

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

ANONYMOUS AUTHOR(S)

Many abstraction tools in functional programming rely heavily 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. We propose a two-stage language to simultaneously get strong guarantees about code generation and strong abstraction features. The object language is a simply-typed first-order language 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 demonstrate two applications of the system. First, we develop monads and monad transformers. Here, abstraction overheads are eliminated by staging and we can reuse almost all definitions from the existing Haskell ecosystem. Second, we develop pull-based stream fusion. Here we make essential use of dependent types to give a concise definition of a `concatMap` operation with guaranteed fusion. We provide an Agda implementation and a typed Template Haskell implementation of these developments.

CCS Concepts: • **Theory of computation** → **Type theory**; • **Software and its engineering** → **Source code generation**.

Additional Key Words and Phrases: two-level type theory, staged compilation

ACM Reference Format:

Anonymous Author(s). 2024. Closure-Free Functional Programming in a Two-Level Type Theory. 1, 1 (February 2024), 22 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

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:

```
f :: Int → Reader Bool Int
f x = do { b ← ask; if x then return (x + 10) else return (x + 20) }
```

With optimizations enabled, GHC compiles this roughly to the code below:

```
f :: Int → Bool → Int
f = λ x b. if b then x + 10 else x + 20
```

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM XXXX-XXXX/2024/2-ART

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

Without optimizations we roughly get:

```
f = λ x. (≫) MonadReaderDict ask (λ b. if b
    then return MonadReaderDict (x + 10)
    else return MonadReaderDict (x + 20))
```

Here, `MonadReaderDict` is a runtime dictionary, containing the methods of the `Monad` instance for `Reader`, and `(≫)` `MonadReaderDict` is a field projection. Here, a runtime closure will be created for the `λ b. ...` function, and `(≫)`, `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. Consider the `mapM` function from the Haskell Prelude:

```
mapM :: Monad m => (a -> m b) -> [a] -> m [b]
```

This is a third-order and second-rank polymorphic function in disguise, because its monad dictionary argument contains the polymorphic second-order `(≫)` method. 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.

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 advanced usage requires knowledge of GHC internals. For example, correctly specifying the *ordering* of certain rule applications is often needed. We also have to care about formal function arities. Infamously, the function composition operator is defined as `(.) f g = λ x -> f (g x)` in the base libraries, instead of as `(.) f g x = f (g x)`, to get better inlining behavior [?]. It remains a common practice in high-performance Haskell programming to manually review GHC's optimized code output.

1.1 Closure-Free 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 robust. 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. In weakly-typed staged systems, code generation might fail *too late* in the pipeline, producing incomprehensible errors. Or, tooling that works for an object language (like debugging, profiling, IDEs) may not work for metaprogramming, or metaprogramming may introduce heavy noise and boilerplate, obscuring the logic of programs and imposing restrictions on code structure.

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 in the downstream 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.

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. 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.

Perhaps 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. 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 interesting to see how far we can go with it.

1.2 Contributions

- In Section ?? we present the two-level type theory CFTT, where the object language is first-order simply-typed and the metalanguage is dependently typed. The object language admits an operational semantics without runtime closures, and can be compiled in a way such that all function calls are statically known. We provide a supplementary Agda formalization of the operational semantics.
- In Section ?? we build a monad transformer library. We believe that this is a good demonstration, because monads and monad transformers are the most widely used effect system in Haskell, and at the same time their compilation to efficient code can be surprisingly difficult. The main idea is to have monads and monad transformers *only* in the metalanguage, and try to express as much as possible at the meta level, and funnel the object-meta interactions through specific *binding-time improvements*. We show how almost all definitions from the Haskell ecosystem of monad transformers can be reused in the meta level of CFTT. We also develop let-insertion, join points and case switching on object-level data, in the generality of monad transformers.
- In Section ?? we build a pull-based stream fusion library. Here, the motivation is to demonstrate essential usage of dependent types to concisely solve a difficult problem, namely guaranteed fusion for unrestricted combinations of `concatMap` and `zip`. We use a state machine representation that is based on *sums-of-products* of object-level values. We show that CFTT is compatible with a *generativity* axiom, which internalizes the fact that metaprograms cannot inspect the structure of object-level terms. We use this to show that the universe of sums-of-products is closed under Σ -types. This in turn enables a very concise definition of `concatMap`.
- We adapt the contents of the paper to a Template Haskell, with some modifications, simplifications and fewer guarantees about generated code. In particular, Haskell does not have

enough dependent types for the simple `concatMap` definition, but we can still work around this limitation. We also provide a more faithful Agda embedding of CFTT and our libraries. Here, the object theory is embedded as a collection of postulated operations, and we can use Agda’s normalization command to print out generated object code.

2 OVERVIEW OF CFTT

In the following we give an overview of CFTT features. We focus on informal explanations and examples here. 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 are orthogonal to staging; they are instead used to increase the strength of the system while maintaining logical consistency.

Throughout this paper we use a mix of Agda and Haskell syntax for CFTT. Dependent functions and implicit arguments follow Agda. A basic example:

$$\begin{aligned} \text{id} &: \{A : \text{MetaTy}\} \rightarrow A \rightarrow A \\ \text{id} &= \lambda x. x \end{aligned}$$

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

$$\begin{aligned} &\text{data Bool}_M : \text{MetaTy where} \\ \text{data Bool}_M : \text{MetaTy} &= \text{True}_M \mid \text{False}_M & \text{True}_M &: \text{Bool}_M \\ & & \text{False}_M &: \text{Bool}_M \end{aligned}$$

Note that we added a $_M$ 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 use Haskell-like newtype notation, such as in `newtype Wrap A = Wrap {unwrap : A}`, and also a similar notation for (dependent) record types, for instance as in

$$\text{data Record} = \text{Record} \{ \text{field1} : A, \text{field2} : B \}.$$

All construction and elimination rules for type formers in **MetaTy** stay within **MetaTy**. For example, induction on meta-level values can only produce meta-level values.

2.2 The Object Level

Ty is the universe of types in the object language. It is itself a metatype, so so we have $\text{Ty} : \text{MetaTy}$. 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 $A : \text{ValTy}$ then also $A : \text{Ty}$. Formally, this is specified as an explicit embedding operation, but we will use implicit subtyping for convenience.

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:

$$- \rightarrow - : \text{ValTy} \rightarrow \text{Ty} \rightarrow \text{CompTy}$$

For instance, if $\text{Bool} : \text{ValTy}$ is defined as an object-level ADT, then $\text{Bool} \rightarrow \text{Bool} : \text{CompTy}$, hence also $\text{Bool} \rightarrow \text{Bool} : \text{Ty}$. However, $(\text{Bool} \rightarrow \text{Bool}) \rightarrow \text{Bool}$ 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 $\text{data Nat} := \text{Zero} \mid \text{Suc Nat}$:

```
add : Nat → Nat → Nat
add := letrec go n m := case n of
  Zero → m;
  Suc n → 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. We also have non-recursive **let**, which can be used to define computations and values as well, and can be used to shadow binders:

```
f : Nat → Nat
f x := let x := x + 10; 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 implementation they are important for guiding 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 : Bool → Nat → Nat
f b := case b of True → (λ x. x + 10); False → (λ x. x * 10)
```

In contrast, computations are call-by-name, and the only way we can compute with functions is to apply them to value arguments. The call-by-name strategy is fairly benign here and does not lead to significant duplication of computation, because functions cannot escape their scope; they cannot be passed as arguments or stored in data constructors. This makes it possible to run object programs without using dynamic closures. This point is not completely straightforward; consider the previous f function which has λ -expressions under a **case**.

However, the call-by-name semantics lets us transform f to $\lambda b x. \text{case } b \text{ of True} \rightarrow x + 10; \text{False} \rightarrow x * 10$, and more generally we can transform programs so that every function call is *saturated*. This means that every function call is of the form $f t_1 t_2 \dots t_n$, where f is a function variable and the definition of f immediately λ -binds n arguments. We do not detail this here. We provide formal syntax and operational semantics of the object language in the Agda supplement. We formalized the specific translation steps that are involved in call saturation, but only specified the full translation informally.

2.2.1 Object-level definitional equality. This is a distinct notion from runtime semantics. Object programs are embedded in CFTT, which is a dependently typed language, so we may need to decide definitional equality of object programs during type checking. The setup is simple: we have no β or η rules for object programs at all, nor any rule for **let**-unfolding. The main reason is the following: we care about the size and efficiency of generated code, and these properties are not stable under $\beta\eta$ -conversion and **let**-unfolding. Moreover, since the object language has general recursion, we do not have a sensible and decidable notion of program equivalence anyway.

2.2.2 Comparison to call-by-push-value. We took inspiration from call-by-push-value (CBPV), and there are similarities, but there are also significant differences. Both systems have a value-computation distinction, with call-by-name computations and call-by-value values. However, our object theory supports variable binding at arbitrary types while CBPV only supports value variables. 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 suit us. In future work we could make a closer look at the connections to CBPV.

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 $A : \text{Ty}$, we have $\uparrow A : \text{MetaTy}$, pronounced as “lift A ”. This is the type of metaprograms that produce A -typed object programs.
- For $A : \text{Ty}$ 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:

$$\text{let } n : \text{Nat} := \sim(\text{id } \langle 10 + 10 \rangle); \dots$$

Here, id is used at type $\uparrow \text{Nat}$. 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:

$$\begin{aligned} &\text{and} : \uparrow \text{Bool} \rightarrow \uparrow \text{Bool} \rightarrow \uparrow \text{Bool} \\ &\text{and } x \ y = \langle \text{case } \sim x \text{ of True} \rightarrow \sim y; \text{False} \rightarrow \text{False} \rangle \end{aligned}$$

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 : ValTy} → (↑A → ↑B) → ↑(List A) → ↑(List B)
map f as = ⟨letrec go as := case as of
           Nil      → Nil;
           Cons a as → Cons ~(f ⟨a⟩) (go as);
           go ~as⟩

```

For example, this can be used as $\text{let } f \text{ as} : \text{List Nat} \rightarrow \text{List Nat} := \sim(\text{map } (\lambda x. \langle \sim x + 10 \rangle) \langle \text{as} \rangle)$. This is staged to a recursive definition where the mapping function is inlined into the Cons case as $\text{Cons } a \text{ as} \rightarrow \text{Cons } (a + 10) \text{ (go as)}$. Note that map has to abstract over value types, since lists can only contain values, not functions. Also, the mapping function has type $\uparrow A \rightarrow \uparrow B$, instead of $\uparrow(A \rightarrow B)$. 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 $f : \uparrow(A \rightarrow B)$, and f is staged to $\langle \lambda x. t \rangle$, then $\sim f u$ is staged to an undesirable “administrative” β -redex $(\lambda x. t) u$.

2.4 Semantics of Staging

We give a short summary of the semantics of staging here. We reuse most of the semantics from [?], also borrowing notation and definitions, and only mention additions and differences. The reader may refer to *ibid.* for technical details.

Staging is given by evaluation of CFTT into a presheaf model, where the base category is given as follows: objects are typing contexts of the object language and morphisms are parallel substitutions. The equality of morphisms is given by strict syntactic equality of object terms. Ty is interpreted as the presheaf which is constantly the set of object-theoretic types. Soundness of staging is proven by a separate logical relations model, internally to the mentioned presheaf model. Soundness means that the output of staging, viewed as a CFTT term, is definitionally equal to the input.

The main addition in CFTT is the assumption of all inductive families in the meta-level. However, based on the results in [?], it is enough to extend the meta-level syntax in [?] with W-types and identity types to derive all inductive families, so we consider these in the following.

2.4.1 W-types. W-types have a standard object-wise inductive definition in the presheaf model [?]. Additionally, we need to give a logical relation interpretation to W-types, as required for the soundness proof for staging. In Section 5.2. of [?] a serialization map is used to define the relation for natural numbers, however, there is no analogous map for W-types because we cannot build finite terms from infinitely branching trees.

Instead, we give an inductive definition for the logical relation at W-types. First, assuming $A : \widehat{\text{Set}}$ and $B : A \rightarrow \widehat{\text{Set}}$, we have the usual constructor for the semantic W-type in the internal language of the presheaf model, as $\text{sup} : (a : A) \rightarrow (B a \rightarrow W A B) \rightarrow W A B$. Second, assuming $A : \text{Ty } \Gamma$ and $B : \text{Ty } (\Gamma \triangleright A)$ in CFTT, we internalize the sup constructor which constructs a term of a syntactic W-type from restricted terms:

$$\mathbb{R}_{\text{sup}} : (\alpha : \mathbb{R} A \gamma) \rightarrow \mathbb{R} (B \Rightarrow (W_{\text{CFTT}} A B)[\text{wk}]) (\gamma, \alpha) \rightarrow \mathbb{R} (W_{\text{CFTT}} A B) \gamma$$

Here, W_{CFTT} denotes the W-type former in CFTT, $B \Rightarrow (W_{\text{CFTT}} A B)[\text{wk}]$ is a non-dependent function type in CFTT, and wk is a weakening substitution.

Then, assuming $\gamma \approx : \Gamma \approx \gamma \gamma'$, we inductively define a relation $W \approx$ between $W(\mathbb{E} A \gamma) (\lambda \alpha. \mathbb{E} B(\gamma, \alpha))$ and $\mathbb{R} (W_{\text{CFTT}} A B) \gamma'$, which relates $\text{sup } \alpha f$ to $\mathbb{R}_{\text{sup}} \alpha' f'$ if α is related to α' at type A and f is related to f' at type $B[\alpha'] \Rightarrow W_{\text{CFTT}} A B$. The relation at the non-dependent function type \multimap is

given as pointwise preservation of relations. With this, we can also give relational interpretation to the syntactic sup rule in CFTT, and we can use induction on W^\approx to interpret the syntactic induction rule for W . Thus, we get soundness of staging with W -types.

2.4.2 Identity type. Our treatment of the identity type is influenced by the following: we would like to assume identities as axioms in CFTT, which are true in the presheaf model (i.e. they are true at staging-time), without blocking staging as a computation. In Section ?? we use such an axiom to good effect. However, axioms do block computation in CFTT up to definitional equality, which implies that the soundness of staging fails. Our solution:

- Have a standard intensional identity type in CFTT, perhaps with uniqueness of identity proofs assumed.
- Before staging, translate CFTT to CFTT plus equality reflection. Operationally, this means erasing all identity proofs and transports.
- Perform staging after erasure of identity proofs.

In the presheaf model, we interpret the extensional identity in a standard way, as strict equality of sections of presheaves. With this, soundness of staging is trivial to show; we can pick the trivial relation for the interpretation of identity. Hence, we weaken the statement of soundness: every CFTT program is convertible to the staging output, up to erasure of identity proofs.

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 Haskell programming, and they also introduce a great amount of abstraction that should be optimized away.

3.1 Binding-Time Improvements

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

$$\begin{array}{ll} \text{up} : \uparrow\!(A \rightarrow B) \rightarrow \uparrow\!A \rightarrow \uparrow\!B & \text{down} : (\uparrow\!A \rightarrow \uparrow\!B) \rightarrow \uparrow\!(A \rightarrow B) \\ \text{up } f \ a = \langle \sim f \ \sim a \rangle & \text{down } f = \langle \lambda a. \sim(f \ \langle a \rangle) \rangle \end{array}$$

We cannot 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 dependencies are unavoidable. Let us look at binding-time-improvement for product types now:

$$\begin{array}{ll} \text{up} : \uparrow\!(A, B) \rightarrow (\uparrow\!A, \uparrow\!B) & \text{down} : (\uparrow\!A, \uparrow\!B) \rightarrow \uparrow\!(A, B) \\ \text{up } x = (\langle \text{fst } \sim x \rangle, \langle \text{snd } \sim x \rangle) & \text{down } (x, y) = \langle (\sim x, \sim y) \rangle \end{array}$$

Here we overload Haskell-style product type notation, both for types and the pair constructor, 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, $\text{down } (\text{up } \langle f \ x \rangle)$ is staged to $\langle (\text{fst } (f \ x), \text{snd } (f \ x)) \rangle$. 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 in MetaTy, and we cannot introduce object binders in meta-level code. Fortunately, there is a principled solution.

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.

```
newtype Gen A = Gen {unGen : {R : Ty} → (A → ↑R) → ↑R}
```

This is a monad in MetaTy in a standard sense:

```
instance Monad Gen where
  return a = Gen $ λ k. k a
  ga >> f = Gen $ λ k. unGen ga (λ a. (f a k))
```

From now on, we reuse Haskell-style type classes and **do**-notation in CFTT. We will use type classes in an informal way, without precisely specifying how they work. However, type classes are used in essentially the same way in the Haskell implementation of the paper, and also in the Agda implementation, with modest technical differences.

From the Monad instance, the Functor and Applicative instances can be also derived. We reuse ($\langle \$ \rangle$) and ($\langle * \rangle$) for applicative notation as well. An inhabitant of Gen A can be viewed as an action whose effect is to produce some object code. Gen can be only “run” when the result type is an object type:

```
runGen : Gen (↑A) → ↑A
runGen ma = unGen ma id
```

We can let-bind object expressions in Gen:

```
gen : ↑A → Gen (↑A)
gen a = Gen $ λ k. ⟨let x : A := ~a; ~ (k ⟨x⟩)⟩
```

And also recursive definitions of computations:

```
genRec : {A : CompTy} → (↑A → ↑A) → Gen (↑A)
genRec f = Gen $ λ k. ⟨letrec x : A := ~ (f ⟨x⟩); ~ (k ⟨x⟩)⟩
```

Now, using do-notation, we may write **do** {x ← gen (10 + 10); y ← gen (20 + 20)}; **return** ⟨x + y⟩, for a Gen (↑Nat) action. Running this with runGen yields ⟨let x := 10 + 10; let y := 20 + 20; x + y⟩. We can also define a “safer” binding-time improvement for products, using let-insertion:

```
up : ↑(A, B) → Gen (↑A, ↑B)          down : Gen (↑A, ↑B) → ↑(A, B)
up x = do x ← gen x                    down x = do (a, b) ← x
return ((fst ~x), (snd ~x))            return ((~a, ~b))
```

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? The implicitness of Gen may make it harder to reason about the size and content of generated code. 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 → 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:

$$(\gg) : \uparrow(\text{Maybe } A) \rightarrow (\uparrow A \rightarrow \uparrow(\text{Maybe } B)) \rightarrow \uparrow(\text{Maybe } B)$$

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. We define a class for this conversion:

```
class MonadGen M ⇒ Improve (F : ValTy → Ty) (M : MetaTy → MetaTy) where
  up   : {A : ValTy} → ↑(F A) → M (↑A)
  down : {A : ValTy} → M (↑A) → ↑(F A)
```

Assume that `MaybeM` is the standard meta-level monad, and `MaybeTM` is the standard monad transformer:

```
newtype MaybeTM M A = MaybeTM {runMaybeTM : M (MaybeM A)}
```

Now, the binding-time improvement of `Maybe` is as follows:

```
instance Improve Maybe (MaybeTM Gen) where
  up x = MaybeTM $ Gen $ λ k.
    ⟨case ~x of Nothing → ~(k NothingM); Just a → ~(k (JustM a))⟩
  down (MaybeTM (Gen ma)) =
    ma (λ x. case x of NothingM → ⟨Nothing⟩; JustM a → ⟨Just ~a⟩)
```

With this, we get the `Monad` instance for free from `MaybeTM` and `Gen`. A small example:

```
let n : Maybe Nat := ~(down $ do {x ← return ⟨10⟩; y ← return ⟨20⟩; return ⟨x + y⟩}); ...
```

Since `MaybeTM` is meta-level, its monadic binding fully computes at staging time. Thus, the above code is staged to

```
let n : Maybe Nat := Just (10 + 20); ...
```

Assume also a `lift : Monad M ⇒ M A → MaybeTM M A` operation which comes from `MaybeTM` being a monad transformer. We can do `let`-insertion in `MaybeTM Gen` by simply lifting:

```
gen' : ↑A → MaybeTM Gen (↑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:

```
class Monad M ⇒ MonadGen M where
  liftGen : Gen A → M A
instance MonadGen Gen where liftGen = id
instance MonadGen M ⇒ MonadGen (MaybeTM M) where liftGen = lift ∘ liftGen
```

The `MonadGen` instance can be defined uniformly for every monad transformer as `liftGen = lift ∘ liftGen`. We also redefine `gen` and `genRec` to work in any `MonadGen`, so from now on we have:

```
gen      : MonadGen M ⇒ ↑A → M (↑A)
genRec   : MonadGen M ⇒ (↑A → ↑A) → M (↑A)
```

3.3.1 Case splitting in monads. We often want to case-split on object-level data inside a monadic action, like in the following:

```
f : Nat → Maybe Nat
f x := ~(down $ case x == 10 of
  True  → return (x + 5);
  False → NothingM)
```

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:

```
data SplitList A = Nil' | Cons' (↑A) (↑(List A))
```

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

```
split : ↑(List A) → Gen (SplitList A)
split as = Gen $ λ k. ⟨case ~as of Nil → ~(k Nil'); Cons a as → ~(k (Cons' ⟨a⟩ ⟨as⟩))⟩
```

Now, in any `MonadGen`, we may write

```
do {sp ← liftGen (split as); (case sp of Nil' → ...; Cons' a as → ...)}
```

This can be generalized to splitting on any object value. In the Agda and Haskell implementations, we overload `split` with a class similar to the following:

```
class Split (A : ValTy) where
  SplitTo : MetaTy
  split    : ↑A → Gen SplitTo
```

In a native implementation of CFTT it may make sense to extend `do`-notation, so that we elaborate `case` on object values to an application of the appropriate `split` function. We adopt this in the rest of the paper, so we will be able to write `case x of ...` in a `do`-block, whenever we work in a `MonadGen` and `x` has a value type.

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:

```
newtype Identity A := Identity {runIdentity : A}
instance Improve Identity Gen where
  up x = Gen $ λ k. k (runIdentity ~x)
  down x = unGen x $ λ a. ⟨Identity ~a⟩
```

We also recover the object-level Maybe as MaybeT Identity, from the following MaybeT:

```
newtype MaybeT (M : ValTy → Ty) (A : ValTy) := MaybeT {runMaybeT : M (Maybe A)}
```

With this, improvement can be generally defined for MaybeT:

```
instance Improve F M ⇒ Improve (MaybeT F) (MaybeTM M) where
  up x = MaybeTM $ do
    ma ← up ⟨runMaybeT ~x⟩
    case ma of Nothing → return NothingM
           Just a    → return (JustM a)
  down (MaybeTM x) = ⟨MaybeT ~ (down $ x >>= λ case
    NothingM → return ⟨Nothing⟩
    JustM a  → return ⟨Just ~a⟩)⟩
```

In the **case** in up, we already use our syntax sugar for matching on a Maybe value inside an M action. This is legal, since we know from the Improve F M assumption that M is a MonadGen. In down we also use λ **case** ... to shorten λ x. **case** x **of**

In the meta level, we can reuse essentially all definitions from Haskell’s monad transformer library mtl. 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 StateT_M as the standard meta-level definition. The object-level StateT has type (S : ValTy)(F : ValTy → Ty)(A : ValTy) → ValTy; the state parameter S has to be a value type, since it is an input to an object-level function.

```
instance Improve F M ⇒ Improve (StateT S F) (StateTM (↑S) M) where
  up x = StateTM $ λ s. do
    as ← up ⟨runStateT ~x ~s⟩
    case as of (a, s) → return (a, s)
  down x = ⟨StateT (λ s. ~ (down $ do
    (a, s) ← runStateTM x ⟨s⟩
    return ⟨(~a, ~s)⟩)⟩
```

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 Improve F M ⇒ Improve (ReaderT R F) (ReaderTM (↑R) M) where
  up x = ReaderTM $ λ r. up ⟨runReaderT ~x ~r⟩
  down x = ⟨ReaderT (λ r. ~ (down (runReaderTM x ⟨r⟩)))⟩
```

3.4.1 State and Reader operations. If we try to use the modify function that we already have for State_M, a curious thing happens. The meaning of modify (λ x. ⟨~x + ~x⟩) is to replace the current state x, as an object expression, with the expression ⟨~x + ~x⟩, 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 $ do {modify (λ x. ⟨~x + ~x⟩); modify (λ x. ⟨~x + ~x⟩); return ⟨()⟩}
```

is staged to

```
⟨λ x. ⟨(), (x + x) + (x + x)⟩⟩
```

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 $\uparrow()$ from actions instead of $()$ – the former is more convenient, because the down operation can be immediately used on it.

$$\text{put}' : (\text{MonadState } (\uparrow S) M, \text{MonadGen } M) \Rightarrow \uparrow S \rightarrow M (\uparrow())$$

$$\text{put}' s = \text{do } \{s \leftarrow \text{gen } s; \text{put } s; \text{return } \langle () \rangle\}$$

$$\text{modify}' : (\text{MonadState } (\uparrow S) M, \text{MonadGen } M) \Rightarrow (\uparrow S \rightarrow \uparrow S) \rightarrow M (\uparrow())$$

$$\text{modify}' f = \text{do } \{s \leftarrow \text{get}; \text{put}' (f s)\}$$

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

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

Now,

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

is staged to

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

3.5 Joining Control Flow in Monads

There is a deficiency in our library so far. Assuming $b : \text{Bool}$ as an object value, consider:

$$\text{down } \$ \text{do}$$

$$\text{case } b \text{ of}$$

$$\text{True} \rightarrow \text{put}' \langle 10 \rangle$$

$$\text{False} \rightarrow \text{put}' \langle 20 \rangle$$

$$\text{modify}' (\lambda x. \langle \sim x + 10 \rangle)$$

The $\text{modify}' (\lambda x. \langle \sim x + 10 \rangle)$ action gets inlined to 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 both 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:

$$\text{down } \$ \text{do}$$

$$x \leftarrow \text{gen } \$ \text{down } \$ \text{case } b \text{ of}$$

$$\text{True} \rightarrow \text{put}' \langle 10 \rangle$$

$$\text{False} \rightarrow \text{put}' \langle 20 \rangle$$

$$\text{up } x$$

$$\text{modify}' (\lambda x. \langle \sim x + 10 \rangle)$$

This removes code duplication by round-tripping through an object-level **let**. This solution is fairly good in a state monad, since down only introduces a runtime pair constructor, and it is feasible to compile object-level pairs as unboxed data, without overheads. However, 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, for example, we would like to produce object code which looks like the following:

```

let joinJust n := ...
let joinNothing () := ...
case x == 10 of
  True  → joinJust (x + 10)
  False → 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 meta-level finite sum of value types. Recall that Gen A is defined as $\{R : \text{ValTy}\} \rightarrow (A \rightarrow \uparrow R) \rightarrow \uparrow R$. Here, if A is a finite sum, we can rearrange $A \rightarrow \uparrow R$ to a finite product of functions.

We could proceed with finite sums, but we will need finite *sums-of-products* (SOP in short) later in Section in ??, so we develop SOP-s. We view SOP-s as a Tarski-style universe consisting of a type of descriptions, and a way of interpreting descriptions into MetaTy (“El” for “elements” of the type).

```

USOP : MetaTy          ElSOP : USOP → MetaTy
USOP = List (List ValTy)  ElSOP A = ...

```

A type description is a list of lists of value types. We decode this to a sum of products of value types. We omit the definition here. SOP is closed under value types, finite product types and finite sum types. For instance, we have $\text{Either}_{\text{SOP}} : \text{SOP} \rightarrow \text{SOP} \rightarrow \text{SOP}$ together with $\text{Left}_{\text{SOP}} : \text{El}_{\text{SOP}} A \rightarrow \text{El}_{\text{SOP}} (\text{Either}_{\text{SOP}} A B)$, $\text{Right}_{\text{SOP}} : \text{El}_{\text{SOP}} B \rightarrow \text{El}_{\text{SOP}} (\text{Either}_{\text{SOP}} A B)$ and a case splitting operation. It is more convenient to work with type formers in MetaTy, and only convert to SOP representations when needed, so we define a class for the representable types:

```

class IsSOP (A : MetaTy) where
  Rep : USOP
  rep  : A ≈ ElSOP Rep

```

Above, $A \approx \text{Rep}$ denotes a record containing a pair of back-and-forth functions, together with proofs (as propositional equalities) that they are inverses. We will overload rep as the forward conversion function with type $A \rightarrow \text{Rep}$, and write rep^{-1} for its inverse. Now we define an isomorphic representation of $\text{El}_{\text{SOP}} A \rightarrow \uparrow R$ as a product of object-level functions:

```

FunSOP↑ : USOP → Ty → MetaTy
FunSOP↑ Nil      R = ()
FunSOP↑ (Cons A B) R = (↑(foldr (→) R A), FunSOP↑ B R)

tabulate : (ElSOP A → ↑R) → FunSOP↑ A R
index    : FunSOP↑ A R → (ElSOP A → ↑R)

```

We omit here the definitions of tabulate and index. We will also need to let-bind all functions in a FunSOP^\uparrow :

```

genFunSOP↑ : {A : USOP} → FunSOP↑ A R → FunSOP↑ A R
genFunSOP↑ {Nil}      () = return ()
genFunSOP↑ {Cons _ A} (f, fs) = (.,) <$> gen f <*> genFunSOP↑ {A} fs

```


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

```
class Monad M  $\Rightarrow$  MonadJoin M where join : IsSOP A  $\Rightarrow$  M A  $\rightarrow$  M A
```

The most interesting instance is the “base case” for Gen:

```
instance MonadJoin Gen where
  join ma = Gen $  $\lambda$  k. runGen $ do
    joinPoints  $\leftarrow$  genFunSOP $\uparrow$  (tabulate (k  $\circ$  rep-1))
    a  $\leftarrow$  ma
  return $ index joinPoints (rep a)
```

Here we first convert $k : A \rightarrow \uparrow R$ to a product of join points and let-bind each one of them. Then we generate code that returns to the appropriate join point for each return value. Other MonadJoin instances are straightforward. In StateT_M, we need the extra IsSOP S constraint because S is returned as a result value.

```
instance (MonadJoin M)  $\Rightarrow$  MonadJoin (MaybeTM M) where
  join (MaybeTM ma) = MaybeTM (join ma)
instance (MonadJoin M)  $\Rightarrow$  MonadJoin (ReaderTM R M) where
  join (ReaderTM ma) = ReaderTM (join  $\circ$  ma)
instance (MonadJoin M, IsSOP S)  $\Rightarrow$  MonadJoin (StateTM S M) where
  join (StateTM ma) = StateTM (join  $\circ$  ma)
```

Now, whenever the return value of a case-splitting action is a sum of products of values (this includes just returning an object value, which is by far the most common situation), we can use `join $ case x of ...` to eliminate code duplication, without creating runtime sum types.

bigger example, non-tail-recursive list function with some effects, input, output, mentioned tail rec

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 pattern match on object-level values in monadic code, insert object-level `let`-s with `gen` and avoid code duplication with `join`.
- 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`. Monadic operations are 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 some 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 require that stages of `let`-definitions are always specified (as we do here). Moreover, `up` and `down` should be also effectively inferable, using bidirectional elaboration. With stage inference and binding-time improvement inference, monadic code in CFTT would look only modestly more complicated than in Haskell.

Second, in CFTT we cannot store computations (e.g. functions or State actions) in runtime data structures, nor can we have computations in State state or in Reader environments. However, it would be possible to extend CFTT with a closure type former that converts computations to values, in which case there is no such limitation anymore. Here, closure-freedom would be still available; we would be able to pick where to use or avoid the closure type former.

3.7 Agda & Haskell Implementations

We implemented everything in this section in both Agda and typed Template Haskell. We summarize features and differences:

- The Haskell implementation can be used to generate code that can be further compiled by GHC; here the object language is taken to be Haskell itself. Since Haskell does not distinguish value and computation types, we do not track them in the library, and we do not get guaranteed closure-freedom from GHC. We do get high quality code though, typically.
- In Agda, we postulate all types and terms of the object theory in a faithful way (i.e. equivalently to the CFTT syntax presented here), and take Agda itself to be the metalanguage. Here, we can test “staging” by running Agda programs which compute object expressions. However, we can only inspect staging output and cannot compile or run object programs.
- For sums-of-products in Haskell, we make heavy use of *singleton types* [?] to emulate dependent types. This adds significant noise. Also, in *IsSOP* instances we can only define the conversion functions and cannot prove that they are inverses, because Haskell does not have enough support for dependent types.

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. Stream fusion 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 [?]. 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 streams.

The reason is that pull streams have been historically more difficult to efficiently compile, and we can demonstrate significant improvements in CFTT. We also use dependent types in a more essential way than in the previous section.

4.1 Streams

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

```
data Step S A = Stop | Skip S | Yield A S
data Pull (A : MetaTy) : MetaTy where
  Pull : (S : MetaTy) → IsSOP S ⇒ Gen S → (S → Gen (Step S A)) → Pull A
```

In the Pull constructor, *S* is the type of the internal state which is required to be a sum-of-products of value types by the *IsSumVS S* constraint. Borrowing terminology from Coutts [?], we call each product-of-value-types in the state a *state shape*.

The next field with type *Gen S* is the initial state, while transitions are represented by the *S → Gen (Step S A)* field. 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).

Let us see some operations on streams now. Pull is evidently a Functor. It is not quite a “zippy” applicative functor, because its application operator requires an extra `IsSOP` constraint:

```
repeat : A → Pull A
repeat a = Pull () (return ()) (λ_. return $ Yield a ())
```

```
(<*>Pull) : IsSOP A ⇒ Pull (A → B) → Pull A → Pull B
(<*>Pull) (Pull S seed step) (Pull S' seed' step') =
  Pull (S, S', Maybe A) ((,) <$> seed<*>((, Nothing) <$> seed')) $ λ case
    (s, s', Just a) → step s >>= λ case
      Stop      → return Stop
      Skip s     → return $ Skip (s, s', Just a)
      Yield f s  → return $ Yield (f a) (s, s', Nothing)
    (s, s', Nothing) → step' s' >>= λ case
      Stop      → return Stop
      Skip s'    → return $ Skip (s, s', Nothing)
      Yield a s' → return $ Skip (s, s', Just a)
```

In `repeat`, the state is the unit type, while in `(<*>Pull)`, we take the product of states, and also add another `Maybe A` that is used to buffer a single `A` value while we are stepping the `Pull (A → B)` stream. The `IsSOP` constraints for the new machine states are implicitly dispatched by instance resolution; this works in the Agda and Haskell versions too. We can derive `zip` and `zipWith` from `(<*>Pull)` and `(<$>)` for streams, with the restriction that the zipped streams must all produce `IsSOP` values.

`Pull` is also a monoid with stream appending as the binary operation.

```
empty : Pull A
empty = Pull () (return ()) (λ_. return Stop)
```

```
(<>) : Pull A → Pull A → Pull A
(<>) (Pull S seed step) (Pull S' seed' step') = Pull (Either S S') (Left <$> seed) λ case
  Left s  → step s >>= λ case
    Stop      → (Skip ∘ Right) <$> seed'
    Skip s     → return $ Skip (Left s)
    Yield a s  → return $ Yield a (Left s)
  Right s' → step s' >>= λ case
    Stop      → return Stop
    Skip s'    → return $ Skip (Right s')
    Yield a s' → return $ Yield a (Right s')
```

These definitions are standard for streams; compared to the unstaged definitions in previous literature, the only additional noise is just the `Gen` monad in the initial states and the transitions. Likewise, we can give standard definitions for usual stream functions such as `filter`, `take` or `drop`.

4.2 Running Streams with Mutual Recursion

How do we generate object code from streams? The `S` state is given as a finite sums of products, but the sums and the products are on the meta level, so we cannot directly use `S` in object code. Similarly

as in the the treatment of join points, we tabulate the $S \rightarrow \text{Gen Step } S \ A$ transition function to a product of functions. However, these functions need to be *mutually recursive*, since it is possible to transition from any state to any other state, and each such transition is represented as a function call. This problem of generating well-typed mutual blocks was addressed by Yallop and Kiselyov in [?]. Unlike *ibid.*, which used control effects and mutable references in MetaOCaml, we present a solution that does not use side effects in the metalanguage.

The solution is to extend CompTy with finite products of computations, i.e. assume $() : \text{CompTy}$ and $(,) : \text{CompTy} \rightarrow \text{CompTy} \rightarrow \text{CompTy}$, together with pairing and projections. These types, like functions, are call-by-name in the runtime semantics, and they also cannot escape the scope of their definition. Hence, we can also “saturate” programs that involve computational product types: every computation definition at type (A, B) can be translated to a pair, and every projection of a let-defined variable can be statically matched up with a pairing in the variable definition. Thus, a recursive let-definition at type $(A \rightarrow B, A \rightarrow B)$ can be always compiled to a pair of mutual functions.

We redefine the previous $\text{Fun}_{\text{SOP}\uparrow}$ to return a computation type instead:

$$\begin{aligned} \text{Fun}_{\text{SOP}\uparrow} &: \text{U}_{\text{SOP}} \rightarrow \text{Ty} \rightarrow \text{CompTy} \\ \text{Fun}_{\text{SOP}\uparrow} \text{ Nil} & \quad R = () \\ \text{Fun}_{\text{SOP}\uparrow} (\text{Cons } A \ B) \ R &= (\text{foldr } (\rightarrow) \ R \ A, \text{Fun}_{\text{SOP}\uparrow} \ B \ R) \\ \text{tabulate} &: (\text{El}_{\text{SOP}} \ A \rightarrow \uparrow R) \rightarrow \uparrow (\text{Fun}_{\text{SOP}\uparrow} \ A \ R) \\ \text{index} &: \uparrow (\text{Fun}_{\text{SOP}\uparrow} \ A \ R) \rightarrow (\text{El}_{\text{SOP}} \ A \rightarrow \uparrow R) \end{aligned}$$

Even for join points, this is just as efficient as the previous $\text{Fun}_{\text{SOP}\uparrow}$ version, since a definition of a product of functions will get compiled to a sequence of function definitions. In addition, we can use it to generate object code from streams. There are several choices for this, but in CFTT the foldr function is as good as we can get.

$$\begin{aligned} \text{foldr} &: \{A : \text{MetaTy}\} \{B : \text{Ty}\} \rightarrow (A \rightarrow \uparrow B \rightarrow \uparrow B) \rightarrow \uparrow B \rightarrow \text{Pull } A \rightarrow \uparrow B \\ \text{foldr } \{A\} \{B\} (\text{Pull } S \ \text{seed} \ \text{step}) \ f \ b &= \langle \\ &\quad \text{letrec } fs : \text{Fun}_{\text{SOP}\uparrow} (\text{Rep } \{S\}) \ B := \sim(\text{tabulate } \$ \lambda s. \text{unGen } (\text{step } (\text{rep}^{-1} \ s)) \$ \lambda \text{case} \\ &\quad \text{Stop} \quad \rightarrow b \\ &\quad \text{Skip } s \quad \rightarrow \text{index } \langle fs \rangle (\text{rep } s) \\ &\quad \text{Yield } a \ s \rightarrow f \ a \ (\text{index } \langle fs \rangle (\text{rep } s))); \\ &\quad \sim(\text{unGen } \text{seed } \$ \lambda s. \text{index } \langle fs \rangle (\text{rep } s)) \rangle \end{aligned}$$

This foldr is quite flexible because we can compute any object value from it, including functions. For instance, we can define foldl from foldr :

$$\begin{aligned} \text{foldl} &: \{A : \text{MetaTy}\} \{B : \text{ValTy}\} \rightarrow (\uparrow B \rightarrow A \rightarrow \uparrow B) \rightarrow \uparrow B \rightarrow \text{Pull } A \rightarrow \uparrow B \\ \text{foldl } f \ b \ \text{as} &= \sim(\text{foldr } (\lambda a \ g. \langle \lambda b. \sim g \sim (f \ \langle b \rangle \ \sim a) \rangle) \ \langle \lambda b. b \rangle \ \text{as}) \sim b \end{aligned}$$

Note that since we abstract B in a runtime function, it must be a value type. Here, each Stop and Yield in the transition function gets interpreted as a λ -expression in the output. However, those λ -s will be lifted out in the scope, yielding a proper mutually tail-recursive definition with an accumulator for B . Contrast this to the GHC base library, where foldl for lists is also defined from foldr , to enable push-based fusion, but where a substantial *arity analysis* is required in the compiler to eliminate the intermediate closures [?].

4.3 concatMap for Streams

Can we have a list-like Monad instance for streams, with singleton streams for **return** and **concatMap** for binding? This is not possible. Aiming at

$$\text{concatMap} : (A \rightarrow \text{Pull } B) \rightarrow \text{Pull } A \rightarrow \text{Pull } B,$$

the $A \rightarrow \text{Pull } B$ function can contain an infinite number of different machine state types, which cannot be represented in a finite amount of object code. Here by “infinite” we mean the notion that is internally available in the meta type theory. For instance, we can define a $\text{Nat}_M \rightarrow \text{Pull } B$ function which for each $n : \text{Nat}_M$ produces a concatenation of n streams. Hence, we shall have the following function instead:

$$\text{concatMap} : \text{IsSOP } A \Rightarrow (A \rightarrow \text{Pull } B) \rightarrow \text{Pull } A \rightarrow \text{Pull } B,$$

The idea is the following: if U_{SOP} is closed under dependent sum types, we can directly define this **concatMap**, by taking the appropriate dependent sum of the $A \rightarrow \text{U}_{\text{SOP}}$ family of machine states, which we extract from the $A \rightarrow \text{Pull } B$ function. Let us write $\Sigma A B : \text{MetaTy}$ for dependent sums, for $A : \text{MetaTy}$ and $B : A \rightarrow \text{MetaTy}$, and reuse $(.)$ for pairing. We also use the following field projection functions: $\text{projS} : \text{Pull } A \rightarrow \text{MetaTy}$, $\text{projSeed} : (\text{as} : \text{Pull } A) \rightarrow \text{projS as}$, and $\text{projStep} : (\text{as} : \text{Pull } A) \rightarrow \text{projS as} \rightarrow \text{Gen } (\text{Step } (\text{projS as}) A)$. For now, we assume that the following instance exists:

$$\text{instance } (\text{IsSOP } A, \{a : A\} \rightarrow \text{IsSOP } (B a)) \Rightarrow \text{IsSOP } (\Sigma A B)$$

Above, $\{a : A\} \rightarrow \text{IsSOP } (B a)$ is a universally quantified instance constraint; this is available in Agda and, less generally, in Haskell [?]. The definition of **concatMap** is as follows.

```
concatMap : IsSOP A => (A -> Pull B) -> Pull A -> Pull B
concatMap {A} {B} f (Pull S seed step) =
  Pull (S, Maybe (Σ A (projS ∘ f))) ((.) <$> seed <*> return Nothing) $ λ case
    (s, Nothing) -> step s >>= λ case
      Stop      -> return Stop
      Skip s    -> return $ Skip (s, Nothing)
      Yield as  -> do {s' ← projSeed (f a); return $ Skip (s, Just (a, s'))}
    (s, Just (a, s')) -> projStep (f a) s' >>= λ case
      Stop      -> return $ Skip (s, Nothing)
      Skip s'   -> return $ Skip (s, Just (a, s'))
      Yield b s' -> return $ Yield b (s, Just (a, s'))
```

Here, **Nothing** marks the states where we are in the “outer” loop, running the **Pull A** stream until we get its next value. **Just** marks the states of the “inner” loop, where we have a concrete $a : A$ value and we run the $(f a)$ stream until it stops. In the inner loop, the machine state type depends on the $a : A$ value, hence the need for Σ .

How do we get **IsSOP** for Σ ? The key observations are:

- Metaprograms cannot inspect the structure of object terms.
- Object types do not depend on object terms.

Hence, we expect that during staging, every $f : \uparrow A \rightarrow \text{ValTy}$ function has to be constant. This is actually true in the staging semantics of CFTT. In the semantics, all metafunctions are stable

under object-level parallel substitutions. Also, object types are untouched by substitution. Hence, a straightforward unfolding of the definitions in the staging semantics validates the constancy of f .

Generally speaking, every function whose domain is a product of object types and whose codomain is a constant presheaf in the semantics, is a constant function in the the semantics. We may call these constancy statements *generativity axioms*, since they reflect the inability of metaprograms to inspect terms, and “generativity” often refers to this property in the staging literature [?].

Let us write $U_P = \text{List ValTy}$ and $El_P : U_P \rightarrow \text{MetaTy}$ for a universe of finite products of value types. Concretely in CFTT and the Agda implementation, we assume the following:

Axiom (generativity). *Every $f : El_P A \rightarrow U_{SOP}$ is constant.*

We remark that there is no risk of staging getting *stuck* on this axiom, because propositional equality proofs get *erased* in the staging semantics, as we noted in Section ??.

We also remark that although generativity is inconsistent with inspecting the structure of object terms, it is consistent with inspecting the structure of object types.

From generativity, we derive $\Sigma_{SOP} : (A : U_{SOP}) \rightarrow (B : El_{SOP} A \rightarrow U_{SOP}) \rightarrow U_{SOP}$ as follows. First, for each $A : U_P$, we define $loop_A : El_P A$ as a product of non-terminating object programs. This is only needed to get arbitrary inhabitants with which we can call $El_P A \rightarrow U_{SOP}$ functions. Conveniently, every object type is inhabited by infinite loops.

Then, $\Sigma_{SOP} A B$ is defined as the concatenation of $A_i \times B (\text{inject}_i \text{loop}_{A_i})$ for each $A_i \in A$, where (\times) is the product type former in U_{SOP} and inject_i has type $El_P A_i \rightarrow El_{SOP} A$. This is similar to the definition of non-dependent products in U_{SOP} , except that we have to get rid of the type dependency by instantiating B with looping programs.

Then, we can show using the generativity axiom that Σ_{SOP} supports projections, pairing and the $\beta\eta$ -rules. These are all needed when we define the $lsSOP$ instance, when we have to prove that encoding via rep is an isomorphism. Concretely, assuming $lsSOP A$ and $\{a : A\} \rightarrow lsSOP (B a)$, we define Rep for $\Sigma A B$ as follows:

$$Rep \{\Sigma A B\} = \Sigma (Rep \{A\}) (\lambda x. Rep \{B (rep^{-1} x)\})$$

Then, the definition of encoding for Σ is only well-typed up to the fact that $rep^{-1} \circ rep = id$:

$$rep\{\Sigma A B\} (a, b) = (rep \{A\} a, rep \{B (rep^{-1} (rep x))\} b)$$

This is, in fact, our reason for including the isomorphism equations as well in $lsSOP$. This, in turn, necessitates using SOP instead of finite sums. We can define a product type former for finite sums of value types, by taking the pairwise products of components. However, we can only take the *object-level* products here, and since object-level products have no $\beta\eta$ -rules, we cannot prove $\beta\eta$ for the derived product type in finite sums, and likewise for the derived Σ -type. When we use SOP instead, we do not have to use object-level products and this issue does not appear.

We do not detail the full definition of Σ_{SOP} here. The reader may refer to the SOP module in the Agda implementation, which is altogether 260 lines with all $lsSOP$ instances.

4.4 Let-Insertion & Case Splitting in Monadic Style

While $Pull$ is not a monad, and hence also not a $MonadGen$, we can still use a monadic style of stream programming with good ergonomics. First, we need singleton streams for “returning”:

`single : A → Pull A`

`single a = Pull BoolM True $ λ b. return $ if b then Yield a False else Stop`

Now, this operation should be avoided when possible, since it has two states and can contribute to a blow-up of state size. For example, `concatMap` over `single` introduces at least a doubling of the number of state shapes. Although let-insertion and case splitting could be derived as `concatMap` over `single`, to avoid the size explosion we give more specialized definitions instead. First, we define a helper function for binding for a single `lsSOP` value.

```

bindsingle : lsSOP A' ⇒ ↑A → (↑A → Gen A') → (A' → Pull B) → Pull B
bindsingle a f g =
  Pull (Σ A' (projS ∘ g)) (do {a' ← f a; s ← projSeed (g a'); return (a', s)}) $ λ case
    (a', s) → projStep (g a') s ≍ λ case
      Stop    → return Stop
      Skip s   → return $ Skip (a', s)
      Yield b s → return $ Yield b (a', s)

```

We recover let-insertion and case splitting as follows:

```

genPull : ↑A → (↑A → Pull B) → Pull B
genPull a = bindsingle a gen

casePull : (Split A, lsSOP (SplitTo {A})) ⇒ ↑A → (SplitTo {A} → Pull B) → Pull B
casePull a = bindsingle a split

```

Let us look at small example. We write `forEach` for `concatMap` with flipped arguments.

```

forEach (take ⟨100⟩ (countFrom ⟨0⟩)) $ λ x.
  genPull (∼x * 2) $ λ y.
  casePull (∼x < 50) $ λ case
    True  → take y (countFrom x)
    False → single y

```

Here, in every `forEach` iteration, `genPull (∼x * 2)` evaluates the given expression and saves the result to the machine state. Then, `casePull` branches on an object-level Boolean expression. If we use `foldr` to generate object code from this definition, we get *four* mutually defined functions. We get two state shapes from `single y` and one from `take y (countFrom x)`; we sum these to get three for the `casePull`, then we get an extra shape from `forEach` which introduces an additional `Maybe`.

4.5 Discussion

Our stream library has a fairly strong support for programming in a monadic style, even though `Pull` is not literally a monad. We can bind object values with `concatMap`, and we can also do let-insertion and case splitting for them. We also get guaranteed well-typing, closure-freedom, and arbitrary mixing of zipping and `concatMap`.

We highlight the usage of the generativity axiom as well. Previously in staged compilation, intensional analysis (i.e. the ability to analyze object code) has been viewed as a desirable feature that increases the expressive power of the system. To our knowledge, our work is the first one that exploits the *lack* of intensional analysis in metaprogramming. This is a bit similar to parametricity in type theories, where the inability to analyze types has a payoff in the form of “free theorems” [?].

Regarding the practical application of our stream library, we think that it would make sense to support both push and pull fusion in a realistic implementation, and allow users to benefit from the

strong points of both. Push streams, which we do not present in this paper, have proper `Monad` and `MonadGen` instances and are often more convenient to use in CFTT. They are also better for deep traversals of structures where they can utilize unbounded stack allocations, while pull streams need heap allocations for unbounded space in the machine state.

4.6 Agda & Haskell Implementations

In Agda, to avoid computation getting stuck on the generativity axiom, we use the `primTrustMe` built-in [?] to automatically erase the axiom when the sides of the equation are definitionally equal. Otherwise the implementation closely follows this section.

In Haskell there are some limitations. First, in `concatMap`, the projection function `projS` cannot be defined, because Haskell is not dependently typed, and the other field projections are also out of reach. We only have a weaker “positive” recursion principle for existential types. It might be the case that a strongly typed `concatMap` is possible with only weak existentials, but we attempted this and found that it introduces too much technical complication.

So instead of giving a single generic definition for `concatMap` for `IsSOP` types, we define `concatMap` just for object types,¹ and define `casePull` separately for each object type, as an overloaded class method. In each of these definitions, we only need to deal with a concrete finite number of machine state types, which is feasible to handle with weak existentials.

Also, the generativity axiom is *false* in Template Haskell, since it is possible to look inside quoted expressions. Instead, we use type coercions that can fail at staging time. If library users do not violate generativity, these coercions disappear and the staging output will not contain unsafe coercions.

5 RELATED WORK

REFERENCES

¹Recall that we do not distinguish value and computation types in Haskell.