

## 1 Closure-Free Functional Programming in a Two-Level Type 2 Theory 3

4 ANONYMOUS AUTHOR(S)  
5

6 Many abstraction tools in functional programming rely heavily on general-purpose compiler optimization  
7 to achieve adequate performance. For example, monadic binding is a higher-order function which yields  
8 runtime closures in the absence of sufficient compile-time inlining and beta-reductions, thereby significantly  
9 degrading performance. In current systems such as the Glasgow Haskell Compiler, there is no strong guarantee  
10 that general-purpose optimization can eliminate abstraction overheads, and users only have indirect and  
11 fragile control over code generation through inlining directives and compiler options. We propose a two-stage  
12 language to simultaneously get strong guarantees about code generation and strong abstraction features. The  
13 object language is a simply-typed first-order language which can be compiled without runtime closures. The  
14 compile-time language is a dependent type theory. The two are integrated in a two-level type theory.  
15

16 We demonstrate two applications of the system. First, we develop monads and monad transformers. Here,  
17 abstraction overheads are eliminated by staging and we can reuse almost all definitions from the existing  
18 Haskell ecosystem. Second, we develop pull-based stream fusion. Here we make essential use of dependent  
19 types to give a concise definition of a concatMap operation with guaranteed fusion. We provide an Agda  
20 implementation and a typed Template Haskell implementation of these developments.  
21

22 CCS Concepts: • Theory of computation → Type theory; • Software and its engineering → Source  
23 code generation.  
24

25 Additional Key Words and Phrases: two-level type theory, staged compilation  
26

27 ACM Reference Format:  
28

29 Anonymous Author(s). 2024. Closure-Free Functional Programming in a Two-Level Type Theory. 1, 1 (Febru-  
30 ary 2024), 27 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>  
31

### 32 1 INTRODUCTION

33 Modern functional programming supports many convenient abstractions. These often come with  
34 significant runtime overheads. Sometimes the overheads are acceptable, but in other cases compiler  
35 optimization is crucial. Monads in Haskell is an example for the latter. Even the Reader monad,  
36 which is one of the simplest in terms of implementation, yields large overheads when compiled  
37 without optimizations. Consider the following:  
38

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

41 With optimizations enabled, GHC compiles this roughly to the code below:  
42

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

---

45 Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee  
46 provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the  
47 full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored.  
48 Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires  
49 prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).  
50 © 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
51 ACM XXXX-XXXX/2024/2-ART  
52 <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

## 5 Closure-Free Functional Programming in a Two-Level Type 6 Theory 7

8 FIRST LAST, Institution, Country  
9

10 Many abstraction tools in functional programming rely heavily on general-purpose compiler optimization  
11 to achieve adequate performance. For example, monadic binding is a higher-order function which yields  
12 runtime closures in the absence of sufficient compile-time inlining and beta-reductions, thereby significantly  
13 degrading performance. In current systems such as the Glasgow Haskell Compiler, there is no strong guarantee  
14 that general-purpose optimization can eliminate abstraction overheads, and users only have indirect and  
15 fragile control over code generation through inlining directives and compiler options. We propose a two-stage  
16 language to simultaneously get strong guarantees about code generation and strong abstraction features. The  
17 object language is a simply-typed first-order language which can be compiled without runtime closures. The  
18 compile-time language is a dependent type theory. The two are integrated in a two-level type theory.  
19

20 We demonstrate two applications of the system. First, we develop monads and monad transformers. Here,  
21 abstraction overheads are eliminated by staging and we can reuse almost all definitions from the existing  
22 Haskell ecosystem. Second, we develop pull-based stream fusion. Here we make essential use of dependent  
23 types to give a concise definition of a concatMap operation with guaranteed fusion. We provide an Agda  
24 implementation and a typed Template Haskell implementation of these developments.  
25

26 CCS Concepts: • Theory of computation → Type theory; • Software and its engineering → Source  
27 code generation.  
28

29 Additional Key Words and Phrases: two-level type theory, staged compilation  
30

31 ACM Reference Format:  
32

33 First Last. 2024. Closure-Free Functional Programming in a Two-Level Type Theory. 1, 1 (June 2024), 34 pages.  
34 <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>  
35

### 36 1 INTRODUCTION

37 Modern functional programming supports many convenient abstractions. These often come with  
38 significant runtime overheads. Sometimes the overheads are acceptable, but in other cases compiler  
39 optimization is crucial. Monads in Haskell is an example for the latter. Even the Reader monad,  
40 which is one of the simplest in terms of implementation, yields large overheads when compiled  
41 without optimizations. Consider the following:  
42

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

45 With optimizations enabled, GHC compiles this roughly to the code below:  
46

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

---

49 Author's address: First Last, Institution, Country, City.  
50

51 Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee  
52 provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the  
53 full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored.  
54 Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires  
55 prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).  
56 © 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
57 ACM XXXX-XXXX/2024/6-ART  
58 <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

2

Anon.

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 adept 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 [GHC developers 2024a], but without strong guarantees, 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 [GHC developers 2024b]. It is common practice in high-performance Haskell programming to visually 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.

**Two-level type theory** (2LTT) [Annenkov et al. 2019; Kovács 2022] makes it possible to use expressive dependent type theories for metaprogramming for a wide range of object languages. 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.

2

First Last

Without optimizations we roughly get:

```
f = λ x. (⇒⇒) monadReader (ask monadReaderReader) (λ b. if b
    then return monadReader (x + 10)
    else return monadReader (x + 20))
```

Here, `monadReader` and `monadReaderReader` are runtime dictionaries, respectively for the `Monad` and `MonadReader` instances, and, for example, `(⇒⇒) monadReader` is a field projection from the dictionary. This results from the dictionary-passing elaboration of type classes [Wadler and Blott 1989]. We get a runtime closure from the `λ b. ...` function, and `(⇒⇒)`, `ask` and `return` also produce additional 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 rank-2 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 is applied to.

GHC's optimization efforts are respectable, and it has gotten quite adept 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 [GHC developers 2024a], but without strong guarantees, 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 [GHC developers 2024b]. It is common practice in high-performance Haskell programming to visually 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 work as possible 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.

We use **two-level type theory** (2LTT) [Annenkov et al. 2019; Kovács 2022] to sweeten the deal of staged compilation, aiming for a combination of strong guarantees, good ergonomics, high level

99 We develop a particular two-level type theory for this purpose, which we call **CFTT**, short for  
 100 “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 dependent type theory for the compile-time language. This allows us to recover many features by metaprogramming.

101 Since the object language is first-order, we guarantee that all programs in CFTT can be ultimately  
 102 compiled without any dynamic closures, using only calls and jumps to statically known code.  
 103 Why emphasize closures? They are the foundation to almost all abstraction tools in functional  
 104 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 [Wadler and Blott 1989].
- Functors and first-class modules in OCaml [Leroy et al. 2023] and other ML-s rely on closures.

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

107 Perhaps surprisingly, little practical programming relies essentially on closures. Most of the  
 108 time, programmers use higher-order functions for *abstraction*, such as when mapping over lists,  
 109 where it is expected that the mapping function will be inlined. We note though that our setup is  
 110 compatible with closures as well, and it can support two separate type formers for closure-based  
 111 and non-closure-based “static” functions. Having both of these would be desirable in a practical  
 112 system. In the current work we focus on the closure-free case because it is much less known and  
 113 developed, and it is interesting to see how far we can go with it.

## 114 1.2 Contributions

- In Section 2 we present the two-level type theory CFTT, where the object level is first-order simply-typed and the meta level is dependently typed. The object language supports an operational semantics without runtime closures, and can be compiled with only statically known function calls. We provide a supplementary Agda formalization of the operational semantics of the object language.
- In Section 3 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 Section 4 we build a pull-based stream fusion library. Here, we demonstrate essential usage of dependent types, in providing guaranteed fusion for arbitrary 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  $\Sigma$ -types. This in turn enables a very concise definition of concatMap.

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 dependent 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 emphasize 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 [Wadler and Blott 1989].
- Functors and first-class modules in OCaml [Leroy et al. 2023] 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 2 we present the two-level type theory CFTT, where the object level is first-order simply-typed and the meta level is dependently typed. The object language supports an operational semantics without runtime closures, and can be compiled with only statically known function calls. We provide a supplementary Agda formalization of the operational semantics of the object language.
- In Section 3 we build a staged library for monad transformers [Liang et al. 1995]. 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 continuation monad is well-known in staged compilation [Bondorf 1992], and staged state monads have also been used [Carette and Kiselyov 2011; Kiselyov et al. 2004; Swadi et al. 2006]. These works used specific monads as tools for domain-specific code generation. In contrast, we propose ubiquitous staging for general-purpose monadic programming, where users can write code that looks similar to monadic code in Haskell but with deterministic and robust compilation to efficient code.
- In Section 4 we build a pull-based stream fusion library. Here, we demonstrate essential usage of dependent types, in providing guaranteed fusion for arbitrary 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  $\Sigma$ -types. This in turn enables a very concise definition of concatMap.

- 148 • We adapt the contents of the paper to typed Template Haskell [Xie et al. 2022], with some  
 149 modifications, simplifications and fewer guarantees about generated code. In particular,  
 150 Haskell does not have enough dependent types for the simple concatMap definition, but we  
 151 can still work around this limitation. We also provide a precise Agda embedding of  
 152 CFTT and our libraries as described in this paper. Here, the object theory is embedded as  
 153 a collection of postulated operations, and we can use Agda’s normalization command to  
 154 print out generated object code.

## 155 2 OVERVIEW OF CFTT

156 In the following we give an overview of CFTT features. We first review the meta-level language,  
 157 then the object-level one, and finally the staging operations which bridge between the two.  
 158

### 159 2.1 The Meta Level

160 **MetaTy** is the universe of types in the compile-time language. We will often use the term “metatype”  
 161 to refer to inhabitants of **MetaTy**, and use “metaprogram” for inhabitants of metatypes. **MetaTy**  
 162 supports dependent functions,  $\Sigma$ -types and indexed inductive types [Dybjer 1994].

163 Formally, **MetaTy** is additionally indexed by **universes** levels (orthogonally to staging), and we  
 164 have  $\text{MetaTy}_i : \text{MetaTy}_{i+1}$ . However, universe levels add noise and they are not too relevant to the  
 165 current paper, so we will omit them.

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

```
168 id : {A : MetaTy} → A → A
 169 id = λ x. x
```

170 Here, the type argument is implicit, and it gets inferred when we use the function. For example,  
 171  $\text{id}\ \text{True}$  is elaborated to  $\text{id}\{\text{Bool}\}\ \text{True}$ , where the braces mark an explicit application for the implicit  
 172 argument. Inductive types can be introduced using a Haskell-like ADT notation, or with a GADT-  
 173 style one:

```
174 data BoolM : MetaTy where
 175   data BoolM : MetaTy = TrueM | FalseM
 176           TrueM : BoolM
 177           FalseM : BoolM
```

178 Note that we added an <sub>M</sub> subscript to the type; when analogous types can be defined both on the  
 179 meta and object levels, we will sometimes use this subscript to disambiguate the meta-level version.

180 We use Haskell-like newtype notation, such as in **newtype** Wrap A = Wrap {unWrap : A}, and  
 181 also a similar notation for (dependent) record types, for instance as in

```
182 data Record = Record {field1 : A, field2 : B}.
```

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

### 185 2.2 The Object Level

186 **Ty** is the universe of types in the object language. It is itself a metatype, so we have  $\text{Ty} : \text{MetaTy}$ .  
 187 All construction and elimination rules of type formers in **Ty** stay within **Ty**. We further split **Ty** to  
 188 two sub-universes.

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

- We adapt the contents of the paper to typed Template Haskell [Xie et al. 2022], 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 precise Agda embedding of CFTT and  
 our libraries as described in this paper. 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 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,  $\Sigma$ -types and indexed inductive types [Dybjer 1994].

Formally, **MetaTy** is additionally indexed by **universe** levels (orthogonally to staging), and we  
 have  $\text{MetaTy}_i : \text{MetaTy}_{i+1}$ . However, universe levels add noise and they are not too relevant to the  
 current paper, so we will omit them.

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

```
id : {A : MetaTy} → A → A
id = λ x. x
```

Here, the type argument is implicit, and it gets inferred when we use the function. For example,  
 $\text{id}\ \text{True}$  is elaborated to  $\text{id}\{\text{Bool}\}\ \text{True}$ , 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:

```
data BoolM : MetaTy where
  data BoolM : MetaTy = TrueM | FalseM
           TrueM : BoolM
           FalseM : BoolM
```

Note that we added an <sub>M</sub> 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

```
data Record = Record {field1 : A, 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 we have  $\text{Ty} : \text{MetaTy}$ .  
All construction and elimination rules of type formers in **Ty** stay within **Ty**. 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  $\text{ValTy}$  is a sub-universe of  $\text{Ty}$ , 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,  $\text{CompTy} : \text{MetaTy}$  is the universe of *computation types*. This is also a sub-universe of  $\text{Ty}$  with implicit coercions. For now, we only specify that  $\text{CompTy}$  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 `data Nat := Zero | Suc Nat`:

```
add : Nat → Nat → Nat
add := letrec go n m := case n of
    Zero → m;
    Suc n → Suc (go n m);
go
```

Recursive definitions are 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 alike, 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  $\text{ValTy}$  and  $\text{CompTy}$ . 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  $\lambda$ -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 becomes *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  $\lambda$ -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

Since  $\text{ValTy}$  is a sub-universe of  $\text{Ty}$ , 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,  $\text{CompTy} : \text{MetaTy}$  is the universe of *computation types*. This is also a sub-universe of  $\text{Ty}$  with implicit coercions. For now, we only specify that  $\text{CompTy}$  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 `data Nat := Zero | Suc Nat`:

```
add : Nat → Nat → Nat
add := letrec go n m := case n of
    Zero → m;
    Suc n → Suc (go n m);
go
```

Recursive definitions are 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 alike, 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  $\text{ValTy}$  and  $\text{CompTy}$ . These are assumed to be erased at runtime. In the Haskell and Agda implementations they are important for guiding `instance` 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  $\lambda$ -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 becomes *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  $\lambda$ -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 sometimes we need to decide definitional equality of object programs 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 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) [Levy 1999], and there are similarities to our object language, 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 our applications. Investigating the relation between CBPV and our object language could be future work.

### 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  $\text{lift } A : \text{MetaTy}$ , 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  $\text{quote } t : \text{lift } A$ , pronounced “quote  $t$ ”. This is the metaprogram which immediately returns  $t$ .
- For  $t : \text{lift } A$ , we have  $\text{splice } t : A$ , 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 borrow this notation from MetaML [Taha and Sheard 2000].
- We have  $\text{splice } (\text{splice } t) \equiv \text{splice } t$  and  $\text{splice } (\text{quote } t) \equiv t$  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.

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:

```
let n : Nat := ~(id (10 + 10));...
```

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

```
and : lift Bool → lift Bool → lift Bool
and x y = {case ~x of True → ~y; False → False}
```

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 often with better runtime performance, avoiding the overhead of

specific translation steps that are involved in call saturation, but only specified the full translation informally.

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.

**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 sometimes we need to decide definitional equality of object programs 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 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) [Levy 1999], and there are similarities to our object language, 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 our applications. Investigating the relation between CBPV and our object language could be future work.

### 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 the following primitives:

- For  $A : \text{Ty}$ , we have  $\text{lift } A : \text{MetaTy}$ , 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  $\text{quote } t : \text{lift } A$ , pronounced “quote  $t$ ”. This is the metaprogram which immediately returns  $t$ .
- For  $t : \text{lift } A$ , we have  $\text{splice } t : A$ , 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 borrow this from MetaML [Taha and Sheard 2000].
- We have  $\text{splice } (\text{splice } t) \equiv \text{splice } t$  and  $\text{splice } (\text{quote } t) \equiv t$  as definitional equalities.

We use *unstaging* to refer to the process of extracting object code from CFTT programs, by evaluating all metaprograms in splices. This term has been sporadically used in the literature, e.g. by Odersky and Rompf [Rompf and Odersky 2012] and also by Choi et al. [Choi et al. 2011] with a somewhat different meaning in a multi-stage context. The main precursor to this paper used “staging” instead of our “unstaging” [Kovács 2022]; we use the latter because the former conflicts with other common usages of “staging”.

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

```
let n : Nat := ~(id (10 + 10));...
```

Here, `id` is used at type  $\text{lift Nat}$ . During *unstaging*, the expression in the splice is evaluated, so we get  $\sim(10 + 10)$ , which is definitionally the same as  $10 + 10$ , which is our *code* output here. Boolean

295 thunking. Consider the map function now:

```
296 map : {A B : ValTy} → (||A → ||B) → ||(List A) → ||(List B)
297 map f as = {letrec go as := case as of
298   Nil      → Nil;
299   Cons a as → Cons ~(f {a}) (go as);
300   go ~as}
301
302
```

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

## 2.4 Semantics of Staging

311 We give a short summary of the semantics of staging here. We reuse most of the semantics from  
 312 [Kovács 2022] and only mention additions and differences. The reader may refer to *ibid.* for technical  
 313 details.

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

320 The main addition in CFTT is the assumption of all inductive families in the meta-level. However,  
 321 based on the results in [Hugunin 2020], it is enough to extend the meta-level syntax in [Kovács  
 322 2022] with W-types and identity types to derive all inductive families. We consider the semantics  
 323 of these in the following.

324 2.4.1 *W-types*. W-types have a standard object-wise inductive definition in the presheaf model  
 325 [Moerdijk and Palmgren 2000]. Additionally, we need to give a logical relation interpretation  
 326 to W-types, as required for the soundness proof for staging. In Section 5.2. of [Kovács 2022] a  
 327 serialization map is used to define the relation for natural numbers, however, there is no analogous  
 328 map for W-types because we cannot build finite terms from infinitely branching trees.

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

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

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

336 Then, assuming  $\gamma^* : \Gamma^* \gamma \gamma'$ , we inductively define a relation  $W^*$  between  $W(\exists A \gamma)(\lambda \alpha. \exists B(\gamma, \alpha))$   
 337 and  $R(W_{\text{CFTT}} A B) \gamma'$ , which relates  $\text{sup } \alpha f$  to  $R_{\text{sup}} \alpha' f'$  if  $\alpha$  is related to  $\alpha'$  at type  $A$  and  $f$  is

short-circuiting is another basic use-case:

```
and : ||Bool → ||Bool → ||Bool
and x y = {case ~x of True → ~y; False → False}
```

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 often with 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 `let f as : List Nat → List Nat := ~(map (λ x. {~x + 10}) (as))`. This is  
 unstaged to a recursive definition where the mapping function is inlined into the `Cons` case  
 as `Cons a as → 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  $\|A \rightarrow \|B$ , instead of  
 $\|(A \rightarrow B)$ . The former type is often preferable to the latter; the former is a metafunction with useful  
 computational content, while the latter is merely a black box that computes object code. If we have  
 $f : \|(A \rightarrow B)$ , and  $f$  computes to  $(\lambda x. t)$ , then  $\sim f u$  is unstaged to an undesirable “administrative”  
 $\beta$ -redex  $(\lambda x. t) u$ .

2.3.1 *Comparison to defunctionalization*. Defunctionalization is a program translation which  
 represents higher-order functions using only first-order functions and inductive data [Danvy and  
 Nielsen 2001; Reynolds 1998]. The idea is to represent each closure as a constructor of an inductive  
 type, and implement closure application as a top-level first-order function which switches on all  
 possible closure constructors. Hence, like our unstaging algorithm, defunctionalization produces  
 first-order code from higher-order code. We summarize the differences.

- Defunctionalization mainly supports *compiler optimizations*, by making closures transparent to analysis. It is not an optimization by itself; it does not change the amount of dynamic control flow or dynamic allocations in a program. A dynamic closure call becomes a dynamic case-switch and a heap-allocated “functional” closure becomes a heap-allocated inductive constructor.
- Staging mainly supports optimization by *programmers* by providing control over code generation. Unstaging does not translate higher-order functions to anything in particular. Instead, unstaging is the execution of higher-order metaprograms which produce first-order programs.

## 2.4 Staged Semantics of CFTT

To obtain an unstaging algorithm, together with its correctness, the main task is to extend [Kovács 2022] with identity types and inductive families in the metalanguage, since we use those types in this paper. For inductive types, it is the most sensible to add W-types, since all inductive families can be faithfully derived from them in our setting [Hugunin 2020].

This however requires a substantial amount of technical background from [Kovács 2022], so we relegate it to Appendix A. Here we only give an overview.

related to  $f'$  at type  $B[\alpha'] \Rightarrow W_{\text{CFTT}} A B$ . The relation at the non-dependent function type  $\dashv \dashv$  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  $W^\ast$  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 equations 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 4.3 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.
- Before staging, translate the intensional identity to extensional identity. This amounts to erasing all identity proofs and transports.
- Perform staging after erasure of identity proofs, i.e. model extensional identity in the semantics.

In the presheaf model, we interpret the extensional identity in a standard way, as strict equality of sections of presheaves, and soundness for extensional identity becomes trivial to show. Hence, the soundness property for our overall staging pipeline becomes the following: the staging output is definitionally equal to the identity-erased version of the staging input.

### 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  $\parallel A \rightarrow \parallel B$  is usually preferable to  $\parallel(A \rightarrow B)$ . 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:

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

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 [Jones et al. 1993]. 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:

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

Here we overload Haskell-style product type notation, both for types and the pair constructor, at both levels. There is a problem with this conversion though: *up* uses  $x : \parallel(A, B)$  twice, which can increase code size and duplicate runtime computations. For example, *down* (*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  $\parallel(A, B)$ , and then only use projections of the newly bound variable. This is called *let-insertion* in staged compilation.

- The input of unstaging is a CFTT term with an object-level type, which only depends on object-level free variables. The output is a term in the object theory. In the other direction, there is an evident embedding of object-theoretic terms as CFTT terms.
- *Soundness* means that unstaging followed by embedding is the identity map, up to conversion. In other words, unstaging respects conversion of CFTT terms.
- *Stability* means that embedding followed by unstaging is the identity map. In other words, unstaging has no action on terms which are already purely object-level (i.e. contain no splices).

In [Kovács 2022], there is one more property, called *strictness*, which means that unstaging does not perform  $\beta\eta$ -conversion in the object theory, but this holds trivially in our case since we do not have such conversion rules.

However, for CFTT we slightly depart from the above notion of soundness, for the following reason. We would like to assume certain propositional identities as *axioms* in the metalanguage, without blocking unstaging as a computation. We use such an axiom to good effect in Section 4.3.

Axioms clearly block computation, so they are incompatible with the mentioned notion of soundness. But we are only interested in equations which hold definitionally during unstaging. Hence, we can do the following: first, we erase all identity proofs and their transports from CFTT programs, then we proceed with unstaging as in [Kovács 2022]. This erasure can be formalized as a syntactic translation [Boulier et al. 2017] from CFTT to CFTT extended with the equality reflection principle, which says that propositional equality implies definitional equality (see e.g. [Hofmann 1995]). Thus, soundness of staging for CFTT holds 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 do a preliminary overview before getting to monads. 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 [Jones et al. 1993, Chapter 12].

Translation from  $\parallel A \rightarrow \parallel B$  to  $\parallel(A \rightarrow B)$  is a basic binding-time improvement, as an  $\eta$ -expansion [Danvy et al. 1996]. The two types are equivalent during unstaging, up to the runtime equivalence of object programs, and we can convert back and forth in CFTT, the same way as in MetaML [Taha and Sheard 2000, Section 9]:

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

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.

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:

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

393 But it is impossible to use let-insertion in `up` because the return type is in `MetaTy`, and we cannot  
 394 introduce object binders in meta-level code. Fortunately, there is a principled solution.

### 3.2 The Code Generation Monad

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

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

402 This is a monad in `MetaTy` in a standard sense:

```
403 instance Monad Gen where
 404   return a = Gen $ λ k. k a
 405   ga ≫= f = Gen $ λ k. unGen ga (λ a. unGen (f a) k))
```

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

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

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

418 We can let-bind object expressions in `Gen`:

```
420 gen : {A : Ty} → ↑A → Gen (↑A)
421 gen a = Gen $ λ k. {let x : A := ~a; ~(k (x))}
```

423 And also recursive definitions of computations:

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

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

<code>up : ↑(A, B) → Gen (↑A, ↑B)</code> <code>up x = do x ← gen x</code> <code>return ((fst ~x), (snd ~x))</code>	<code>down : Gen (↑A, ↑B) → ↑(A, B)</code> <code>down x = do (a, b) ← x</code> <code>return ((~a, ~b))</code>
--	---

435 Working in `Gen` is convenient, since we can freely generate object code and also have access  
 436 to the full metalanguage. Also, the point of staging is that eventually all metaprograms will be  
 437 used for the purpose of code generation, so we eventually run `Gen` all of our actions. So why not  
 438 just always work in `Gen`? The implicit nature of `Gen` may make it harder to reason about the size  
 439 and content of generated code. This is a bit similar to the `IO` monad in Haskell, where eventually  
 440 everything needs to run in `IO`, but we may not want to write all of our code in `IO`.

Here we overload Haskell-style product type notation, both for types and the pair constructor, at both levels. There is a problem with this conversion though: `up` uses `x : ↑(A, B)` twice, which can increase code size and duplicate runtime computations. For example, `down (up {fx})` is staged to `{fst (f x), snd (f x)}`. It would be safer to first let-bind an expression with type `↑(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.

### 3.2 The Code Generation Monad

A principled solution to the previous issue is to write code generators in continuation-passing style, as first proposed by Bondorf [Bondorf 1992]. This can be structured as a continuation monad, and later work used explicit monadic notation for it [Carette and Kiselyov 2011; Kiselyov et al. 2004; Swadi et al. 2006]. Our definition is as follows:

```
newtype Gen (A : MetaTy) = 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. unGen (f a) k))
```

In [Kiselyov et al. 2004], [Carette and Kiselyov 2011] and [Swadi et al. 2006], the answer type is a parameter to the continuation monad, while we have it as polymorphic. We have found that polymorphic answer types are more convenient, because we do not have to anticipate the type of the code output when are defining a code generator. Hence, we can “run” actions as follows:

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

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 Agda and Haskell implementations, with modest technical differences. From the Monad instance, the Functor and Applicative instances can be also derived. We reuse (`<$>`) and (`<*>`) for applicative notation as well.

We can let-bind object expressions in `Gen`:

```
gen : {A : Ty} → ↑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:

<code>up : ↑(A, B) → Gen (↑A, ↑B)</code> <code>up x = do x ← gen x</code> <code>return ((fst ~x), (snd ~x))</code>	<code>down : Gen (↑A, ↑B) → ↑(A, B)</code> <code>down x = do (a, b) ← x</code> <code>return ((~a, ~b))</code>
--	---

### 442 3.3 Monads

443 Let us start with the `Maybe` monad. We have `data Maybe A := Nothing | Just A`, and `Maybe` itself  
 444 is available as a `ValTy → ValTy` metafunction. However, we cannot directly fashion a monad out  
 445 of `Maybe`, since we do not have enough type formers in `ValTy`. We could try to use the following  
 446 type for binding:  
 447

```
(>=) : ⌞(Maybe A) → (⌞A → ⌞(Maybe B)) → ⌞(Maybe B)
```

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

451 Instead, our strategy is to only use proper monads in `MetaTy`, and convert between object types  
 452 and meta-monads when necessary, as a form of binding-time improvement. We define a class for  
 453 this conversion:

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

457 Assume that `MaybeM` is the standard meta-level monad, and `MaybeTM` is the standard monad  
 458 transformer, defined as follows:

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

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

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

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

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

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

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

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

```
473   gen' : ⌞A → MaybeTM Gen (⌞A)
 474   gen' a = lift (gen a)
```

475 However, it is more convenient to proceed in the style of Haskell’s monad transformer library `mtl`  
 476 [`mtl` developers 2024], and have a class for monads that support code generation:

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

Working in `Gen` is convenient, since we can freely generate object code and also have access to  
 the full metalanguage. Also, the point of staging is that eventually all metaprograms will be used  
 for the purpose of code generation, so we eventually run `Gen` all of our actions. So why not just  
 always work in `Gen`? The implicit nature 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 [Jones 2001], where  
 eventually everything needs to run in `IO`, but we may not want to write all of our code in `IO`.

### 3.3 Monads

First, following the style of Haskell’s monad transformer library `mtl` [mtl developers 2024], we  
 define a class for monads that support code generation.

```
class Monad M ⇒ MonadGen M where
  liftGen : Gen A → M A
```

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)
```

In the rest of this paper, we shall only mention `ReaderT`, `MaybeT` and `StateT` as monad transformers.  
 For all of these, we can define the `MonadGen` instance simply by `liftGen = lift ∘ liftGen`.

Let us look at the `Maybe` monad now. 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 polymorphism in `ValTy`, nor can we store functions inside  
`Maybe`. We could try to use the following type for binding:

```
(>=) : ⌞(Maybe A) → (⌞A → ⌞(Maybe B)) → ⌞(Maybe B)
```

This would work, but the definition would necessarily use 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, defined as follows:

```
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 ma = MaybeTM $ Gen $ λ k.
    (case ~ma 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))
```

491 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:

```
492 gen   : MonadGen M ⇒ □A → M (□A)
493 genRec : MonadGen M ⇒ (□A → □A) → M (□A)
```

494 3.3.1 *Case splitting in monads.* We often want to case-split on object-level data inside a monadic action, and perform different actions in different branches. At first, this may seem problematic, because we cannot directly compute metaprograms from object-level case splits. Fortunately, with a little bit more work, this is actually possible in any `MonadGen`, for any value type.

495 We demonstrate this for lists. 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:

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

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

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

500 Now, in any `MonadGen`, assuming `as : □(List A)`, we may write

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

502 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:

```
503 class Split (A : ValTy) where
504   SplitTo : MetaTy
505   split   : □A → Gen SplitTo
```

506 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 when working in a `MonadGen`, we can write `case as of Nil → ...; Cons a as → ...`, binding `a : □A` and `as : □(List A)` in the `Cons` case.

### 3.4 Monad Transformers

507 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:

```
508 newtype Identity A := Identity {runIdentity : A}
509 instance Improve Identity Gen where
510   up x   = return {runIdentity ~x}
511   down x = unGen x $ λ a. {Identity ~a}
```

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

```
513 newtype MaybeT (M : ValTy → Ty) (A : ValTy) := MaybeT {runMaybeT : M (Maybe 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 unstaging time. Thus, the above code is computed to

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

We can also use `gen` for let-insertion in `MaybeT Gen`, since it has a `MonadGen` instance.

3.3.1 *Case splitting in monads.* We often want to case-split on object-level data inside a monadic action, and perform different actions in different branches. At first, this may seem problematic, because we cannot directly compute metaprograms from object-level case splits. Fortunately, with a little bit more work, this is possible in any `MonadGen`, for any value type, using a variation of the so-called “trick” [Danvy et al. 1996].

We demonstrate this for lists. 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`, assuming `as : □(List A)`, 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 `splitting` function. We adopt this in the rest of the paper, so when working in a `MonadGen`, we can write `case as of Nil → ...; Cons a as → ...`, binding `a : □A` and `as : □(List A)` in the `Cons` case.

### 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   = return {runIdentity ~x}
  down x = unGen x $ λ a. {Identity ~a}
```

We 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 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  $\lambda$  `case` ... to shorten  $\lambda x. \text{case } x \text{ of } \dots$

In the meta level, we can reuse essentially all definitions from `mt1`. From there, 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 only present `StateT` and `ReaderT`.

We start with `StateT`. We assume `StateTM` as the standard meta-level definition. The object-level `StateT` has type  $(S : \text{ValTy})(F : \text{ValTy} \rightarrow \text{Ty})(A : \text{ValTy}) \rightarrow \text{ValTy}$ ; 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 use the `modify` function that we already have in `StateM`, 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  $(\sim x + \sim x)$ , 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

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 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  $\lambda$  `case` ... to shorten  $\lambda x. \text{case } x \text{ of } \dots$

On the meta level, we can borrow essentially all definitions from `mt1`. 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 only present `StateT` and `ReaderT`.

We start with `StateT`. We assume `StateTM` as the standard meta-level definition. The object-level `StateT` has type  $(S : \text{ValTy})(F : \text{ValTy} \rightarrow \text{Ty})(A : \text{ValTy}) \rightarrow \text{ValTy}$ ; 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 use the `modify` function that we already have in `StateM`, 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  $(\sim x + \sim x)$ , and this happens at unstagin 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 unstaged to

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

which duplicates the evaluation of  $x + x$ . The duplication can be avoided by let-binding the result in the object language. A similar phenomenon happens with the local function in `Reader`. So we

in Reader. So we define “stricter” versions of these operations. We also return  $\mathbb{I}()$  from actions instead of  $()$  – the former is more convenient, because the down operation can be immediately used on it.

```

589 put' : (MonadState (||S) M, MonadGen M) => ||S -> M (||())
590 put's = do {s <- gen s; put s; return ()}
591
592 modify' : (MonadState (||S) M, MonadGen M) => (||S -> ||S) -> M (||())
593 modify' f = do {s <- get; put' (f s)}
594
595
596 local' : (MonadReader (||R) M, MonadGen M) => (||R -> ||R) -> MA -> MA
597 local' f ma = do {r <- ask; r <- gen (f r); local (λ _ . r) ma}
598
599 Now,
600     down $ do {modify' (λ x. (¬x + ¬x)); modify' (λ x. (¬x + ¬x))}

601 is staged to
602     (λ x. let x := x + x; let x := x + x; ((), x))

603 3.5 Joining Control Flow in Monads
604 There is a deficiency in our library so far. Assuming b : ||Bool, consider:
605
606     do (case b of True -> put' 10; False -> put' 20)
607         modify' (λ x. (¬x + 10))

608 The modify' (λ x. (¬x + 10)) 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 both branches with the same action. If we have multiple Boolean case splits sequenced after each other, that yields exponential code size. Right now, we can fix this by let-binding the problematic action:
```

```

609     do action <- gen $ down (case b of True -> put' 10; False -> put' 20)
610         up action
611         modify' (λ x. (¬x + 10))

612 This removes code duplication by round-tripping through an object-level let. This solution is fairly
613 good in a state monad, where down only introduces a runtime pair constructor, and it is feasible
614 to compile object-level pairs as unboxed data, without overheads. However, for Maybe, down
615 introduces a runtime Just or Nothing, and up introduces a runtime case split. A better solution
616 would be to introduce two let-bound join points before the offending case, one for returning a Just
617 and one for returning Nothing, but fusing away the actual runtime constructors. So, for example,
618 we would like to produce object code which looks like the following:
619
620     let joinJust n := ...
621     let joinNothing () := ...
622     case x == 10 of
623         True -> joinJust (x + 10)
624         False -> joinNothing ()
```

define “stricter” versions of these operations. We also return  $\mathbb{I}()$  from actions instead of  $()$  – the former is more convenient, because the down operation can be immediately used on it.

```

589 put' : (MonadState (||S) M, MonadGen M) => ||S -> M (||())
590 put's = do {s <- gen s; put s; return ()}
591
592 modify' : (MonadState (||S) M, MonadGen M) => (||S -> ||S) -> M (||())
593 modify' f = do {s <- get; put' (f s)}
594
595
596 local' : (MonadReader (||R) M, MonadGen M) => (||R -> ||R) -> MA -> MA
597 local' f ma = do {r <- ask; r <- gen (f r); local (λ _ . r) ma}

598 Now,
599     down $ do {modify' (λ x. (¬x + ¬x)); modify' (λ x. (¬x + ¬x))}

600 is unstaged to
601     (λ x. let x := x + x; let x := x + x; ((), x)).

602 3.5 Joining Control Flow in Monads
603 There is a deficiency in our library so far. Consider:
604
605     f : Bool -> StateT Nat (MaybeT Identity) ()
606     f b := ¬(down $ do
607         case (b) of True -> put' 10; False -> put' 20)
608         modify' (λ x. (¬x + 10))

609 This is unstaged to the following (omitting newtype wrappers):
610
611     f b s := case b of
612         True -> let s := s + 10; Just (((), s))
613         False -> let s := s + 20; Just (((), s))

614 Notice that the final state modification gets inlined into both branches. This duplication follows
615 from the definition of monadic binding in Gen and the split function in the desugaring of case.
616 Code generation is continued in both branches with the same action. If we have multiple Boolean
617 case splits sequenced after each other, that yields exponential code size. A possible fix is to let-bind
618 the branching action:
619
620     f b := ¬(down $ do
621         act <- gen $ down $ case (b) of True -> put' 10; False -> put' 20)
622         up act
623         modify' (λ x. (¬x + 10))
```

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  $\{R : \text{ValTy}\} \rightarrow (A \rightarrow \parallel R) \rightarrow \parallel R$ . Here, if A is a finite sum, we can rearrange  $A \rightarrow \parallel R$  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 4, 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).

```
686 USOP : MetaTy           ElSOP : USOP → MetaTy
687 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. USOP is closed under value types, finite product types and finite sum types. For instance, we have  $\text{Either}_{\text{SOP}} : \text{USOP} \rightarrow \text{USOP} \rightarrow \text{USOP}$  together with  $\text{Left}_{\text{SOP}} : \text{ElSOP A} \rightarrow \text{ElSOP} (\text{Either}_{\text{SOP}} A B)$ ,  $\text{Right}_{\text{SOP}} : \text{ElSOP B} \rightarrow \text{ElSOP} (\text{Either}_{\text{SOP}} A B)$  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:

```
688 class IsSOP (A : MetaTy) where
689   Rep : USOP
690   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  $A \rightarrow \text{ElSOP Rep}$ , and write  $\text{rep}^{-1}$  for its inverse. Now we define an isomorphic presentation of  $\text{ElSOP A} \rightarrow \parallel R$  as a product of object-level functions:

```
691 FunSOPT : USOP → Ty → MetaTy
692 FunSOPT Nil      R = ()
693 FunSOPT (Cons A B) R = (parallel(foldr (→) RA), FunSOPT B R)
694
695 tabulate : (ElSOP A → ∥R) → FunSOPT A R
696 index    : FunSOPT 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  $\text{FunSOP}_{\text{T}}$ :

```
697 genFunSOPT : {A : USOP} → FunSOPT A R → FunSOPT A R
698 genFunSOPT {Nil}     ()    = return ()
699 genFunSOPT {Cons _A} {f, fs} = (.) <$> gen f <> genFunSOPT {A} fs
```

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

```
703 class Monad M ⇒ MonadJoin M where
704   join : IsSOP A ⇒ MA → MA
```

However, this yields suboptimal code:

```
f b s :=
let act s := case b of True → let s := 10; Just(((), s))
                           False → let s := 20; Just(((), s));
case act s of
  Just (_ , s) → let s := s + s; Just (((), s))
  Nothing → Nothing
```

This solution is fairly good when we only have State or Reader effects, where 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:

```
f b s :=
let joinNothing s := Nothing;
let joinJust    s := let s := s + s; Just (((), s));
case b of
  True → let s := 10; joinJust s
  False → let s := 20; joinJust s
```

Here, joinNothing happens to be dead code, but it is easy to clean that up in downstream compilation. 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  $\{R : \text{ValTy}\} \rightarrow (A \rightarrow \parallel R) \rightarrow \parallel R$ . Here, if A is a finite sum, we can rearrange  $A \rightarrow \parallel R$  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 4, so we develop SOP-s. SOPs have been used in generic and staged programming; see e.g. [de Vries and Löh 2014a] and [Pickering et al. 2020]. 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. For example,  $\text{ElSOP} [[\text{Bool}, \text{Bool}], []]$  (using Haskell-like list notation) is isomorphic to  $\text{Either}(\parallel \text{Bool}, \parallel \text{Bool})()$ . USOP is closed under value types, finite product types and finite sum types. For instance, we have  $\text{Either}_{\text{SOP}} : \text{USOP} \rightarrow \text{USOP} \rightarrow \text{USOP}$  together with  $\text{Left}_{\text{SOP}} : \text{ElSOP A} \rightarrow \text{ElSOP} (\text{Either}_{\text{SOP}} A B)$ ,  $\text{Right}_{\text{SOP}} : \text{ElSOP B} \rightarrow \text{ElSOP} (\text{Either}_{\text{SOP}} A B)$  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
```

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

```
688 instance MonadJoin Gen where
689   join ma = Gen $ λk. runGen $ do
690     joinPoints ← genFunSOP (tabulate (k ∘ rep-1))
691     a ← ma
692     return $ index joinPoints (rep a)
```

693 Here we first convert  $k : A \rightarrow \mathbb{I}R$  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 694 MonadJoin instances are straightforward. In  $\text{StateT}_M$ , we need the extra IsSOP S constraint because 695 S is returned as a result value.

```
696 instance (MonadJoin M) ⇒ MonadJoin (MaybeTM M) where
697   join (MaybeTM ma) = MaybeTM (join ma)
698 instance (MonadJoin M) ⇒ MonadJoin (ReaderTM R M) where
699   join (ReaderTM ma) = ReaderTM (join ∘ ma)
700 instance (MonadJoin M, IsSOP S) ⇒ MonadJoin (StateTM S M) where
701   join (StateTM ma) = StateTM (join ∘ ma)
```

702 Now, whenever the return value of a case-splitting action is a sum of products of values (this 703 includes just returning an object value, which is by far the most common situation), we can use 704 join \$ case x of ... to eliminate code duplication, without creating runtime sum types.

### 3.6 Code Example

705 We look at a slightly larger code example. We define annotated binary trees as

```
706 data Tree A := Leaf | Node A (Tree A) (Tree A).
```

707 We write fail : M A for returning Nothing in any monad transformer stack that contains MaybeT. 708 We define a function which traverses a tree and replaces values with values taken from a list. If the 709 tree contains 0, we throw an error. If we run out of list elements, we leave values unchanged.

```
710 letrec f : Tree Nat → StateT (List Nat) (MaybeT Identity) (Tree Nat)
711   ft := ~(down $ case t of
712     Leaf → return (Leaf)
713     Node n | r → do
714       case (~n == 0) of True → fail
715                         False → return ()
716       ns ← get
717       n ← join $ case ns of Nil → return n
718                           Cons n ns → do {put ns; return n}
719       l ← up (f ~l)
720       r ← up (f ~r)
721       return (Node ~n ~l ~r))
```

Above,  $A = \text{Rep}$  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  $A \rightarrow \text{El}_{\text{SOP}} \text{Rep}$ , and write  $\text{rep}^{-1}$  for its inverse.

In 2LTT literature, a metatype A is called *cofibrant* if  $A \rightarrow B$  is isomorphic to an object type whenever B is isomorphic to an object type [Annenkov et al. 2019]. We conjecture that our IsSOP types are exactly the cofibrant types in this sense. However, “cofibrant” is not very descriptive in our setting, so we shall keep writing IsSOP.

Accordingly, we define an isomorphic presentation of  $\text{El}_{\text{SOP}} A \rightarrow \mathbb{I}R$  as a product of object-level functions:

```
FunSOP : Usop → Ty → MetaTy
FunSOP Nil      R = ⊥
FunSOP (Cons A B) R = (⊗(foldr (→) RA), FunSOP B R)

tabulate : (ElSOP A → I R) → FunSOP A R
index    : FunSOP A R → (ElSOP A → I R)
```

We omit here the definitions of tabulate and index. We will also need to let-bind all functions in a  $\text{Fun}_{\text{SOP}}$ :

```
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 ⇒ MonadJoin M where
  join : IsSOP A ⇒ MA → MA
```

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

```
instance MonadJoin Gen where
  join ma = Gen $ λk. runGen $ do
    joinPoints ← genFunSOP (tabulate (k ∘ rep-1))
    a ← ma
    return $ index joinPoints (rep a)
```

Here we first convert  $k : A \rightarrow \mathbb{I}R$  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 734 MonadJoin instances are straightforward. In  $\text{StateT}_M$ , we need the extra IsSOP S constraint because 735 S is returned as a result value.

```
instance (MonadJoin M) ⇒ MonadJoin (MaybeTM M) where
  join (MaybeTM ma) = MaybeTM (join ma)
instance (MonadJoin M) ⇒ MonadJoin (ReaderTM R M) where
  join (ReaderTM ma) = ReaderTM (join ∘ ma)
instance (MonadJoin M, IsSOP S) ⇒ MonadJoin (StateTM S M) where
  join (StateTM ma) = StateTM (join ∘ ma)
```

738 Here, all of the `case` matches are done on object-level values, so they are all desugared to applications  
 739 of `split`.

- Note the join: without it, the recursive Node traversal code would be inlined in both case branches.
- We do not use `join` when checking  $\{\sim n == 0\}$ . Here, it happens to be superfluous, since the fail “destroys” all subsequent code generation in the branch, by short-circuiting at the meta-level.
- Also, we can use `put` instead of the “stricter” `put’` because the state modification immediately gets sequenced by the enclosing join point.
- To make object-level recursive calls, we just need to wrap them in up.

740 Here, we intentionally tuned the definition for a nice-looking staging output. However, we could also just default to “safe” choices everywhere: we could use a joint point in every `case` with two or more branches, and use the strict state modifications everywhere. This would only add a small amount of administrative noise that can be easily cleaned up by the downstream compiler.

741 Let us look at the staging output. We omit newtype wrappers and use nested pattern matching on pairs, but otherwise this is exactly the code that we get in our Agda implementation.

```
742 letrec f : Tree Nat → StateT (List Nat) (MaybeT Identity) (Tree Nat)
743   f := λ ns. case t of
744     Leaf → Just (Leaf, ns)
745     Node n l r → case (n == 0) of
746       True → Nothing
747       False →
748         let joinNothing := λ _ . Nothing;
749         let joinJust := λ n ns. case f l ns of
750           Nothing → Nothing
751           Just (l, ns) → case f r ns of
752             Nothing → Nothing
753             Just (r, ns) → Just (Node n l r, ns);
754         case ns of
755           Nil → joinJust n ns
756           Cons n' ns' → joinJust n' ns'
```

757 The only flaw in this code is the `joinNothing`, which is never called, because we joined an action  
 758 that never throws an error. But this is also very easy to clean up after staging.

### 3.7 Discussion

772 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 matching 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.

```
letrec f : Tree Nat → StateT (List Nat) (MaybeT Identity) (Tree Nat)
  f t := ~(down $ case t of
    Leaf → return (Leaf)
    Node n l r → do
      case (n == 0) of
        True → fail
        False → return ()
      ns ← get
      n ← join $ case ns of
        Nil → return n
        Cons n' ns' → do { put ns; return n' }
      l ← up {f ~l}
      r ← up {f ~r}
      return (Node ~n ~l ~r))
```

Fig. 1. Monadic source code

Now, the following definition yields the previously seen code output with `joinNothing` and `joinJust`.

```
f b := ~(down $ do
  join $ case (b) of
    True → put' {10}; False → put' {20}
    modify' (λ x. {~x + ~x}))
```

It might make sense to also have a simple “desugaring” rule that inserts a `join` whenever we have a object-level case split with more than one branch. At worst, this generates some noise and dead code that is easy to remove by conservative code optimization.

### 3.6 Code Example

We look at a slightly larger code example in Figure 1. We define annotated binary trees as

```
data Tree A := Leaf | Node A (Tree A) (Tree A).
```

We write `fail : M A` for returning `Nothing` in any monad transformer stack that contains `MaybeT`. We define a function which traverses a tree and replaces values with values taken from a list. If the tree contains 0, we throw an error. If we run out of list elements, we leave values unchanged. Here, all of the `case` matches are done on object-level values, so they are all desugared to applications of `split`.

- Note the join: without it, the recursive Node traversal code would be inlined in both case branches.
- We do not use `join` when checking  $\{\sim n == 0\}$ . Here, it happens to be superfluous, since the fail “destroys” all subsequent code generation in the branch, by short-circuiting at the meta-level.
- Also, we can use `put` instead of the “stricter” `put’` because the state modification immediately gets sequenced by the enclosing join point.
- To make object-level recursive calls, we just need to wrap them in up.

Here, we intentionally tuned the definition for a nice-looking staging output. However, we could also just default to “safe” choices everywhere: we could use a joint point in every `case` with two or

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 noise can be mitigated very effectively in a native CFTT implementation. It was demonstrated in a 2LTT implementation in [Kovács 2022] that almost all quotes and splices are unambiguously inferable, 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 such 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.

### 3.8 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.
- 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* [Eisenberg and Weirich 2012] to emulate dependent types. This adds significant noise. Also, in IsOP 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; see e.g. [Hinze et al. 2010]. 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) → IsOP S ⇒ Gen S → (S → Gen (Step S A)) → Pull A
```

```
letrec f : Tree Nat → StateT (List Nat) (MaybeT Identity) (Tree Nat)
f := λ t ns. case t of
  Leaf → Just (Leaf, ns)
  Node n l r → case (n == 0) of
    True → Nothing
    False →
      let joinNothing = λ _ . Nothing;
      let joinJust := λ n ns. case f l ns of
        Nothing → Nothing
        Just (l, ns) → case f r ns of
          Nothing → Nothing
          Just (r, ns) → Just (Node n l r, ns);
      case ns of
        Nil → joinJust n ns
        Cons n ns → joinJust n ns
```

Fig. 2. Unstaged monadic code

more branches, and use the strict state modifications everywhere, and leave it up to the downstream compiler to clean up the code.

We show the unstaging output on Figure 2. We omit newtype wrappers and use nested pattern matching on pairs, but otherwise this is exactly the code that we get in our Agda implementation. Again, the only flaw in this code is the dead binding for joinNothing.

### 3.7 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 matching 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 noise can be mitigated very effectively in a native CFTT implementation. It was demonstrated in a 2LTT implementation in [Kovács 2022] that almost all quotes and splices are unambiguously inferable, 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 such inference, monadic code in CFTT would look only slightly 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

834 In the Pull constructor, S is the type of the internal state which is required to be a sum-of-products  
835 of value types by the `IsSumVS S` constraint. Borrowing terminology from Coutts [Coutts 2011], we  
836 call each product in the state a *state shape*.

837 The next field with type Gen S is the initial state, while transitions are represented by the  
838  $S \rightarrow \text{Gen}(\text{Step } S A)$  field. Possible transitions are stopping (Stop), transitioning to a new state  
839 while outputting a value (Yield) and making a silent transition to a new state (Skip).

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

```
842 repeat : A → Pull A
843 repeat a = Pull () (return ()) (λ _ → return $ Yield a ())
```

```
844 ( $\langle\langle\rangle_{\text{Pull}}$ ) : IsSOP A ⇒ Pull (A → B) → Pull A → Pull B
845 ( $\langle\langle\rangle_{\text{Pull}}$ ) (Pull S seed step) (Pull S' seed' step') =
846 Pull (S, S', Maybe A) ((, ) <$> seed <> ((, Nothing) <$> seed') $ λ case
847 (s, s', Just a) → step s ≈≈ λ case
848 Stop → return Stop
849 Skip s → return $ Skip (s, s', Just a)
850 Yield f s → return $ Yield (f a) (s, s', Nothing)
851 (s, s', Nothing) → step' s' ≈≈ λ case
852 Stop → return Stop
853 Skip s' → return $ Skip (s, s', Nothing)
854 Yield a s' → return $ Skip (s, s', Just a)
```

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

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

```
862 empty : Pull A
863 empty = Pull () (return ()) (λ _ → return Stop)
```

```
864 ( $\langle\rangle$ ) : Pull A → Pull A → Pull A
865 ( $\langle\rangle$ ) (Pull S seed step) (Pull S' seed' step') = Pull (Either S S') (Left <$> seed) λ case
866 Left s → step s ≈≈ λ case Stop → (Skip ∘ Right) <$> seed'
867 Skip s → return $ Skip (Left s)
868 Yield a s → return $ Yield a (Left s)
869 Right s' → step s' ≈≈ λ case Stop → return Stop
870 Skip s' → return $ Skip (Right s')
871 Yield a s' → return $ Yield a (Right s')
```

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.

### 3.8 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.
- 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 [Eisenberg and Weirich 2012] 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; see e.g. [Hinze et al. 2010]. 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 `IsSOP S` constraint. The next field with type Gen S is the initial state, while transitions are represented by the  $S \rightarrow \text{Gen}(\text{Step } S A)$  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).

Our definition roughly follows the non-staged stream fusion setup of Coutts [Coutts 2011]; the difference is the addition of IsSOP and Gen in our version. Also, borrowing terminology from Coutts, we call each product in the state a *state shape*.

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 state `S` 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 treatment of join points, we tabulate the `S → Gen Step SA` 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 [Yallop and Kiselyov 2019]. 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 `ComPly` with finite products of computations, i.e. assume `() : CompTy` and `(,) : CompTy → CompTy → 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 → B, A → B)` can be always compiled to a pair of mutual functions.

We redefine the previous `FunSOP†` to return a computation type instead:

```

FunSOP† : USOP → Ty → CompTy
FunSOP† Nil      R = ()
FunSOP† (Cons A B) R = (foldr (→) RA, FunSOP† BR)

tabulate : (ElSOP A → ∏R) → ∏(FunSOP† A R)
index   : ∏(FunSOP† A R) → (ElSOP A → ∏R)

```

Even for join points, this is just as efficient as the previous `FunSOP†` 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.

```

foldr : {A : MetaTy}{B : Ty} → (A → ∏B → ∏B) → ∏B → Pull A → ∏B
foldr {A} {B} (Pull S seed step) f b =
  letrec fs : FunSOP† (Rep {S}) B := ~(tabulate $ λ s. unGen (step (rep-1 s)) $ λ case
  Stop    → b
  Skip s → index {fs} (rep s)
  Yield a s → f a (index {fs} (rep s)));
  ~{unGen seed $ λ s. index {fs} (rep s)}}

```

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

```

foldl : {A : MetaTy}{B : ValTy} → (∏B → A → ∏B) → ∏B → Pull A → ∏B
foldl f b as = ~(foldr (λ a g. (λ b. ~g ~{f (b) ~a})) (λ b. b) as) ~b

```

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 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 non-staged 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 state `S` 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

932 Note that since we abstract B in a runtime function, it must be a value type. Here, each Stop and  
 933 Yield in the transition function gets interpreted as a  $\lambda$ -expression in the output. However, those  
 934  $\lambda$ -s will be lifted out in the scope, yielding a proper mutually tail-recursive definition with an  
 935 accumulator for B. Contrast this to the GHC base library, where foldl for lists is also defined from  
 936 foldr, to enable push-based fusion, but where a substantial *arity analysis* is used in the compiler to  
 937 (incompletely) eliminate the intermediate closures [Breitner 2014].

### 4.3 concatMap for Streams

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

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

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

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

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

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

954 Above,  $\{a : \text{A}\} \rightarrow \text{IsSOP}(\text{B a})$  is a universally quantified instance constraint; this can be also  
 955 represented in Agda. The definition of concatMap is as follows.

```
956 concatMap : IsSOP A ⇒ (A → Pull B) → Pull A → Pull B
957 concatMap {A} {B} f (Pull S seed step) ≡
958   Pull (S, Maybe (Σ A (projS ∘ f))) (λ _ . <$> seed <*> return Nothing) $ λ case
959     (s, Nothing) → step s ≈≈ λ case
960       Stop → return Stop
961       Skip s → return $ Skip (s, Nothing)
962       Yield a s → do {s' ← projSeed (f a); return $ Skip (s, Just (a, s'))}
963     (s, Just (a, s')) → projStep (f a) s' ≈≈ λ case
964       Stop → return $ Skip (s, Nothing)
965       Skip s' → return $ Skip (s, Just (a, s'))
966       Yield b s' → return $ Yield b (s, Just (a, s'))
```

967 Here, Nothing marks the states where we are in the “outer” loop, running the Pull A stream until  
 968 we get its next value. Just marks the states of the “inner” loop, where we have a concrete  $a : \text{A}$   
 969 value and we run the  $(f a)$  stream until it stops. In the inner loop, the machine state type depends  
 970 on the  $a : \text{A}$  value, hence the need for  $\Sigma$ . How do we get IsSOP for  $\Sigma$ ? The key observations are:

- Metaprograms cannot inspect the structure of object terms.

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 [Yallop and Kiselyov 2019]. In contrast to this work, 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  $(\text{A}, \text{B})$  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  $(\text{A} \rightarrow \text{B}, \text{A} \rightarrow \text{B})$  can be always compiled to a pair of mutual  
 functions.

We redefine the previous  $\text{FunSOP}$  to return a computation type instead:

```
FunSOP : USOP → Ty → CompTy
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)
```

Even for join points, this is just as efficient as the previous  $\text{FunSOP}$  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.

```
foldr : {A : MetaTy}{B : Ty} → (A → ||B → ||B) → ||B → Pull A → ||B
foldr {A} {B} f b (Pull S seed step) = {
  letrec fs : FunSOP (Rep {S}) B = ~ (tabulate $ λ s. unGen (step (rep-1 s)) $ λ case
    Stop → b
    Skip s → index {fs} (rep s)
    Yield a s → f a (index {fs} (rep s));
    ~ (unGen seed $ λ s. index {fs} (rep s)))
  in
    foldr f b (Pull S seed step)
```

This foldr is quite flexible because it can produce object terms with computation types. For instance,  
 we can define foldl from foldr:

```
foldl : {A : MetaTy}{B : ValTy} → (||B → A → ||B) → ||B → Pull A → ||B
foldl f b as = ~ (foldr (λ a g. (λ b. ~g ~ (f (b) ~a))) (λ b. b) as) ~b
```

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  $\lambda$ -expression in the output. However, those  
 $\lambda$ -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 used in the compiler to  
 (incompletely) eliminate the intermediate closures [Breitner 2014].

- Object types do not depend on object terms.

Hence, we expect that during staging, every  $f : \mathbb{A} \rightarrow \text{ValTy}$  function has to be constant. This is actually true in the staging semantics of CFTT. There, all metafunctions are stable under object-level parallel substitutions. Also, object types are untouched by substitution. Hence, a straightforward unfolding of semantic definitions validates the constancy of  $f$ .

Generally speaking, in the semantics, every function whose domain is a product of object types and whose codomain is a constant presheaf, is a constant function. 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 staged compilation literature.

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

**Axiom (generativity).** For every  $f : \text{Elp } A \rightarrow \text{Usop}$  and  $a, a' : \text{Elp } A$ , we have that  $f a = f a'$ .

We remark that there is no risk of staging getting stuck on this axiom, because propositional equality proofs get erased, as we described in Section 2.4.

From generativity, we derive  $\Sigma_{\text{SOP}} : (\text{A} : \text{Usop}) \rightarrow (\text{B} : \text{El}_{\text{SOP}} \text{A} \rightarrow \text{Usop}) \rightarrow \text{Usop}$  as follows. First, for each  $\text{A} : \text{Up}$ , we define  $\text{loop}_A : \text{Elp } A$  as a product of non-terminating object programs. This is only needed to get arbitrary inhabitants with which we can call  $\text{Elp } A \rightarrow \text{Usop}$  functions.

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

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

$$\begin{aligned} \text{Rep } \{\Sigma \text{A B}\} &= \Sigma (\text{Rep } \{\text{A}\}) (\lambda x. \text{Rep } \{\text{B}(\text{rep}^{-1} x)\}) \\ \text{rep } \{\Sigma \text{A B}\} (a, b) &= (\text{rep } \{\text{A}\} a, \text{rep } \{\text{B}(\text{rep}^{-1} (\text{rep } x))\} b) \end{aligned}$$

Note that  $\text{rep}$  is only well-typed up to the fact that  $\text{rep}^{-1} \circ \text{rep} = \text{id}$ ; the Agda definition contains an additional transport that we omit here. This is, in fact, our reason for including the isomorphism equations in  $\text{IsSOP}$ . This, in turn, necessitates using  $\text{SOP}$  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  $\Sigma$ -type. When we use  $\text{SOP}$  instead, we do not have to use object-level products and this issue does not appear.

We skip the full definition of  $\Sigma_{\text{SOP}}$  here. The reader may refer to the  $\text{SOP}$  module in the Agda implementation, which is altogether 260 lines with all  $\text{IsSOP}$  instances.

#### 4.4 Let-Insertion & Case Splitting in Monadic Style

While  $\text{Pull}$  is not a monad, and hence also not a  $\text{MonadGen}$ , 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
```

#### 4.3 concatMap for Streams

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

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

the  $\text{A} \rightarrow \text{Pull B}$  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:

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

The idea is the following: if  $\text{Usop}$  is closed under  $\Sigma$ -types, we can directly define this  $\text{concatMap}$ , by taking the appropriate dependent sum of the  $\text{A} \rightarrow \text{Usop}$  family of machine states, which we extract from the  $\text{A} \rightarrow \text{Pull B}$  function. Let us write  $\Sigma \text{A B} : \text{MetaTy}$  for dependent sums, for  $\text{A} : \text{MetaTy}$  and  $\text{B} : \text{A} \rightarrow \text{MetaTy}$ , 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  $\text{projStep} : (\text{as} : \text{Pull A}) \rightarrow \text{projS as} \rightarrow \text{Gen}(\text{Step}(\text{projS as}) \text{A})$ . For now, we assume that the following instance exists:

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

Above,  $\{a : \text{A}\} \rightarrow \text{IsSOP } (\text{B } a)$  is a universally quantified instance constraint; this can be also represented in Agda. The definition of  $\text{concatMap}$  is as follows.

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

$\text{concatMap } \{A\} \{B\} f (\text{Pull S seed step}) =$

$\text{Pull } (S, \text{Maybe } (\Sigma A (\text{projS } \circ f))) ((, ) \triangleleft \text{seed} \triangleleft \text{return Nothing}) \$ \lambda \text{case}$

$(s, \text{Nothing}) \rightarrow \text{step s} \triangleleft \lambda \text{case}$

$\text{Stop} \rightarrow \text{return Stop}$

$\text{Skip s} \rightarrow \text{return } \$ \text{Skip } (s, \text{Nothing})$

$\text{Yield as} \rightarrow \text{do } s' \leftarrow \text{projSeed } (f a); \text{return } \$ \text{Skip } (s, \text{Just } (a, s'))$

$(s, \text{Just } (a, s')) \rightarrow \text{projStep } (f a) s' \triangleleft \lambda \text{case}$

$\text{Stop} \rightarrow \text{return } \$ \text{Skip } (s, \text{Nothing})$

$\text{Skip s'} \rightarrow \text{return } \$ \text{Skip } (s, \text{Just } (a, s'))$

$\text{Yield b s'} \rightarrow \text{return } \$ \text{Yield b } (s, \text{Just } (a, s'))$

Here,  $\text{Nothing}$  marks the states where we are in the “outer” loop, running the  $\text{Pull A}$  stream until we get its next value.  $\text{Just}$  marks the states of the “inner” loop, where we have a concrete  $a : \text{A}$  value and we run the  $(f a)$  stream until it stops. In the inner loop, the machine state type depends on the  $a : \text{A}$  value, hence the need for  $\Sigma$ . How do we get  $\text{IsSOP}$  for  $\Sigma$ ? 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  $f : \mathbb{A} \rightarrow \text{ValTy}$  function has to be constant. This is actually true in the staging semantics of CFTT. There, all metafunctions are stable under object-level parallel substitutions. Also, object types are untouched by substitution. Hence, a straightforward unfolding of semantic definitions validates the constancy of  $f$ .

Now, we should use this operation with care, since it has two states and can contribute to a code size blow-up. 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 a single IsSOP value.

```
bindsingle : IsSOP 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, IsSOP (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, gen<sub>Pull</sub> (~x \* 2) evaluates the given expression and saves the result to the machine state. Then, case<sub>Pull</sub> 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 case<sub>Pull</sub>, then finally we get an extra shape from forEach which introduces an additional Maybe.

#### 4.5 Discussion

Our stream library has a 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 a 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” [Wadler 1989].

Generally speaking, in the semantics, every function whose domain is a product of object types and whose codomain is a constant presheaf, is a constant function. 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 staged compilation literature.

Let us write Up = List ValTy and El<sub>P</sub> : Up → MetaTy for a universe of finite products of value types. Concretely, in CFTT and the Agda implementation, we assume the following:

**Axiom (generativity).** For every f : El<sub>P</sub> A → Usop and a, a' : El<sub>P</sub> A, we have that f a = f a'.

We remark that there is no risk of staging getting stuck on this axiom, because propositional equality proofs get erased, as we described in Section 2.4.

From generativity, we derive Σ<sub>SOP</sub> : (A : Usop) → (B : El<sub>SOP</sub> A → Usop) → Usop as follows. First, for each A : Up, we define loop<sub>A</sub> : El<sub>P</sub> A as a product of non-terminating object programs. This is only needed to get arbitrary inhabitants with which we can call El<sub>P</sub> A → Usop functions.

Then, Σ<sub>SOP</sub> A B is defined as the concatenation of A × B (inject; loop<sub>A</sub>) for each A<sub>i</sub> ∈ A, where (×) is the product type former in Usop and inject; has type El<sub>P</sub> A<sub>i</sub> → El<sub>SOP</sub> A. This is similar to the definition of non-dependent products in Usop, except that we have to get rid of the type dependency by instantiating B with arbitrary programs.

Then, we can show using the generativity axiom that Σ<sub>SOP</sub> supports projections, pairing and the βη-rules. These are all needed when we define the IsSOP instance, when we have to prove that encoding via rep is an isomorphism. Concretely, assuming IsSOP A and {a : A} → IsSOP (B a), we define Rep and rep for Σ A B as follows:

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

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

Note that rep is only well-typed up to the fact that  $\text{rep}^{-1} \circ \text{rep} = \text{id}$ ; the Agda definition contains an additional transport that we omit here. This is, in fact, our reason for including the isomorphism equations in IsSOP.

This, in turn, necessitates using SOP instead of finite sums, for the following reason. We can define a product type former for finite sums of value types, by taking the pairwise products of components, but in this case the products would be object-level products, which do not support βη rules. This implies that we cannot prove βη-rules for the derived product type in finite sums. In contrast, when we define products for SOP-s, we take meta-level pairwise products of components, which does support β and η.

We skip the full definition of Σ<sub>SOP</sub> here. The reader may refer to the SOP module in the Agda implementation, which is altogether 260 lines with all IsSOP 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, we should use this operation with care, since it has two states and can contribute to a code size blow-up. 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 an embedding

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 [Agda developers 2024] 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,<sup>1</sup> 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

*Two-level type theory.* Two-level lambda calculi were first developed in the context of abstract interpretation and partial evaluation [Nielson 1984; Nielson and Nielson 1992]. This line of research is characterized by simple types, having the same language features at different stages, and an emphasis on *binding time analysis*, i.e. automatically inferring stage annotations as part of a pipeline for partial evaluation and program optimization.

Later and independently, the same notion of level appeared in homotopy type theory, first in Voevodsky’s Homotopy Type System [Voevodsky 2013], and subsequently developed as 2LTT in [Annenkov et al. 2019]. Here, dependent types are essential, and we have different theories at the two stages. The application of 2LTT to staged compilation was developed in [Kovács 2022]. Binding-time analysis has not been developed in this setting; 2LTT-s have been meant to be used as surface languages to directly work in.

We build on both traditions. CFTT itself is a 2LTT and makes heavy use of dependent type theory. However, we are mostly concerned with code generation and optimization and borrow concepts from the partial evaluation literature.

*Tracking function arities and closures in types.* Downen et al.’s intermediate language  $\mathcal{IL}$  has similar motivations as our object language [Downen et al. 2020]. In both systems, function types are distinguished from closure types, enjoy universal  $\eta$ -conversion and have an explicit calling arity. In fact, our object language can be almost viewed as a small simply-typed fragment of  $\mathcal{IL}$ , with only letrec missing from  $\mathcal{IL}$ .

<sup>1</sup>Recall that we do not distinguish value and computation types in Haskell.

of Gen operations in streams.

```
bindGen : IsSOP A ⇒ Gen A → (A → Pull B) → Pull B
bindGen ga f =
  Pull (Σ A (projS ∘ f)) (do {a ← ga; s ← projSeed (f a); return (a, s)}) $ λ case
    (a, s) → projStep (f 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 = bindGen (gen a)

casePull : (Split A, IsSOP (SplitTo {A})) ⇒ ↑A → (SplitTo {A} → Pull B) → Pull B
casePull a = bindGen (split a)
```

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 finally we get an extra shape from `forEach` which introduces an additional `Maybe`.

#### 4.5 Discussion

Our stream library has a 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 a 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” [Wadler 1989].

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.

1128 *Sums-of-products.* SOP was proposed for generic Haskell programming by De Vries and Löh  
1129 [de Vries and Löh 2014]. Our own SOP implementation in Haskell is mostly the same as in *ibid*.  
1130 The SOP-of-object-code representation appeared as well in typed Template Haskell in [Pickering  
1131 et al. 2020].

1132 *CPS and binding-time improvement.* Our Gen monad can be viewed as a variation on the continuation-  
1133 passing-style that has been often used in staged programming. Flanagan et al.’s ANF translation  
1134 algorithm uses let-insertion in CPS that works the same way as our gen [Flanagan et al. 1993].  
1135 Kiselyov et al. likewise uses essentially the same CPS, in MetaOCaml [Kiselyov et al. 2017]. Jones et  
1136 al. discuss CPS and binding-time improvement in partial evaluation [Jones et al. 1993].  $\eta$ -expansion  
1137 for object-level finite sums is known as “the trick” [Danvy et al. 1996].

1138 *Join points.* The use join points in GHC’s core language is partly motivated by avoiding code  
1139 duplication in case-of-case transformations [Maurer et al. 2017]. In our monad library, case-of-case  
1140 is implicitly and eagerly computed during staging, and we similarly use join points to avoid code  
1141 duplication.

1142 *Stream fusion.* The staged stream fusion library *strymonas* by Kiselyov et al. is the primary prior art  
1143 [Kiselyov et al. 2017]. Here, streams are represented as a sum type of push and pull representations,  
1144 allowing both zipping and concatMap-ing. However, fusion is not guaranteed for all combinations;  
1145 zipping two concatMap-s reifies one of the streams in a runtime closure. In newer versions of the  
1146 library, fusion is reported to be complete and guaranteed [Kobayashi and Kiselyov 2024], however,  
1147 an exposition of this has not yet been published. Another notable difference is that *strymonas* relies  
1148 on mutable references in the object language, while we do not.

1149 *Machine fusion* [Robinson and Lippmeier 2017] supports splitting streams to multiple streams  
1150 while avoiding duplicated computation; this is not possible in *strymonas* or our library. Machine  
1151 fusion also supports concatMap and zipping. However, its state machine representation is sig-  
1152 nificantly more complicated than ours, and its fusion algorithm is not guaranteed to succeed on  
1153 arbitrary stream programs.

1154 Coutts developed pull stream fusion in depth in [Coutts 2011]. We borrowed the basic design and  
1155 the basic stream combinator definitions from there. This account is also close to existing stream  
1156 implementations in Haskell. *Ibid.* characterizes fusibility in a non-staged setting, in terms of “good  
1157 consumers” and “good producers”. This scheme does not cover concatMap.

## 6 CONCLUSIONS AND FUTURE WORK

We believe that a programming language in the style of this paper would be highly useful, especially in high-performance functional programming and domain-specific programming. In particular, the pipeline for using and compiling monads could look like the following.

- Users write definitions in a style similar to Haskell, without quotes, splices, up-s and down-s, only marking stages in type annotations and in let-definitions, with := and =. Storing monadic actions at runtime requires an explicit boxing operation that yields a closure type.
- Type- and stage-directed elaboration adds the missing operations, and also desugars case splitting using join and the underlying splitting function.
- In the downstream compiler, fast & conservative optimization passes are sufficient, since abstraction overheads are already eliminated by staging. Eliminating dead code and unused arguments would be important, since we have not yet addressed this through staging.
- A systematic way to de-duplicate code would be also needed, probably already during staging.

## 4.6 Agda & Haskell Implementations

In Agda, to avoid computation getting stuck on the generativity axiom, we use the `primTrustMe` built-in [Agda developers 2024] 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,<sup>1</sup> 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 unstaging time. If library users do not violate generativity, these coercions disappear and the code output will not contain unsafe coercions.

## 5 RELATED WORK

*Two-level calculi.* Two-level lambda calculi were first developed in the context of abstract interpretation and partial evaluation [Nielson 1984; Nielson and Nielson 1992]. This line of research is characterized by simple types, having the same language features at different stages, and an emphasis on *binding time analysis*, i.e. automatically inferring stage annotations as part of a pipeline for partial evaluation and program optimization.

Later and independently, the same notion of level appeared in homotopy type theory, first in Voevodsky’s Homotopy Type System [Voevodsky 2013], and subsequently developed as 2LTT in [Annenkov et al. 2019]. Here, dependent types are essential, and we have different theories at the two stages. The application of 2LTT to staged compilation was developed in [Kovács 2022]. Binding-time analysis has not been developed in this setting; 2LTT-s have been meant to be used as surface languages to directly work in.

*Higher-order abstract syntax and logical frameworks.* Hofmann’s work on the semantics of higher-order abstract syntax (HOAS) is an important precursor to our work [Hofmann 1999]. It anticipates many of the later developments in logical frameworks and two-level type theories. In particular, we build on Hofmann’s presheaf model for our 2LTT semantics, and his sketch of adequacy for HOAS corresponds to our definition of unstaging and its soundness proof. In general, there is a close correspondence between logical frameworks that represent object languages in HOAS style (e.g. [Harper et al. 1993]) and two-level type theories, and in some cases the two kinds of presentations are inter-derivable [Kovács 2023, Section 3.3]. In fact, our Agda implementation of CFTT also uses a HOAS embedding of the object theory.

Cocon is a logical framework that has a dependently typed meta level (“computation” level) and an object language that is configured by a signature [Pientka et al. 2019]. Cocon, like 2LTTs, can be used as a dependently typed metaprogramming language, but there are several differences. First, Cocon supports and emphasizes intensional analysis (i.e. the ability to do structural induction on object code), which is made possible through a contextual modality [Nanevski et al. 2008]. In contrast, intensional analysis is missing from 2LTTs so far. However, in Cocon, we can only

<sup>1</sup>Recall that we do not distinguish value and computation types in Haskell.

1177 Going a bit further, we believe that a 2LTT-based language could be a good design point for  
 1178 programming in general, buying us plenty of control over code generation for a modest amount of  
 1179 extra complexity.

1180 Also, rebuilding known abstractions in staged programming is valuable because it provides a  
 1181 semantic explanation of how abstractions can be compiled, and provides insights that could be  
 1182 reused even in non-staged settings. For instance, in this paper we demonstrated that compiling  
 1183 monadic code can be done by using *the same monads* in the metalanguage, extended with code  
 1184 generation as an effect.

1185 Continuing this line of thought, it could be interesting to also adapt to 2LTT the style of  
 1186 *binding-time analysis* that is well-known in partial evaluation. For example, we might do program  
 1187 optimization in the following way. First, start with monadic code in a non-staged language. Second,  
 1188 try to infer stages, inventing quotes, splices, up-s and down-s, thereby translating definitions to a  
 1189 2LTT. Third, perform staging and proceed from there. This would be more fragile than unambiguous  
 1190 stage inference, but we would still benefit from shifting a lot of machinery into staging (which is  
 1191 deterministic and efficient).

1192 We also find it interesting how little impact closure-freedom had on the developments in this  
 1193 paper. This provides some evidence that in staged functional programming, closures can be an  
 1194 opt-in feature instead of the pervasive default.

1195 In this paper we only briefly looked at two applications. Both monad transformers and streams  
 1196 could be further developed, and many other programming abstractions could be revisited in the  
 1197 staged setting.

- We did not discuss the issue of *tail calls* in monadic code. For example, a tail call in `Maybe`  
 1199 should not case split on the result, while in our current library, making a call with `up` always  
 1200 does so. We did develop an abstraction for tail calls in the Agda and Haskell implementations,  
 1201 but omitted it here partly for lack of space, and partly because we have not yet explored  
 1202 much of the design space.
- In a practical stream library, it would be useful to further try to minimize the size of the  
 1204 machine state. In the Agda and Haskell versions, we additionally track whether streams can  
 1205 `Skip`, to provide more efficient zipping for non-skipping “synchronous” streams. However,  
 1206 this could be refined in many ways, and we could also try to represent the full transition  
 1207 graph in an observable way, thereby enabling merging or deleting states.
- It seems promising to look for synergies between push streams, pull streams and monad  
 1209 transformers. It seems that pull streams could be generalized from `Gen` to different monads,  
 1210 in the representation of transitions and initial states. This would allow putting a `Pull` on the  
 1211 top of a monad transformer stack. On the other hand, push streams form a proper monad,  
 1212 but they cannot be used to “transform” other monads, so they can be placed at the bottom  
 1213 of a transformer stack.
- We could try to lean more heavily on datatype-generic programming and generalize ab-  
 1215 stractions to arbitrary object types. For example, push and pull streams could be generalized  
 1216 to inductive and coinductive Church-codings of arbitrary value types.

## REFERENCES

- Agda developers. 2024. Agda documentation. <https://agda.readthedocs.io/en/v2.6.4.2/>
- Daniil Ananenkov, Paolo Capriotti, Nicolai Kraus, and Christian Sattler. 2019. Two-Level Type Theory and Applications. *ArXiv e-prints* (may 2019). <http://arxiv.org/abs/1705.03307>
- Joachim Breitner. 2014. Call Arity. In *Trends in Functional Programming - 15th International Symposium, TFP 2014, Soesterberg, The Netherlands, May 26-28, 2014. Revised Selected Papers (Lecture Notes in Computer Science, Vol. 8843)*. Jurriaan Hage and Jay McCarthy (Eds.). Springer, 34–50. [https://doi.org/10.1007/978-3-319-14675-1\\_3](https://doi.org/10.1007/978-3-319-14675-1_3)

manipulate object terms at the meta level by specifying their concrete object-level contexts (i.e.  
 possible free variables), while in 2LTTs the object-level contexts are implicit and are only computed  
 during code generation. The latter style is more convenient for staged compilation applications,  
 where explicit context tracking would add a significant amount of bureaucracy.

Nonetheless, it seems to be possible and useful to have a system which simultaneously supports  
 intensional analysis through contextual types and the 2LTT-style implicit contexts.

*CPS and monads in staged programming.* Writing code generators in continuation-passing style  
 was first proposed by Bondorf [Bondorf 1992]. Flanagan et al.’s ANF translation algorithm uses  
 CPS for let-insertion as well [Flanagan et al. 1993], similarly to our `gen` function. These sources do  
 no use explicit monadic notation though Jones et al. discuss CPS and binding-time improvement  
 in partial evaluation in [Jones et al. 1993]. In [Kiselyov et al. 2004], [Swadi et al. 2006] and [Carette  
 and Kiselyov 2011], the composition of state and continuation monads was used, similarly to our  
`StateT S Gen`, but without polymorphic answer types. The definition of let-insertion here is again  
 similar to ours.

*Tracking function arities and closures in types.* Downen et al.’s intermediate language  $\mathcal{IL}$  has similar  
 motivations as our object language [Downen et al. 2020]. In both systems, function types are  
 distinguished from closure types, enjoy universal  $\eta$ -conversion and have an explicit calling arity.  
 In fact, our object language can be almost viewed as a small simply-typed fragment of  $\mathcal{IL}$ , with  
 only `letrec` missing from  $\mathcal{IL}$ .

*Sums-of-products.* SOP was proposed for generic Haskell programming by De Vries and Löh  
 [de Vries and Löh 2014b]. Our own SOP implementation in Haskell is mostly the same as in ibid.  
 The SOP-of-object-code representation appeared as well in typed Template Haskell in [Pickering  
 et al. 2020].  $\eta$ -expansion for object-level finite sums is known as “the trick” [Danvy et al. 1996].

*join points.* The use join points in GHC’s core language is partly motivated by avoiding code  
 duplication in case-of-case transformations [Maurer et al. 2017]. In our monad library, case-of-case  
 is implicitly and eagerly computed during staging, and we similarly use join points to avoid code  
 duplication.

*Stream fusion.* The staged stream fusion library `strymonas` by Kiselyov et al. is the primary prior art  
 [Kiselyov et al. 2017]. Here, streams are represented as a sum type of push and pull representations,  
 allowing both zipping and concatMap-ing. However, fusion is not guaranteed for all combinations;  
 zipping two concatMap-s reifies one of the streams in a runtime closure. In newer versions of  
 the library, fusion is complete [Kobayashi and Kiselyov 2024], however, an exposition of this has  
 not yet been published. Additionally, `strymonas` heavily relies on mutable references in the object  
 language, and also uses some switching on object-level data to implement control flow. This style of  
 code output can be still reliably optimized by downstream compilers (C or OCaml), but our solution  
 is more conservative in that it does not use any mutation or runtime switching. This “purity” of our  
 streams also makes them easier to generalize, e.g. by using arbitrary `MonadGen` monads instead of  
`Gen` in stream internals.

*Machine fusion* [Robinson and Lippmeier 2017] supports splitting streams to multiple streams  
 while avoiding duplicated computation; this is not possible in `strymonas` or our library. Machine  
 fusion also supports concatMap and zipping. However, its state machine representation is signifi-  
 cantly more complicated than ours, and its fusion algorithm is not guaranteed to succeed on  
 arbitrary stream programs.

Coutts developed pull stream fusion in depth in [Coutts 2011]. We borrowed the basic design and  
 the basic stream combinator definitions from there. This account is also close to existing stream

- 1226 Duncan Coutts. 2011. *Stream fusion : practical shortcut fusion for conductive sequence types*. Ph.D. Dissertation. University  
1227 of Oxford, UK. <http://ora.ox.ac.uk/objects/uuid:b4971f57-2b94-4fdf-a5e0-98d6935a44da>
- 1228 Olivier Danvy, Karoline Malmkjær, and Jens Palsberg. 1996. Eta-Expansion Does The Trick. *ACM Trans. Program. Lang.  
Syst.* 18, 6 (1996), 730–751. <https://doi.org/10.1145/236114.236119>
- 1229 Edsko de Vries and Andres Löh. 2014. True sums of products. In *Proceedings of the 10th ACM SIGPLAN workshop on Generic  
programming, WGP 2014, Gothenburg, Sweden, August 31, 2014*, José Pedro Magalhães and Tiark Roßmann (Eds.). ACM,  
83–94. <https://doi.org/10.1145/2633628.2633634>
- 1230 Paul Downen, Zena M. Ariola, Simon Peyton Jones, and Richard A. Eisenberg. 2020. Kinds are calling conventions. *Proc.  
ACM Program. Lang.* 4, ICFP (2020), 104:1–104:29. <https://doi.org/10.1145/3408986>
- 1231 Peter Dybjer. 1994. Inductive Families. *Formal Aspects Comput.* 6, 4 (1994), 449–465. <https://doi.org/10.1007/BF01211308>
- 1232 Richard A. Eisenberg and Stephanie Weirich. 2012. Dependently typed programming with singletons. In *Proceedings of the  
5th ACM SIGPLAN Symposium on Haskell, Haskell 2012, Copenhagen, Denmark, 13 September 2012*, Janis Voigtländer (Ed.).  
ACM, 117–130. <https://doi.org/10.1145/2364506.2364522>
- 1233 Cormac Flanagan, Amr Sabry, Bruce F. Duba, and Matthias Felleisen. 1993. The Essence of Compiling with Continuations. In  
*Proceedings of the ACM SIGPLAN'93 Conference on Programming Language Design and Implementation (PLDI), Albuquerque,  
New Mexico, USA, June 23–25, 1993*, Robert Cartwright (Ed.). ACM, 237–247. <https://doi.org/10.1145/155090.155113>
- 1234 GHC developers. 2024a. GHC documentation. [https://downloads.haskell.org/ghc/9.8.2/docs/users\\_guide/](https://downloads.haskell.org/ghc/9.8.2/docs/users_guide/)
- 1235 GHC developers. 2024b. GHC Base source. <https://hackage.haskell.org/package/base-4.19.1.0/docs/src/GHC.Base.html>
- 1236 Ralf Hinze, Thomas Harper, and Daniel W. H. James. 2010. Theory and Practice of Fusion. In *Implementation and Application  
of Functional Languages - 22nd International Symposium, IFL 2010, Alphen aan den Rijn, The Netherlands, September 1–3,  
2010, Revised Selected Papers (Lecture Notes in Computer Science, Vol. 6647)*, Jurriaan Hage and Marco T. Morazán (Eds.).  
Springer, 19–37. [https://doi.org/10.1007/978-3-642-24276-2\\_2](https://doi.org/10.1007/978-3-642-24276-2_2)
- 1237 Jasper Hugunin. 2020. Why Not W?. In *26th International Conference on Types for Proofs and Programs, TYPES 2020, March  
2–5, 2020, University of Twente, Italy (LIPIcs, Vol. 188)*, Ugo de'Liguoro, Stefano Berardi, and Thorsten Altenkirch (Eds.).  
Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 8:1–8:9. <https://doi.org/10.4230/LIPICS.TYPES.2020.8>
- 1238 Neil D. Jones, Carsten K. Gomard, and Peter Sestoft. 1999. *Partial evaluation and automatic program generation*. Prentice  
Hall.
- 1239 Oleg Kiselyov, Aggelos Biboudis, Nick Palladinos, and Yannis Smaragdakis. 2017. Stream fusion, to completeness. In  
*Proceedings of the 44th ACM SIGPLAN Symposium on Principles of Programming Languages, POPL 2017, Paris, France,  
January 18–20, 2017*, Giuseppe Castagna and Andrew D. Gordon (Eds.). ACM, 285–299. <https://doi.org/10.1145/3009837.3009880>
- 1240 Tomonori Kobayashi and Oleg Kiselyov. 2024. Complete Stream Fusion for Software-Defined Radio. In *Proceedings of the  
2024 ACM SIGPLAN International Workshop on Partial Evaluation and Program Manipulation, PEPM 2024, London, UK, 16  
January 2024*, Gabriele Keller and Meng Wang (Eds.). ACM, 57–69. <https://doi.org/10.1145/3635800.3636962>
- 1241 András Kovács. 2022. Staged compilation with two-level type theory. *Proc. ACM Program. Lang.* 6, ICFP (2022), 540–569.  
<https://doi.org/10.1145/3547641>
- 1242 András Kovács. 2022. Demo implementation for the paper "Staged Compilation With Two-Level Type Theory". <https://doi.org/10.5281/zenodo.6757737>
- 1243 Xavier Leroy, Damien Doligez, Alain Frisch, Jacques Garrigue, Didier Rémy, Kc Sivaramakrishnan, and Jérôme Vouillon.  
2023. The OCaml system release 5.1: Documentation and user's manual. <https://v2.ocaml.org/manual/>
- 1244 Paul Blain Levy. 1999. Call-by-Push-Value: A Subsuming Paradigm. In *Typed Lambda Calculi and Applications, 4th Interna-tional Conference, TLCA 99, L'Aquila, Italy, April 7–9, 1999, Proceedings (Lecture Notes in Computer Science, Vol. 1581)*,  
Jean-Yves Girard (Ed.). Springer, 228–242. [https://doi.org/10.1007/3-540-48959-2\\_17](https://doi.org/10.1007/3-540-48959-2_17)
- 1245 Luke Maurer, Paul Downen, Zena M. Ariola, and Simon L. Peyton Jones. 2017. Compiling without continuations. In  
*Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2017,  
Barcelona, Spain, June 18–23, 2017*, Albert Cohen and Martin T. Vechev (Eds.). ACM, 482–494. <https://doi.org/10.1145/3062341.3062380>
- 1246 Ieke Moerdijk and Erik Palmgren. 2000. Wellfounded trees in categories. *Ann. Pure Appl. Log.* 104, 1–3 (2000), 189–218.  
[https://doi.org/10.1016/S0168-0072\(00\)00012-9](https://doi.org/10.1016/S0168-0072(00)00012-9)
- 1247 ntl developers. 2024. ntl documentation. <https://hackage.haskell.org/package/ntl>
- 1248 Flemming Nielson. 1984. *Abstract interpretation using domain theory*. Ph.D. Dissertation. University of Edinburgh, UK.  
<https://ethos.bl.uk/OrderDetails.do?uin=uk.bl.ethos.350060>
- 1249 Flemming Nielson and Hanne Riis Nielson. 1992. *Two-level functional languages*. Cambridge tracts in theoretical computer  
science, Vol. 34. Cambridge University Press.
- 1250 Matthew Pickering, Andres Löh, and Nicolas Wu. 2020. Staged sums of products. In *Proceedings of the 13th ACM SIGPLAN  
International Symposium on Haskell, Haskell@ICFP 2020, Virtual Event, USA, August 7, 2020*, Tom Schrijvers (Ed.). ACM,  
122–135. <https://doi.org/10.1145/3406088.3409021>

implementations in Haskell. We find it interesting that our staged definitions are very close to the non-staged versions there; essentially all of the additional complexity is compartmentalized in our SOP module. Coutts characterized fusibility in a non-staged setting, in terms of “good consumers” and “good producers”, but this scheme does not cover concatMap.

## 6 CONCLUSIONS AND FUTURE WORK

We believe that a programming language in the style of this paper would be highly useful, especially in high-performance functional programming and domain-specific programming. In particular, the pipeline for using and compiling monads could look like the following.

- Users write definitions in a style similar to Haskell, without quotes, splices, up-s and down-s, only marking stages in type annotations and in let-definitions, with := and =. Storing monadic actions at runtime requires an explicit boxing operation that yields a closure type.
- Type- and stage-directed elaboration adds the missing operations, and also desugars case splitting using join and the underlying splitting function.
- In the downstream compiler, fast & conservative optimization passes are sufficient, since abstraction overheads are already eliminated by staging. Eliminating dead code and unused arguments would be important, since we have not yet addressed this through staging.
- A systematic way to de-duplicate code would be also needed, probably already during staging.

Going a bit further, we believe that a 2LTT-based language could be a good design point for programming in general, buying us plenty of control over code generation for a modest amount of extra complexity.

Also, rebuilding known abstractions in staged programming is valuable because it provides a semantic explanation of how abstractions can be compiled, and provides insights that could be reused even in non-staged settings. For instance, in this paper we demonstrated that compiling monadic code can be done by using the same monads in the metalanguage, extended with code generation as an effect.

Continuing this line of thought, it could be interesting to also adapt to 2LTT the style of binding-time analysis that is well-known in partial evaluation. For example, we might do program optimization in the following way. First, start with monadic code in a non-staged language. Second, try to infer stages, inventing quotes, splices, up-s and down-s, thereby translating definitions to a 2LTT. Third, perform staging and proceed from there. This would be more fragile than unambiguous stage inference, but we would still benefit from shifting a lot of machinery into staging (which is deterministic and efficient).

We also find it interesting how little impact closure-freedom had on the developments in this paper. This provides some evidence that in staged functional programming, closures can be an opt-in feature instead of the pervasive default.

In this paper we only briefly looked at two applications. Both monad transformers and streams could be further developed, and many other programming abstractions could be revisited in the staged setting.

- We did not discuss the issue of tail calls in monadic code. For example, a tail call in Maybe should not case split on the result, while in our current library, making a call with up always does so. We did develop an abstraction for tail calls in the Agda and Haskell implementations, but omitted it here partly for lack of space, and partly because we have not yet explored much of the design space.
- In a practical stream library, it would be useful to further try to minimize the size of the machine state. In the Agda and Haskell versions, we additionally track whether streams can skip, to provide more efficient zipping for non-skipping “synchronous” streams. However,

- 1275 Amos Robinson and Ben Lippmeier. 2017. Machine fusion: merging merges, more or less. In *Proceedings of the 19th*  
 1276 *International Symposium on Principles and Practice of Declarative Programming, Nmara, Belgium, October 09 - 12, 2017*,  
 1277 Wim Vanhoof and Brigitte Pientka (Eds.). ACM, 139–150. <https://doi.org/10.1145/3131851.3131865>
- 1278 Walid Taha and Tim Sheard. 2000. MetaML and multi-stage programming with explicit annotations. *Theor. Comput. Sci.* 248,  
 1279 1-2 (2000), 211–242. [https://doi.org/10.1016/S0304-3975\(00\)00053-0](https://doi.org/10.1016/S0304-3975(00)00053-0)
- 1280 Vladimir Voevodsky. 2013. A simple type system with two identity types. (2013). Unpublished note.
- 1281 Philip Wadler. 1989. Theorems for free!. In *Functional Programming Languages and Computer Architecture*. ACM Press,  
 1282 347–359.
- 1283 Philip Wadler and Stephen Blott. 1989. How to Make ad-hoc Polymorphism Less ad-hoc. In *Conference Record of the Sixteenth*  
 1284 *Annual ACM Symposium on Principles of Programming Languages, Austin, Texas, USA, January 11-13, 1989*. ACM Press,  
 1285 69–76. <https://doi.org/10.1145/75277.75283>
- 1286 Ningning Xie, Matthew Pickering, Andres Löh, Nicolas Wu, Jeremy Yallop, and Meng Wang. 2022. Staging with class: a  
 1287 specification for typed template Haskell. *Proc. ACM Program. Lang.* 6, POPL (2022), 1–30. <https://doi.org/10.1145/3498723>
- 1288 Jeremy Yallop and Oleg Kiselyov. 2019. Generating mutually recursive definitions. In *Proceedings of the 2019 ACM SIGPLAN*  
 1289 *Workshop on Partial Evaluation and Program Manipulation, PEPM@POPL 2019, Cascais, Portugal, January 14-15, 2019*,  
 1290 Manuel V. Hermenegildo and Atsushi Igarashi (Eds.). ACM, 75–81. <https://doi.org/10.1145/3294032.3294078>
- 1291  
 1292  
 1293  
 1294  
 1295  
 1296  
 1297  
 1298  
 1299  
 1300  
 1301  
 1302  
 1303  
 1304  
 1305  
 1306  
 1307  
 1308  
 1309  
 1310  
 1311  
 1312  
 1313  
 1314  
 1315  
 1316  
 1317  
 1318  
 1319  
 1320  
 1321  
 1322  
 1323

this could be refined in many ways, and we could also try to represent the full transition graph in an observable way, thereby enabling merging or deleting states.

- It seems promising to look for synergies between push streams, pull streams and monad transformers. It seems that pull streams could be generalized from Gen to different monads, in the representation of transitions and initial states. This would allow putting a Pull *on the top* of a monad transformer stack. On the other hand, push streams form a proper monad, but they cannot be used to “transform” other monads, so they can be placed at the *bottom* of a transformer stack.
- We could try to lean more heavily on datatype-generic programming and generalize abstractions for more object types. For example, push and pull streams could be generalized to inductive and coinductive Church-codings of arbitrary value types.