) \rangle$$

Here we overload Haskell-style product type notation for types and pairing, at both levels (products are definable as a value type at the object level). There is a problem with this conversion though: $up$ uses $x : \Uparrow(A,\ B)$ twice, which can increase code size and duplicate runtime computations. For example, $down\ (up\ \langle f\ x \rangle)$ is staged to $\langle (\mathsf{fst}\ (f\ x),\ \mathsf{snd}\ (f\ x)) \rangle$. It would be safer to first **let**-bind an expression with type $\Uparrow(A,\ B)$, and then only use projections of the newly bound variable. This is called **let-insertion** in staged compilation. But it is impossible to use let-insertion in $up$ because the return type is a product in MetaTy, and we cannot introduce object binders in meta-level code.

### 3.2  The Code Generation Monad

Our solution is to use a monad which extends MetaTy with the ability to freely generate object-level code and introduce object binders.

$$\textbf{newtype}\ \mathsf{Gen}\ A = \mathsf{Gen}\ \{ unGen : \{R : \mathsf{Ty}\} \to (A \to \Uparrow R) \to \Uparrow R \}$$

This is a monad in MetaTy in a very standard sense, and we reuse Haskell type classes and **do** - notation.

$$\textbf{instance } \text{Monad Gen } \textbf{where}$$
$$\textbf{return } a = \text{Gen} \, (\lambda \, k. \, k \, a)$$
$$ga \ggg f = \text{Gen} \, (\lambda \, k. \, unGen \, ga \, (\lambda \, a. \, (f \, a \, k)))$$

From the Monad instance, the Functor and Monad instances can be also derived. We will reuse (<\$>) and (<\*>) for applicative notation as well. An inhabitant of Gen $A$ can be viewed as an action whose effect is to produce some object code and returns an $A$ value that can depend on produced object binders. Gen can be only "run" when the result type is an object type:

$$runGen : \text{Gen} \, (\Uparrow A) \to \Uparrow A$$
$$runGen \, ma = unGen \, ma \, id$$

We can let-bind object expressions in Gen:

$$gen : \Uparrow A \to \text{Gen} \, (\Uparrow A)$$
$$gen \, a = \text{Gen} \, (\lambda \, k. \, \langle \textbf{let } x : A := {\sim}a; {\sim}(k \, \langle x \rangle) \rangle)$$

And also recursive function definitions:

$$genRec : \{A : \text{CompTy}\} \to (\Uparrow A \to \Uparrow A) \to \text{Gen} \, (\Uparrow A)$$
$$genRec \, f = \text{Gen} \, (\lambda \, k. \, \langle \textbf{letrec } x : A := {\sim}(f \, \langle x \rangle); {\sim}(k \, \langle x \rangle) \rangle)$$

Now, using do-notation, we may write **do** $\{x \leftarrow gen \, \langle 10 + 10 \rangle; y \leftarrow gen \, \langle 20 + 20 \rangle\}; \textbf{return } \langle x + y \rangle\}$, for a Gen $(\Uparrow \text{Nat})$ action. Running this with $runGen$ yields $\langle \textbf{let } x := 10 + 10; \textbf{let } y := 20 + 20; x + y \rangle$. We can also define a "safer" binding-time improvement for products, using let-insertion:

$$\textbf{instance } \text{Improve } (A, \, B) \, (\text{Gen} \, (\Uparrow A, \, \Uparrow B)) \, \textbf{where}$$
$$up \, x \quad = \textbf{do} \, \{x \leftarrow gen \, x; \textbf{return } (\langle \text{fst } {\sim}x \rangle, \, \langle \text{snd } {\sim}x \rangle)\}$$
$$down \, x = runGen \, \$ \, \textbf{do} \, \{(a, \, b) \leftarrow x; \textbf{return } \langle ({\sim}a, \, {\sim}b) \rangle\}$$

Working in Gen is convenient, since we can freely generate object code and also have access to the full metalanguage. Also, the whole point of staging is that *eventually* all metaprograms will be used for the purpose of code generation, so ultimately we always want to *runGen* our actions. So why not just always work in Gen? We may not want to do that, because we would like to reason about the size and makeup of the code we are generating, and the implicit code generation in Gen is not very transparent. This is a bit similar to the IO monad in Haskell, where eventually everything needs to run in IO, but we may prefer to not write most of our program in IO.

### 3.3 Monads

Let us start with the Maybe monad. We have **data** Maybe $A :=$ Nothing | Just $A$, and Maybe itself is available as a ValTy $\to$ ValTy metafunction. However, we cannot directly fashion a monad out of Maybe, since we do not have enough type formers in ValTy. We could try to use the following type for binding:

$$({\ggg}) : \Uparrow (\text{Maybe} \, A) \to (\Uparrow A \to \Uparrow (\text{Maybe} \, B)) \to \Uparrow (\text{Maybe} \, B)$$

This works, but the definition necessarily uses runtime case splits on Maybe values, many of which could be optimized away during staging. Also, not having a "real" monad is inconvenient for the purpose of code reuse.

Instead, our strategy is to only use proper monads in MetaTy, and convert between object types and meta-monads when necessary, as a form of binding-time improvement. Assume that $\text{Maybe}_M$ is a meta-level monad, and $\text{MaybeT}_M$ is the monad transformer with the standard definition:

$$\textbf{newtype } \text{MaybeT}_M \, M \, A = \text{MaybeT}_M \, \{runMaybeT_M : M \, (\text{Maybe}_M \, A)\}$$

Then the object-level Maybe is binding-time improved as follows:

$$\textbf{instance } \text{Improve} \, (\text{Maybe} \, A) \, (\text{MaybeT}_M \, \text{Gen} \, (\Uparrow A)) \, \textbf{where}$$

$$up \, x = \text{MaybeT}_M \, \$ \, \text{Gen} \, \$ \, \lambda k.$$

$$\langle \textbf{case } {\sim}x \textbf{ of } \text{Nothing} \to {\sim}(k \, \text{Nothing}_M); \text{Just } a \to {\sim}(k \, (\text{Just}_M \langle a \rangle)) \rangle$$

$$down \, (\text{MaybeT}_M \, (\text{Gen} \, ma)) =$$

$$ma \, (\lambda x. \textbf{case } x \textbf{ of } \text{Nothing}_M \to \langle \text{Nothing} \rangle; \text{Just}_M \, a \to \langle \text{Just } {\sim}a \rangle)$$

With this, we get the Monad instance for free from $\text{MaybeT}_M$ and Gen. A small example:

$$\textbf{let } n : \text{Maybe Nat} := {\sim}(down \, \$ \, \textbf{do } \{x \leftarrow \textbf{return } \langle 10 \rangle; y \leftarrow \textbf{return } \langle 20 \rangle; \textbf{return } \langle x + y \rangle)\}); \, \dots$$

Since $\text{MaybeT}_M$ is meta-level, its binding always fully computes at staging time. Thus, the above code is staged to

$$\textbf{let } n : \text{Maybe Nat} := \text{Just} \, (10 + 20); \, \dots$$

Assume also a $lift : \text{Monad } M \Rightarrow M \, A \to \text{MaybeT}_M \, M \, A$ operation, as in Haskell. We shall assume this for all monad transformers henceforth. We can do let-insertion in $\text{MaybeT}_M$ Gen by simply lifting:

$$gen' : \Uparrow A \to \text{MaybeT}_M \, \text{Gen} \, (\Uparrow A)$$

$$gen' \, a = lift \, (gen \, a)$$

However, it is more convenient to proceed in the style of Haskell's monad transformer library [?], and have a class for monads that can do code generation:

$$\textbf{class } \text{Monad } M \Rightarrow \text{MonadGen } M \textbf{ where}$$

$$liftGen : \text{Gen } A \to M \, A$$

$$\textbf{instance } \text{MonadGen Gen } \textbf{where } liftGen = id$$

$$\textbf{instance } \text{MonadGen } M \Rightarrow \text{MonadGen } (\text{MaybeT}_M \, M) \textbf{ where } liftGen = lift \circ liftGen$$

We also redefine $gen$ and $genRec$ to work in any MonadGen, so from now on we have:

$$gen \quad : \text{MonadGen } M \Rightarrow \Uparrow A \to M \, (\Uparrow A)$$

$$genRec : \text{MonadGen } M \Rightarrow (\Uparrow A \to \Uparrow A) \to M \, (\Uparrow A)$$

*3.3.1 Case splitting in monads.* Trying to write more substantial code, we quickly bump into an issue. We want to case-split on object-level data inside a monadic action, like in the following:

$$f : \text{Nat} \to \text{Maybe Nat}$$

$$f \, x := {\sim}(down \, \$ \, \textbf{case } x == 10 \textbf{ of}$$

$$\text{True} \; \to \textbf{return } \langle x + 5 \rangle;$$

$$\text{False} \to \text{Nothing}_M)$$

This is ill-typed as written, since we cannot compute meta-level actions from an object-level case split. Fortunately, with a little bit more work, case splitting on object values is actually possible in any MonadGen, on any value type.

We demonstrate this for lists first. An object-level **case** on lists introduces two points where code generation can continue. We define a metatype which gives us a "view" on these points:

$$\textbf{data } \text{SplitList } A = \text{Nil}' \mid \text{Cons}' (\Uparrow A) (\Uparrow(\text{List } A))$$

We can generate code for a case split, returning a view on it:

$$split : \Uparrow(\text{List } A) \rightarrow \text{Gen } (\text{SplitList } A)$$

$$split \; as = \text{Gen } \$ \; \lambda k. \; \langle \textbf{case } {\sim} as \textbf{ of } \text{Nil} \rightarrow {\sim}(k \; \text{Nil}'); \text{Cons } a \; as \rightarrow {\sim}(k \; (\text{Cons}' \; \langle a \rangle \; \langle as \rangle)) \rangle$$

Now, in any MonadGen, we may write

$$\textbf{do } \{ sp \leftarrow liftGen \; (split \; as); (\textbf{case } sp \textbf{ of } \text{Nil}' \rightarrow ...; \text{Cons}' \; a \; as \rightarrow ...) \}$$

This can be generalized to any splitting on object values. In the supplemental Haskell implementation, we overload *split* with a type class. However, in a native implementation of CFTT it might make sense to simply extend **do** -notation with an elaboration of **case** on $\Uparrow A$-typed data to an application of the appropriate *split* function, so we can simply write a **case** *as* **of** ... in a **do** -block. We adopt this in the rest of this paper.

## 3.4 Monad Transformers

At this point, it makes sense to aim for a monad transformer library where binding-time improvement is defined compositionally, by recursion on the transformer stack. The base case is the following:

$$\textbf{newtype } \text{Identity } A := \text{Identity } \{ runIdentity : A \}$$

$$\textbf{instance } \text{Improve } (\text{Identity } A) (\text{Gen } (\Uparrow A)) \textbf{ where}$$

$$up \; x = \text{Gen } \$ \; \lambda k. \; k \; \langle runIdentity \; {\sim}x \rangle$$

$$down \; x = unGen \; x \; \$ \; \lambda a. \; \langle \text{Identity } {\sim}a \rangle$$

We also switch from the object-level Maybe to a "transformer" form as well, and recover it as MaybeT Identity.

$$\textbf{newtype } \text{MaybeT } M \; A := \text{MaybeT } \{ runMaybeT : M \; (\text{Maybe } A) \}$$

Generally, we want to improve an operator $F : \text{ValTy} \rightarrow \text{ValTy}$ to some $M : \text{MetaTy} \rightarrow \text{MetaTy}$, where $M$ is a MonadGen. We define a *quantified class constraint* [?] for this, that we will often use as an assumption:

$$\text{ImproveF } F \; M = \text{MonadGen } M \times ((A : \text{ValTy}) \rightarrow \text{Improve } (F \; A) (M \; (\Uparrow A)))$$

With this, improvement can be generally defined for MaybeT:

$$\textbf{instance } \text{ImproveF } F \; M \Rightarrow \text{Improve } (\text{MaybeT } F \; A) (\text{MaybeT}_M \; M \; (\Uparrow A)) \textbf{ where}$$

$$up \; x = \text{MaybeT}_M \; \$ \; \textbf{do}$$

$$ma \leftarrow up \; \langle runMaybeT \; {\sim}x \rangle$$

$$\textbf{case } ma \textbf{ of } \text{Nothing} \rightarrow \textbf{return } \text{Nothing}_M$$

$$\qquad\qquad \text{Just } a \quad \rightarrow \textbf{return } (\text{Just}_M \; a)$$

$$down \; (\text{MaybeT}_M \; x) = \langle \text{MaybeT } {\sim}(down \; \$ \; x \ggg \lambda x. \textbf{ case } x \textbf{ of}$$

$$\text{Nothing}_M \rightarrow \textbf{return } \langle \text{Nothing} \rangle$$

$$\text{Just}_M \; a \quad \rightarrow \textbf{return } \langle \text{Just } {\sim}a \rangle) \rangle$$

In the **case** in *up*, we already use our syntax sugar for matching on a Maybe inside an M action. This is legal, since we know from the ImproveF $F \; M$ assumption that $M$ is a MonadGen.

In the meta level, we can reuse essentially all definitions from Haskell's monad transformer library mtl. The additional work that we have to do in the CFTT setting:

- Adding MonadGen instances, but these are trivially defined as *liftGen = lift ∘ liftGen* everywhere.
- Defining the object-level counterparts of monad transformers, and also the Improve instances.

From mtl, only the continuation monad transformer fails to support binding-time-improvement in CFTT, because of the obvious need for dynamic closures. In the following we present only StateT and ReaderT. Starting with StateT, we assume $\text{StateT}_M$ as the standard meta-level definition. The object-level StateT has type $(S : \text{ValTy})(F : \text{ValTy} \to \text{ValTy})(A : \text{ValTy}) \to \text{ValTy}$; the state parameter $S$ has to be a value type, since it is an input to an object-level function.

> **instance** ImproveF $F\ M$ ⇒ Improve (StateT $S\ F\ A$) ($\text{StateT}_M$ ($⇑S$) $M$ ($⇑A$)) **where**
>
>     $up\ x = \text{StateT}_M\ \$\ \lambda\ s.\ \textbf{do}$
>
>       $as \leftarrow up\ \langle runStateT \sim x \sim s \rangle$
>
>       **case** $as$ **of** $(a,\ s) \to \textbf{return}\ (a,\ s)$
>
>     $down\ x = \langle \text{StateT}\ (\lambda\ s.\ \sim(down\ \$\ \textbf{do}$
>
>       $(a,\ s) \leftarrow runStateT_M\ x\ \langle s \rangle$
>
>       $\textbf{return}\ \langle (\sim a,\ \sim s) \rangle)))\rangle$

Like before in MaybeT, we rely on object-level case splitting in the definition of *up*. For Reader, the environment parameter also has to be a value type, and we define improvement as follows.

> **instance** ImproveF $F\ M$ ⇒ Improve (ReaderT $R\ F\ A$) ($\text{ReaderT}_M$ ($⇑R$) $M$ ($⇑A$)) **where**
>
>     $up\ x = \text{ReaderT}_M\ \$\ \lambda\ r.\ up\ \langle runReaderT \sim x \sim r \rangle$
>
>     $down\ x = \langle \text{ReaderT}\ (\lambda\ r.\ \sim(down\ (runReaderT_M\ x\ \langle r \rangle))))\rangle$

*3.4.1 State and Reader operations.* Now, if we try to use the *modify* function that we already have for State$_M$, a curious thing happens. The meaning of *modify* $(\lambda\ x.\ \langle \sim x + \sim x \rangle)$ is to replace the current state $x$, as an object expression, with the expression $\langle \sim x + \sim x \rangle$, and this happens at staging time. This behaves as an "inline" modification which replaces every subsequent mention of the state with a different expression. For instance, ignoring newtype wrappers for now,

$$down\ \$\ \textbf{do}\ \{modify\ (\lambda\ x.\ \langle \sim x + \sim x \rangle);\ modify\ (\lambda\ x.\ \langle \sim x + \sim x \rangle);\ \textbf{return}\ \langle () \rangle\}$$

is staged to

$$\langle \lambda\ x.\ ((),\ (x + x) + (x + x)) \rangle$$

which duplicates the evaluation of $x + x$. The solution is to force the evaluation of the new state in the object language, by let-insertion. A similar phenomenon happens with the *local* function in Reader. So we define "stricter" versions of these operations. We also return $⇑()$ from actions instead of () − the former is more convenient, because it can be immediately lowered with *down*.

> $put' : (\text{MonadState}\ (⇑S)\ M,\ \text{MonadGen}\ M) \Rightarrow ⇑S \to M\ (⇑())$
>
> $put'\ s = \textbf{do}\ \{s \leftarrow gen\ s;\ put\ s;\ \textbf{return}\ \langle () \rangle\}$
>
>
> $modify' : (\text{MonadState}\ (⇑S)\ M,\ \text{MonadGen}\ M) \Rightarrow (⇑S \to ⇑S) \to M\ (⇑())$
>
> $modify'\ f = \textbf{do}\ \{s \leftarrow get;\ put'\ (f\ s)\}$

$$local' : (\text{MonadReader } (\Uparrow R)\ M,\ \text{MonadGen } M) \Rightarrow (\Uparrow R \to \Uparrow R) \to M\,A \to M\,A$$

$$local'\ f\ ma = \textbf{do}\ \{r \leftarrow ask; r \leftarrow gen\,(f\ r); local\,(\lambda\_.\ r)\ ma\}$$

Now,

$$down\,\$\,\textbf{do}\ \{modify'\,(\lambda x.\,\langle{\sim}x + {\sim}x\rangle); modify'\,(\lambda x.\,\langle{\sim}x + {\sim}x\rangle)\}$$

is staged to

$$\langle\lambda x.\,\textbf{let}\ x := x + x; \textbf{let}\ x := x + x; ((),\ x)\rangle$$

## 3.5 Joining Control Flow in Monads

There is a deficiency in our library so far. Assuming $b$ : Bool, consider:

$$down\,\$\,\textbf{do}$$
$$\textbf{case}\ b\ \textbf{of}$$
$$\text{True}\ \to put'\,\langle 10\rangle$$
$$\text{False}\ \to put'\,\langle 20\rangle$$
$$modify'\,(\lambda x.\langle{\sim}x + 10\rangle)$$

The $modify'\,(\lambda x.\langle{\sim}x + 10\rangle)$ action gets inlined in both case branches during staging! This follows from the definition of monadic binding in Gen and the *split* function in the desugaring of **case**. Code generation is continued in two branches with the same action. If we have multiple **case** splits sequenced after each other, that yields exponential code size. Right now, we can fix this by let-binding the problematic action:

$$down\,\$\,\textbf{do}$$
$$x \leftarrow gen\,\$\,down\,\$\,\textbf{case}\ b\ \textbf{of}$$
$$\text{True}\ \to put'\,\langle 10\rangle$$
$$\text{False}\ \to put'\,\langle 20\rangle$$
$$up\ x$$
$$modify'\,(\lambda x.\langle{\sim}x + 10\rangle)$$

This removes code duplication by round-tripping through an object-level **let**. However, while this solution is fairly good in a state monad, where we can reasonably expect that product types can be unboxed by the object compiler, it introduces unnecessary runtime constructors for Maybe and other monads containing sum types. For Maybe, *down* introduces a runtime Just or Nothing, and *up* introduces a runtime case split. A better solution would be to introduce two let-bound *join points* before the offending **case**, one for returning a Just and one for returning Nothing, but fusing away the actual runtime constructors. So, we would like to produce object code for a **case** which looks like the following:

$$\textbf{let}\ joinJust\ n := \dots$$
$$\textbf{let}\ joinNothing\,() := \dots$$
$$\textbf{case}\ x == 10\ \textbf{of}$$
$$\text{True}\ \to joinJust\,(x + 10)$$
$$\text{False}\ \to joinNothing\,()$$

Such fused returns are possible whenever we have a Gen $A$ action at the bottom of the transformer stack, such that $A$ is isomorphic to a finite meta-level sum of value types. Recall that Gen $A$ is defined as $\{R : \text{ValTy}\} \to (A \to \Uparrow R) \to \Uparrow R$. Here, the $A \to \Uparrow R$ continuation represents all possible

code points where we can return to, and if $A$ is a finite sum, we can rearrange $A \to \Uparrow R$ to a finite product of functions. Hence, we need a modest amount of generic programming with sums of values and products of computations.

$$\textbf{data } \mathsf{SumVS} : \mathsf{List\ ValTy} \to \mathsf{MetaTy\ } \textbf{where} \dots$$

$$\textbf{data } \mathsf{ProdCS} : \mathsf{List\ CompTy} \to \mathsf{MetaTy\ } \textbf{where} \dots$$

We use a class for types that have a SumVS representation:

$$\textbf{class } \mathsf{IsSumVS}\ (A : \mathsf{MetaTy})\ \textbf{where}$$

$$Rep \quad : \mathsf{List\ ValTy}$$

$$encode : A \to \mathsf{SumVS}\ Rep$$

$$decode : \mathsf{SumVS}\ Rep \to A$$

We also define the mentioned isomorphism for Gen continuations.

$$tabulate : (\mathsf{SumVS}\ As \to \Uparrow B) \to \mathsf{ProdCS}\ (map\ (\lambda\ A.\ A \to B)\ As)$$

$$index \quad : \mathsf{ProdCS}\ (map\ (\lambda\ A.\ A \to B)\ As) \to (\mathsf{SumVS}\ As \to \Uparrow B)$$

We omit the definitions of *tabulate*, *index* and the IsSumVS instances, since they amount to standard generic programming; see e.g. [?]. The last basic definition we need is to let-bind all computations in a generic product:

$$genProdCS : \mathsf{ProdCS}\ As \to \mathsf{Gen}\ (\mathsf{ProdCS}\ As)$$

$$genProdCS\ \mathsf{Nil} \qquad = \textbf{return}\ \mathsf{Nil}$$

$$genProdCS\ (\mathsf{Cons}\ c\ cs) = \mathsf{Cons}\ \texttt{<\$>}\ gen\ c\ \texttt{<*>}\ genProdCS\ cs$$

We introduce a class for monads that support control flow joining.

$$\textbf{class } \mathsf{MonadJoin}\ M\ \textbf{where}\ join : \mathsf{IsSumVS}\ A \Rightarrow M\ A \to M\ A$$

The most interesting instance is for Gen.

$$\textbf{instance } \mathsf{MonadJoin\ Gen\ } \textbf{where}$$

$$join\ ma = \mathsf{Gen}\ \$\ \lambda\ k.\ runGen\ \$\ \textbf{do}$$

$$joinpoints \leftarrow genProdCS\ (tabulate\ (k \circ decode))$$

$$a \leftarrow ma$$

$$\textbf{return}\ \$\ index\ joinpoints\ (encode\ a)$$

Here we first convert $k : A \to \Uparrow R$ to a product of join points and let-bind each of them. Then we adjust the return code points of *ma* by looking up the corresponding continuation from the list of join points, using *index*. We also need MonadJoin instances for the other transformers in our library, but these are all essentially just lifting *join*. In $\mathsf{StateT_M}$ we must have IsSumVS for $S$ because $S$ appears in the transformer stack's return value.

$$\textbf{instance }\ (\mathsf{MonadJoin}\ M) \Rightarrow \mathsf{MonadJoin}\ (\mathsf{MaybeT_M}\ M)\ \textbf{where}$$

$$join\ (\mathsf{MaybeT_M}\ ma) = \mathsf{MaybeT_M}\ (join\ ma)$$

$$\textbf{instance }\ (\mathsf{MonadJoin}\ M) \Rightarrow \mathsf{MonadJoin}\ (\mathsf{ReaderT_M}\ R\ M)\ \textbf{where}$$

$$join\ (\mathsf{ReaderT_M}\ ma) = \mathsf{ReaderT_M}\ (join \circ ma)$$

$$\textbf{instance }\ (\mathsf{MonadJoin}\ M, \mathsf{IsSumVS}\ S) \Rightarrow \mathsf{MonadJoin}\ (\mathsf{StateT_M}\ S\ M)\ \textbf{where}$$

$$join\ (\mathsf{StateT_M}\ ma) = \mathsf{StateT_M}\ (join \circ ma)$$

687 Now, whenever the return value of a **case**-splitting action is a sum of values (this includes just
688 returning an object value, which is by far the most common situation), we can use *join* $ **case** *x* **of** ...
689 to block inlining into the branches, without creating any runtime sum types.

690 bigger example

## 3.6 Discussion

So far, we have a monad transformer library with the following features:

- Almost all definitions from the well-known Haskell ecosystem of monads and monad transformers can be directly reused, in the meta level.
- We can quite liberally pattern match on object-level values in monadic code, insert object-level **let**-s with *gen* and control code size by *join*-s.
- In monadic code, object-level data constructors are only ever created by *down*, and switching on object-level data is only created by *split* and *up*. Everything else is fully fused, and all function calls can be compiled to statically known saturated calls.

As to potential weaknesses, first, the system as described in this section has quite a bit of syntactic noise and requires extra attention from programmers. We believe that the code noise can be mitigated very effectively in a native CFTT implementation. Kovács [?] demonstrated in a prototype elaborator that almost all quotes and splices are unambiguously inferable in 2LTT programs, if we make the concession that stages of let-definitions are always specified (as we do here). Moreover, *up* and *down* should be also effectively inferable, using bidirectional elaboration, but it would make sense to limit inferred binding-time improvement to monads only. With stage inference and binding-time improvement inference, code in CFTT would look only modestly more complicated than in Haskell, and we would get far more robust and controllable code generation.

Second, in CFTT we cannot store most monadic actions in runtime data structures, those that are implemented using functions, nor can we have functions in State state, or in Reader environments. However:

- Essential usage of this is rare in practice.
- It is possible to extend CFTT with a closure type former, converting computations to values, in which case there is no such limitation anymore. Here, closure-freedom is still available; we can choose to use or abstain from the explicit closure type.

## 3.7 Haskell Implementation

We implemented all features in this section in typed Template Haskell [?], using only the strongly typed interface. The organization of the code is essentially the same as here, but some changes are necessary:

- TH cannot properly handle scoping of type variables across quotations, generating ill-typed code in some cases. It is fairly easy to work around this using term-level proxy values.
- We do not bother tracking value and computation types. This could be done, using dummy type classes, but we would not get any extra guarantees from it in downstream compilation.
- When we take the product of generic sums in CFTT, we simply take the object-level product of component types, pairwise. We do this because we assume that the downstream compiler can represent products without overheads, in an unboxed manner, and also erase unit types. In Haskell we do not rely on this, and instead use a slightly more complicated sum-of-products-of-values representation, and also fuse away product types during staging.
- We make heavy use of *singleton types* [?] in generic sums-of-products, to emulate dependent types. This adds significant noise compared to a native dependently typed version.

- We have two separate type classes for binding-time improvement, one for base types and one for monads specifically. This is because the quantified constraints in **??** yield poor type inference for *up* and *down*, by preventing *functional dependencies* [**?**] between the class parameters. We expect that in a native CFTT implementation, the extra staging information in type universes would help inference in this case.

## 4 STREAM FUSION

Stream fusion refers to a collection of techniques for generating efficient code from declarative definitions involving streams of values, where intermediate data structures are eliminated. For example, mapping twice over a list produces an intermediate list in a naive implementation but does not do so in a fused implementation.

Fusion techniques can be broadly grouped into *push* fusion, which is based on Church-encodings of inductive lists, and *pull* fusion, which is based on Church-encodings of coinductive lists, which manifests as *state machines*. The two styles have different trade-offs, and in practical programming it is a good idea to support both, but in this section we focus on pull fusion.

The reason is partly lack of space, but also that pull fusion has been historically the more difficult to reliably compile, and we can demonstrate significant improvements in CFTT, in robustness and expressive power. Our fusion always goes through, and we also support a *concatMap* operation, which has been elusive in guaranteed fusion. We believe that our implementation also provides valuable conceptual clarity.

### 4.1 Streams

A pull stream is a meta-level specification of a state machine:

> **data** Step $S\,A$ = Stop | Yield $A\,S$ | Skip $S$
>
> **data** Pull $(A : \text{MetaTy}) : \text{MetaTy}$ **where**
>
> MkPull : $(S : \text{MetaTy}) \to \text{IsSumVS}\,S \Rightarrow S \to (S \to \text{Gen}\,(\text{Step}\,S\,A)) \to \text{Pull}\,A$

In MkPull, $S$ is the type of the internal state which is required to be a finite sum of values by the IsSumVS $S$ constraint. The next $S$-typed field is the initial state, while $S \to \text{Gen}\,(\text{Step}\,S\,A)$ represents all transitions. Possible transitions are stopping (Stop), transitioning to a new state while outputting a value (Yield) and making a silent transition to a new state (Skip).

## 5 TODOS

Introduce newtypes accordingly

allow arbitrary parameters (including ValTy -> ValTy functions) for object ADTs

extend as I have more stuff

FIGURE OUT: monadic tailrec

FIGURE OUT: Traversals

DEVELOP: mutual letrec based on computational product types

## REFERENCES