



# From Capabilities to Regions: Enabling Efficient Compilation of Lexical Effect Handlers

MARIUS MÜLLER, University of Tübingen, Germany

PHILIPP SCHUSTER, University of Tübingen, Germany

JONATHAN LINDEGAARD STARUP, Aarhus University, Denmark

KLAUS OSTERMANN, University of Tübingen, Germany

JONATHAN IMMANUEL BRACHTHÄUSER, University of Tübingen, Germany

Effect handlers are a high-level abstraction that enables programmers to use effects in a structured way. They have gained a lot of popularity within academia and subsequently also in industry. However, the abstraction often comes with a significant runtime cost and there has been intensive research recently on how to reduce this price.

A promising approach in this regard is to implement effect handlers using a CPS translation and to provide sufficient information about the nesting of handlers. With this information the CPS translation can decide how effects have to be lifted through handlers, *i.e.*, which handlers need to be skipped, in order to handle the effect at the correct place. A structured way to make this information available is to use a calculus with a region system and explicit subregion evidence. Such calculi, however, are quite verbose, which makes them impractical to use as a source-level language.

We present a method to infer the lifting information for a calculus underlying a source-level language. This calculus uses second-class capabilities for the safe use of effects. To do so, we define a typed translation to a calculus with regions and evidence and we show that this lift-inference translation is typability- and semantics-preserving. On the one hand, this exposes the precise relation between the second-class property and the structure given by regions. On the other hand, it closes a gap in a compiler pipeline enabling efficient compilation of the source-level language. We have implemented lift inference in this compiler pipeline and conducted benchmarks which indicate that the approach is indeed working.

CCS Concepts: • **Software and its engineering** → **Control structures; Compilers**; • **Theory of computation** → **Type theory; Control primitives**.

Additional Key Words and Phrases: effect handlers, region inference, lift inference

## ACM Reference Format:

Marius Müller, Philipp Schuster, Jonathan Lindegaard Starup, Klaus Ostermann, and Jonathan Immanuel Brachthäuser. 2023. From Capabilities to Regions: Enabling Efficient Compilation of Lexical Effect Handlers. *Proc. ACM Program. Lang.* 7, OOPSLA2, Article 255 (October 2023), 30 pages. <https://doi.org/10.1145/3622831>

## 1 INTRODUCTION

Languages with effect handlers [Plotkin and Pretnar 2009, 2013] offer a high-level way to structure effectful programs. Effect handlers allow for a combination of various effects by giving meaning to abstract effect operations (such as exceptions, async-await, generators, logic programming, or

---

Authors' addresses: Marius Müller, University of Tübingen, Germany, mari.mueller@uni-tuebingen.de; Philipp Schuster, University of Tübingen, Germany, philipp.schuster@uni-tuebingen.de; Jonathan Lindegaard Starup, Aarhus University, Denmark, jls@cs.au.dk; Klaus Ostermann, University of Tübingen, Germany, klaus.ostermann@uni-tuebingen.de; Jonathan Immanuel Brachthäuser, University of Tübingen, Germany, jonathan.brachthaeuser@uni-tuebingen.de.



This work is licensed under a Creative Commons Attribution 4.0 International License.

© 2023 Copyright held by the owner/author(s).

2475-1421/2023/10-ART255

<https://doi.org/10.1145/3622831>

probabilistic programming) in a composable way. In the past decade, effect handlers have been a hot topic in programming language research and have also gained more and more popularity outside academia. They have been implemented not only in research languages (such as Eff [Bauer and Pretnar 2015], Koka [Leijen 2017], Frank [Lindley et al. 2017], Effekt [Brachthäuser et al. 2020], or Helium [Biernacki et al. 2019]), but also practical general purpose languages such as OCaml [Sivaramakrishnan et al. 2021], Scala<sup>2</sup>, Unison<sup>3</sup>, and WebAssembly<sup>1</sup> are following lead and have started to integrate effects and handlers. However, to make effect handlers practically useful, it is critically important to minimize the runtime cost incurred by this abstraction.

Our goal, thus, is to have a source-level language with effect handlers that can be compiled efficiently. Ideally, the language should be effect-safe and sufficiently expressive without being unnecessarily complex. A recently developed way to strike this balance is using lightweight effect polymorphism. Brachthäuser et al. [2020] show how to design such a system using *second-class capabilities* [Osvald et al. 2016]. They develop the language Effekt which features lexical effect handlers [Biernacki et al. 2019; Brachthäuser et al. 2020; Zhang and Myers 2019], a recent variant of effect handlers, and show how to translate it to System  $\Xi$ , a calculus in explicit capability-passing style. The translation preserves typing and is used to give the semantics for Effekt.

A recently developed way to efficiently implement lexical effect handlers uses iterated continuation-passing style (CPS). It has been shown to yield good performance results [Schuster et al. 2020]. Moreover, Schuster et al. [2022b] have designed a core calculus  $\Lambda_{\text{cap}}$ , which features effect handlers based on *regions*. They show how to translate  $\Lambda_{\text{cap}}$  to pure System F in a typability- and semantics-preserving way.

Figure 1 summarizes the developments mentioned above. To reach our goal and complete the pipeline, we have to close the gap in the middle: we have to show how to translate from a system with second-class capabilities to a system with region-based effects, *i.e.*, we have to understand the relation between these two concepts precisely.

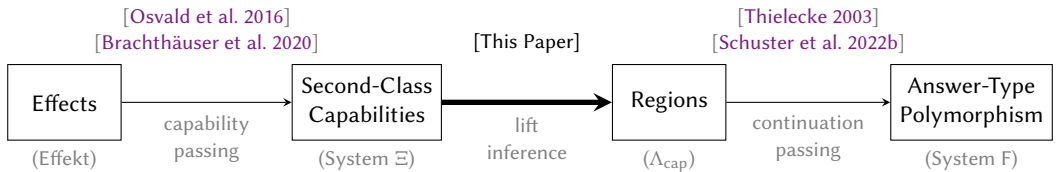


Fig. 1. Overview of this paper in relation to prior work. Nodes are labeled with mechanisms to ensure effect safety (e.g., Effects); below each node one example calculus is listed. Each arrow corresponds to a translation between calculi.

To do so, in this paper we present a typed translation from the calculus of second-class capabilities System  $\Xi$  to the region-based calculus  $\Lambda_{\text{cap}}$ . The typed nature of the translation makes the relation between the two concepts explicit. It connects the lexical scopes of the definition sites of capabilities in System  $\Xi$  with corresponding regions in  $\Lambda_{\text{cap}}$ , in which the capabilities are allowed to be used. This connection also materializes in the definition of sound translation environments (see Definition 3.4), which are maintained during the translation as a key component.

To practically evaluate our approach, we have implemented the translation as a compiler phase in the Effekt language. As our implementation strategy is to translate effects and handlers to CPS

<sup>2</sup><https://github.com/lampepfl/dotty/pull/16739>

<sup>3</sup><https://www.unison-lang.org/learn/fundamentals/abilities>

<sup>1</sup><https://github.com/effect-handlers/wasm-spec>

[Hillerström et al. 2020; Schuster et al. 2022b, 2020], we have also implemented the CPS translation of Schuster et al. [2022b] with Standard ML [Milner et al. 1997] as the target language, hence completing a compiler pipeline for Effekt. To further compile SML we use the MLton<sup>4</sup> compiler.

For a certain restricted class of programs Schuster et al. [2020] prove that all abstractions related to effects and handlers can theoretically be eliminated. Moreover, they demonstrate excellent performance on a number of benchmarks which have independently been reproduced by Karachalias et al. [2021]. However, for their performance evaluation, they wrote programs directly in a core language  $\lambda_{\text{cap}}$ . We, for the first time, can reproduce the performance claims of Schuster et al. [2020] in a realistic source-level language.

Comparing our implementation with other state-of-the-art implementations of effect handlers, we found that our relative performance ranges from 2.1x slowdown to 44.4x speedup compared with OCaml [Sivaramakrishnan et al. 2021], 1.1x slowdown to 87.8x speedup compared with Koka [Leijen 2017; Xie et al. 2020], and 1.2x slowdown to 23.3x speedup compared with Eff [Karachalias et al. 2021; Pretnar et al. 2017]. The results indicate that the compilation technique presented by Schuster et al. [2020] works for a high-level language presented by Brachthäuser et al. [2020] and yields good performance.

By closing the conceptual gap between lexical scoping and regions, we enable efficient compilation of lexical effect handlers, like those found in Effekt or Helium. However, our results do not immediately carry over to dynamic effect handlers, as implemented in OCaml 5, Koka, and WebAssembly and we leave further investigation of this to future work.

Our contributions are the following.

- We formally present a typability- and semantics-preserving translation from System  $\Xi$  to  $\Lambda_{\text{cap}}$ . We refer to this as *lift-inference translation* for reasons to be explained shortly.
- From a theoretical perspective, we hence clarify the precise relation between scope-based reasoning for second-class capabilities and region-based reasoning.
- From a practical perspective, we fill in an important missing link in the compiler pipeline illustrated in Figure 1.
- In  $\Lambda_{\text{cap}}$ , continuation calls have to be scoped [Xie et al. 2020]. No such restriction exists for System  $\Xi$ . To support System  $\Xi$ , we lift the scoped-continuation restriction imposed by Schuster et al. [2022b] and thus generalize  $\Lambda_{\text{cap}}$ . We prove that the generalized system still allows for a translation to iterated CPS satisfying the same properties as the original.
- We have implemented the lift-inference translation and the CPS translation as steps for the efficient compilation of the source-level language Effekt to SML. In addition, we have performed benchmarks indicating that our approach is competitive with or often faster than other state-of-the-art languages featuring effect handlers.

Next, in Section 2, we introduce the main ideas of the lift-inference translation by considering examples. In Section 3, we first recap the two calculi involved and then formally present the lift-inference translation. A discussion of the implementation and the corresponding benchmarks is given in Section 4. In Section 5, we compare to related work. We conclude and outline future work in Section 6.

Due to lack of space, proofs and some definitions have been omitted. These can be found in an extended technical report [Müller et al. 2023b].

## 2 MAIN IDEAS

In this section we introduce the main ideas for our lift-inference translation from the source calculus System  $\Xi$  [Brachthäuser et al. 2020] to the target calculus  $\Lambda_{\text{cap}}$  [Schuster et al. 2022b]. Both calculi

<sup>4</sup><http://mlton.org>

are mostly standard functional languages with multi-arity functions, but additionally feature lexical effect handlers [Biernacki et al. 2019; Brachthäuser et al. 2020; Zhang and Myers 2019]. Effect handlers are called *lexical* if effect operation calls are lexically related to the corresponding effect handler. This is in contrast to the more traditional dynamically scoped handlers, which search the stack for the closest handler at runtime. Both calculi establish this lexical connection between effect and handler by using explicit capability-passing style [Brachthäuser et al. 2020; Zhang and Myers 2019]. This means that effects are bound as capabilities by their handlers and can then be used in the scope of the handled statement.

Both languages are effect safe, that is, they guarantee that all effect operations are eventually handled. In System  $\Xi$ , this is achieved by making functions and capabilities *second-class* [Osvald et al. 2016] so that they cannot be returned or stored in data structures, hence making sure that capabilities cannot leave their defining scope. To make this explicit, these second-class functions and capabilities are called blocks. While blocks in System  $\Xi$  are required to be second-class, they still can be higher-order, *i.e.*, functions can abstract block parameters.

In  $\Lambda_{\text{cap}}$ , there is no such second-class restriction on functions and capabilities. To ensure effect safety,  $\Lambda_{\text{cap}}$  features a region system with subregioning and explicit subregion evidence instead. Here, a *region* denotes the scope of an effect handler and subregion evidence is used to constructively witness how handlers are nested. By enforcing that capabilities can only be called in a subregion of their corresponding handler, this system also ensures that they cannot leave their defining scope.

Since evidence terms precisely witness how handlers are nested, they contain the information of where capabilities have to be *lifted* to when they are called, *i.e.*, how many handlers have to be jumped over until the correct handler is found. This enables efficient compilation of effects and handlers [Schuster et al. 2022b, 2020]: the lifting can often be promoted to compile time which avoids the search for the correct handler at runtime. The goal of our lift-inference translation is thus to infer this lifting information by endowing terms in System  $\Xi$  with correct regions and evidence to obtain valid terms in  $\Lambda_{\text{cap}}$ .

## 2.1 Basic Example

To illustrate the need for lifting, as well as how to perform lift inference, consider the following example in Effekt:

```
effect Yield(i: Int): Int
try {
  def g(i: Int): Int / {} = { do Yield(i) };
  val x = g(1);
  try { g(x) } with Yield { j ⇒ 42 }
} with Yield { j ⇒ resume(j + 1) }
```

We define a local function  $g$ , which uses the `Yield` effect to return an integer. It is annotated to have type `Int` and no observable effects `{}`. This means the (dynamic) call site of  $g$  cannot handle the `Yield` effect and it needs to be handled at the (lexical) definition site of  $g$ . In consequence, running the example will result in the integer 3.

Subfigure 2 (a) shows the result of the type-and-effect directed translation [Brachthäuser et al. 2020] from Effekt (with support for effect inference) to System  $\Xi$  in explicit capability-passing style. Comparing the System  $\Xi$  term to the original program, we notice that handlers explicitly bind capabilities (*e.g.*, `yield1`) and effect calls now directly refer to a capability (*e.g.*, `do yield1(i)`). Capability passing in System  $\Xi$  also makes explicit that the effect call in the body of  $g$  refers to the outer handler.

<pre> <b>try</b> { (yield<sub>1</sub>) ⇒   <b>def</b> g(i : Int) { <b>do</b> yield<sub>1</sub>(i) };   <b>val</b> x = g(1);   <b>try</b> { (yield<sub>2</sub>) ⇒     g(x)   } <b>with</b> { (j, k) ⇒ 42 } } <b>with</b> { (j, k) ⇒ <b>do</b> k(j + 1) } </pre>	<pre> <b>try</b> { [r<sub>1</sub>; n<sub>1</sub> : r<sub>1</sub> ⊆ ⊤](yield<sub>1</sub>) ⇒   <b>def</b> g[r; n : r ⊆ r<sub>1</sub>](i : Int) <b>at</b> r { <b>do</b> yield<sub>1</sub>[n](i) };   <b>val</b> x = g[r<sub>1</sub>; ⊙](1);   <b>try</b> { [r<sub>2</sub>; n<sub>2</sub> : r<sub>2</sub> ⊆ r<sub>1</sub>](yield<sub>2</sub>) ⇒     g[r<sub>2</sub>; n<sub>2</sub>](x)   } <b>with</b> { (j, k) ⇒ 42 } } <b>with</b> { (j, k) ⇒ <b>do</b> k[⊙](j + 1) } </pre>
--	--

(a) Program in System  $\Xi$ . Effect safety is established by treating capabilities as second-class values.

(b) Program in  $\Lambda_{\text{cap}}$ . Effect safety is established by tracking the region of a capability and requiring subregion evidence.

Fig. 2. Simple example illustrating lexical effect handling.

At runtime, we need to make sure to handle the effect operation with the correct handler. Thus, when  $g$  is called the second time, we have to lift the capability in its body through the inner handler, that is, we need to skip over the inner handler to transfer control flow to the outer handler. Our goal is to make this skipping of handlers explicit. Then, our CPS translation can make use of this information to avoid searching for the correct handler at runtime, as it knows how many segments of the stack it has to capture.

One possibility, to make this information explicit, would be to use lifting annotations [Biernacki et al. 2017; Schuster et al. 2020]. The definition of  $g$  (for the second call) would then become

$$\mathbf{def} \ g(i : \text{Int}) \{ \mathbf{do} \ (\mathbf{lift} \ \text{yield}_1)(i) \};$$

But this is only correct when  $g$  is called under the inner handler. When it is called the first time, immediately after its definition, this annotation is incorrect since no handler has to be skipped. It is hence not clear what lifting annotation should be used when defining  $g$ . The lifting information at the definition-site should be correct for any call-site.

A very structured and sufficiently expressive way to deal with this situation is to use regions and subregion evidence instead. They allow us encode the lifting information and also give us the ability to abstract over it at the definition-site. This can be seen in Subfigure 2 (b), which shows what the example looks like in  $\Lambda_{\text{cap}}$ . Each handler now not only binds a capability, but also a fresh region (e.g.,  $r_1$ ) and subregion evidence (e.g.,  $n_1 : r_1 \sqsubseteq \top$ ) witnessing that the fresh region is a subregion of the current one (e.g., the toplevel region  $\top$ ). The basic idea now is to abstract the required lifting information in the form of an evidence parameter  $n$  for  $g$ , which is then provided to the call of the capability  $\text{yield}_1$  in the function body. Capability  $\text{yield}_1$  is bound at the outer handler in region  $r_1$ , so its evidence should witness that the region in which  $\text{yield}_1$  is called is a subregion of  $r_1$  and which subregion it is. To express this, we also abstract over a region parameter  $r$  for  $g$ , which stands for the region at its call-site as is visible in the annotation **at**  $r$ . The evidence parameter  $n$  is thus typed as  $r \sqsubseteq r_1$ , expressing that the call-site region  $r$  must be a subregion of  $g$ 's (and  $\text{yield}_1$ 's) definition region  $r_1$ .

When calling  $g$  under the inner handler, the current region is the region  $r_2$  bound at this handler, so we instantiate  $g$ 's region parameter with  $r_2$ . The evidence passed to  $g$  thus must have type  $r_2 \sqsubseteq r_1$ . This subregion relation is witnessed by evidence  $n_2$  bound at the inner handler. But now we can also call  $g$  immediately after defining it, in which case we instantiate its region parameter with  $r_1$ . For the evidence we then use the trivial one,  $\emptyset : r_1 \sqsubseteq r_1$ , stating that subregioning is reflexive. In either case, the evidence passed to  $\text{yield}_1$  correctly witnesses how the capability must be lifted.

To see how this lifting works, consider the CPS translation of the above program:

```

RESET(
  (Λr1. λn1. λyield1.
    λkg. (λk. k (Λr. λn. λi. n Int (yield1 i)))
      (λg. (λk1. (g r1 (Λa. λm. m) 1)
        (λx. RESET(
          (Λr2. λn2. λyield2. g r2 n2 x)
          (CPS r1 Int)
          LIFT
          (λj. λk. λk0. k0 42)
        )
        k1)))
      kg))
  (Cps Void Int)
  LIFT
  (λj. λk. (Λa. λm. m) Int (λk0. k (j + 1) k0))
)

```

Here we can see that the evidence parameters  $n_1, n_2$  bound at the handlers are eventually instantiated with the function `LIFT` which has the effect of capturing a further delimited continuation. The delimiters for continuations are given by the meta function `RESET` which is wrapped around each handler upon translation and has the effect of applying its argument to an empty continuation. For the definition of `LIFT` and `RESET` we refer to Subsection 3.3. The trivial evidence  $\emptyset$  just becomes the (polymorphic) identity function,  $\Lambda a. \lambda m. m$ .

The function  $g$  is translated to

$$\Lambda r. \lambda n. \lambda i. n \text{ Int } (\text{yield}_1 \ i)$$

In the function body we can see that the application of the capability is translated to an application of the evidence parameter  $n$  to this capability. When looking at the call sites we can then see how the lifting happens concretely. At the outer call site,  $n$  is instantiated with the identity function so that no lifting happens. At the inner call site (under the second `RESET`),  $n$  is instantiated with  $n_2$ , *i.e.*, with `LIFT`, so that here the delimited continuation which is captured does not end at the inner `RESET` but at the outer one.

To sum up this subsection, after lift inference each function-block definition should abstract a fresh region standing for the region in which the function will run, and an evidence parameter witnessing that this call-site region is a subregion of the function's definition-site region.

## 2.2 Higher-Order Functions

The example in the previous subsection only uses first-order functions. However, System  $\Xi$  also supports higher-order functions which make lift inference a bit more complicated. To see why, consider the following variation of the example from the previous subsection:

```

def call { f: Int ⇒ Int / {} } =
  val x = f(1);
  try { f(x) }
  with Yield { j ⇒ 42 }
  try {
    call { (i: Int) ⇒ do Yield(i) }
  } with Yield { j ⇒ resume(j + 1) }

```

As in the previous example, the effect operation is called under two handlers and the result of running it is the same. This time however, the inner handler is installed by a higher-order function `call`. This example motivates, why lexical effect handling can be desirable: as programmers, we



<pre> <b>def</b> call(<math>f : \text{Int} \rightarrow \text{Int}</math>) {   <b>val</b> <math>x = f(1)</math>;   <b>try</b> { <math>(\text{yield}_2) \Rightarrow f(x)</math> }   <b>with</b> { <math>(j, k) \Rightarrow 42</math> } };  <b>try</b> { <math>(\text{yield}_1) \Rightarrow</math>   call { <math>(i : \text{Int}) \Rightarrow \text{do yield}_1(i)</math> } } <b>with</b> { <math>(j, k) \Rightarrow \text{do } k(j + 1)</math> } </pre> <p>(a) Program in System <math>\Xi</math>.</p>	<pre> <b>def</b> call[<math>r_c, r_f; n_c : r_c \sqsubseteq \top, n_f : r_c \sqsubseteq r_f</math>](   <math>f : \forall [r; r \sqsubseteq r_f](\text{Int}) \rightarrow r \text{ Int}</math> ) <b>at</b> <math>r_c</math> {   <b>val</b> <math>x = f[r_c; n_f](1)</math>;   <b>try</b> { <math>[r_2; n_2 : r_2 \sqsubseteq r_c](\text{yield}_2) \Rightarrow f[r_2; n_2 \oplus n_f](x)</math> }   <b>with</b> { <math>(j, k) \Rightarrow 42</math> } };  <b>try</b> { <math>[r_1; n_1 : r_1 \sqsubseteq \top](\text{yield}_1) \Rightarrow</math>   call[<math>r_1, r_1; n_1, \emptyset</math>] { <math>[r_g; n_g](i : \text{Int}) \text{ at } r_g \Rightarrow \text{do yield}_1[n_g](i)</math> } } <b>with</b> { <math>(j, k) \Rightarrow \text{do } k[\emptyset](j + 1)</math> } </pre> <p>(b) Program in <math>\Lambda_{\text{cap}}</math>.</p>
---	--

Fig. 3. Example with higher-order function.

want to reason locally about the relation of the call to `do Yield` and its lexically enclosing handler, without having to be aware of `call`'s implementation details.

Subfigure 3 (a) shows the program in System  $\Xi$ . Again, the capability-passing translation makes the lexical relation of the effect call and the outer handler explicit.

The translation to  $\Lambda_{\text{cap}}$  in Subfigure 3 (b) is now slightly more involved. Handlers are endowed with regions and evidence as before. The block passed to call is the same as `g` in the previous example, it is just anonymous now. Its translation thus is the same as for `g`, it abstracts a fresh region  $r_g$  and evidence  $n_g : r_g \sqsubseteq r_1$  which it then passes to the capability in its body.

As the anonymous block is called indirectly via the parameter `f` of function call, we have to pass an appropriate region and evidence to `f` in the body of call. Since this region and evidence should be correct for any block `f` is instantiated with, we abstract over them in the definition of call. This way, we can provide them at the call-site of call when we know the concrete block argument we pass for `f`.

Therefore, call now abstracts two regions and two evidence parameters. Region  $r_c$  again stands for the region where call will run later and evidence  $n_c$  again witnesses that  $r_c$  is a subregion of call's definition region<sup>5</sup>. The second region  $r_f$  represents the definition region of the block argument passed for `f`. The second evidence  $n_f$  witnesses that  $r_c$  is a subregion of  $r_f$ . Moreover, in the type of `f` we can see that, as any other function block, `f` abstracts a fresh region  $r$  and corresponding evidence witnessing that  $r$  is a subregion of  $r_f$ . When `f` is called the first time in the body of call, immediately at the beginning, the current region is  $r_c$ , so we instantiate the region parameter of `f` with this region. The evidence for `f` hence has to witness the subregion relation  $r_c \sqsubseteq r_f$  which is precisely what the evidence  $n_f$  specifically abstracted for `f` does. When `f` is called the second time under the second handler, the current region is not  $r_c$  anymore but the region  $r_2$  abstracted at that handler, so we instantiate the region parameter of `f` with  $r_2$ . To obtain the correct evidence for `f` we thus have to compose the evidence  $n_f$  with evidence witnessing that  $r_2 \sqsubseteq r_c$ . This relation is precisely witnessed by the evidence  $n_2$  abstracted at the handler. The evidence passed to the second call of `f` thus is the composition  $n_2 \oplus n_f$ . In the CPS translation this composition of evidence becomes function composition, so that multiple LIFT functions can be combined to obtain the overall lift if necessary.

In the application of call, we have to provide appropriate regions and evidence. The first region argument is again the current region, which is now  $r_1$ , and the first evidence is  $n_1$ , because call

<sup>5</sup>For call, this is the toplevel region  $\top$  of which any region is a subregion, so  $n_c$  is not really necessary in this case.

was defined outside of the handler. The second region parameter must be instantiated with the definition-site region of the block argument, which is again  $r_1$ , as the anonymous block is defined in place. The second evidence hence has to witness that  $r_1 \sqsubseteq r_1$ , so we have to pass the trivial evidence  $\emptyset$ . Note that in call this evidence is then passed to the calls of the block argument and eventually to  $\text{yield}_1$ , where for the second call it is composed with  $n_2$  beforehand. Hence,  $\text{yield}_1$  indeed receives the correct evidence in both cases since composing with the trivial evidence eventually has no effect. In the CPS translation it is just composition with the identity function.

In general, after lift inference each function block should abstract an additional region and evidence parameter for each of its block parameters. The region stands for the definition-site region of the block argument and the evidence witnesses that the region in which the whole function block is called must be a subregion of the region for the block parameter.

### 2.3 Summary

Summing up, the guiding principle of our lift-inference translation is that each call-site region of a block is a subregion of its definition-site region. This is facilitated by the second-class property of blocks in System  $\Xi$ . Every function block thus abstracts a fresh region in which it will run later and evidence witnessing the above principle. When a block is called, its region parameter is instantiated with the current region and its evidence parameter with appropriate evidence. Hence, we have to track the current region during the translation and we have to remember the regions in which the blocks have been defined. We also have to keep track of the correct evidence for each block.

Moreover, for each block parameter of a function an additional region and evidence parameter is abstracted. The region represents the definition-site of the instantiation of the block parameter and the evidence witnesses that the whole function is called in a subregion of that definition-site region. As the block parameter cannot be returned, this ensures that its instantiation again satisfies our guiding principle.

## 3 TECHNICAL DEVELOPMENT

In this section, we formally present how the lift-inference translation from our version of System  $\Xi$  to our version of  $\Lambda_{\text{cap}}$  proceeds. This translation infers correct regions and evidence for a well-typed term in System  $\Xi$  to yield a well-typed term in  $\Lambda_{\text{cap}}$ . Before doing so, we recap both languages and detail the changes we have made relative to the original versions of the languages to overcome technical difficulties. Proofs and the parts of the formalization omitted here can be found in the extended technical report [Müller et al. 2023b].

### 3.1 Syntax and Type Systems

We first describe the syntax and type systems of the two calculi. In the following, source calculus means System  $\Xi$ , not to be confused with the source-level language Effekt it underlies.

**3.1.1 Source Calculus System  $\Xi$ .** We start with the source calculus System  $\Xi$ , a calculus with lexical effect handlers in explicit capability-passing style with second-class functions.

*Syntax.* The syntax of System  $\Xi$  is given in Figure 4. The calculus syntactically distinguishes potentially effectful statements from terms that cannot have control effects, *i.e.*, it is in fine-grain call-by-value [Levy et al. 2003]. Non-effectful terms are further divided into values and blocks. Values are either variables or constants, blocks are either variables or anonymous multi-arity functions. It is important to note that only values can be returned but blocks cannot, that is, blocks are second class. Still, blocks can be higher-order, *i.e.*, they cannot only abstract value parameters but also block parameters.



**Syntax of Terms:**

## Statements

$s$	$::=$	<b>val</b> $x = s; s$	sequencing
		<b>return</b> $v$	returning values
		<b>def</b> $f = b; s$	defining blocks
		$b(\bar{v}, \bar{b})$	calling blocks
		<b>do</b> $b(v)$	performing capabilities
		<b>try</b> $\{ (c) \Rightarrow s \}$ <b>with</b> $\{ (x, k) \Rightarrow s \}$	handling effects

## Values

$v$	$::=$	$x$	variables
		$  ()   0   1   \dots   \text{true}   \dots$	constants

## Blocks

$b$	$::=$	$f, k, c   w$
-----	-------	---------------

## Block Values

$w$	$::=$	$\{ (\bar{x} : \bar{\tau}, \bar{f} : \bar{\sigma}) \Rightarrow s \}$
-----	-------	--

**Syntax of Types:**

## Value Types

$\tau$	$::=$	Unit   Int   Bool   ...
--------	-------	-------------------------

## Block Types

$\sigma$	$::=$	$(\bar{\tau}, \bar{\sigma}) \rightarrow \tau$	functions
		<b>Cap</b> $\tau \tau$	capabilities

**Environments:**

## Value Environment

$\Gamma$	$::=$	$\emptyset   \Gamma, x : \tau$
----------	-------	--------------------------------

## Block Environment

$\Delta$	$::=$	$\emptyset   \Delta, f : \sigma$
----------	-------	----------------------------------

**Names:**

Value Variables  $x, y \in x, y$     Block Variables  $f, g, k, c \in f, g, k, \text{Fail}, \text{Choice}, \dots$

Fig. 4. Syntax of System  $\Xi$ .

Statements can be sequenced using **val**  $x = s_1; s_2$  where the result of  $s_1$  is bound to variable  $x$  in  $s_2$ . Defining a local block is done with **def**  $f = b; s$  making block  $b$  available in the scope of statement  $s$  by binding it to variable  $f$ . We distinguish calls of a function block from calling capabilities standing for effect operations. For the latter we add the construct **do**  $b(v)$  and also reflect this on the type level by adding an additional block type for capabilities. Syntactically distinguishing capabilities is a minor technical difference to the original version of System  $\Xi$ , but it simplifies the presentation of the lift-inference translation which treats capabilities and function blocks differently. Capabilities cannot have block parameters as this would allow blocks to leave their defining scopes and therefore break the second-class property. Finally, handling effects is done in explicit capability-passing style. In **try**  $\{ (c) \Rightarrow s_0 \}$  **with**  $\{ (x, k) \Rightarrow s \}$  the capability is bound to  $c$  in the scope of the handled statement  $s_0$ . The implementation of the capability binds its value parameter  $x$  and the continuation  $k$  in the implementation statement  $s$ .

**Value Typing**

$$\boxed{\begin{array}{c} \Gamma \vdash v : \tau \\ \uparrow \quad \uparrow \quad \downarrow \end{array}}$$

$$\frac{\Gamma(x) = \tau}{\Gamma \vdash x : \tau} [\text{VAR}] \quad \frac{}{\Gamma \vdash 3 : \text{Int}} [\text{LIT}]$$

**Block Typing**

$$\boxed{\begin{array}{c} \Gamma \mid \Delta \vdash b : \sigma \\ \uparrow \quad \uparrow \quad \uparrow \quad \downarrow \end{array}}$$

$$\frac{\Delta(f) = \sigma}{\Gamma \mid \Delta \vdash f : \sigma} [\text{BLOCKVAR}] \quad \frac{\Gamma, \overline{x : \tau} \mid \Delta, \overline{f : \sigma} \vdash s_0 : \tau_0}{\Gamma \mid \Delta \vdash \{ (\overline{x : \tau}, \overline{f : \sigma}) \Rightarrow s_0 \} : (\overline{\tau}, \overline{\sigma}) \rightarrow \tau_0} [\text{BLOCK}]$$

**Statement Typing**

$$\boxed{\begin{array}{c} \Gamma \mid \Delta \vdash s : \tau \\ \uparrow \quad \uparrow \quad \uparrow \quad \downarrow \end{array}}$$

$$\frac{\Gamma \mid \Delta \vdash s_0 : \tau_0 \quad \Gamma, x : \tau_0 \mid \Delta \vdash s : \tau}{\Gamma \mid \Delta \vdash \mathbf{val} \ x = s_0; s : \tau} [\text{VAL}]$$

$$\frac{\Gamma \vdash v : \tau}{\Gamma \mid \Delta \vdash \mathbf{return} \ v : \tau} [\text{RET}] \quad \frac{\Gamma \mid \Delta \vdash b : \sigma \quad \Gamma \mid \Delta, f : \sigma \vdash s : \tau}{\Gamma \mid \Delta \vdash \mathbf{def} \ f = b; s : \tau} [\text{DEF}]$$

$$\frac{\Gamma \mid \Delta \vdash b_0 : (\overline{\tau}, \overline{\sigma}) \rightarrow \tau_0 \quad \overline{\Gamma \vdash v : \tau} \quad \overline{\Gamma \mid \Delta \vdash b : \sigma}}{\Gamma \mid \Delta \vdash b_0(\overline{v}, \overline{b}) : \tau_0} [\text{APP}]$$

$$\frac{\Gamma \mid \Delta \vdash b : \mathbf{Cap} \ \tau_1 \ \tau_2 \quad \Gamma \vdash v : \tau_1}{\Gamma \mid \Delta \vdash \mathbf{do} \ b(v) : \tau_2} [\text{DO}]$$

$$\frac{\Gamma \mid \Delta, c : \mathbf{Cap} \ \tau_1 \ \tau_2 \vdash s_0 : \tau \quad \Gamma, x : \tau_1 \mid \Delta, k : \mathbf{Cap} \ \tau_2 \ \tau \vdash s : \tau}{\Gamma \mid \Delta \vdash \mathbf{try} \ \{ (c) \Rightarrow s_0 \} \ \mathbf{with} \ \{ (x, k) \Rightarrow s \} : \tau} [\text{TRY}]$$

Fig. 5. Type system of System  $\Xi$ .

*Typing rules.* Figure 5 defines the typing rules for System  $\Xi$ . Typing for values is entirely standard. Note that values as well as statements are typed against value types. In contrast, blocks are typed against block types, that is, they have either function type or capability type. Moreover, there are two kinds of environments in the typing judgment of blocks, namely  $\Gamma$  for value bindings and  $\Delta$  for block bindings. The same is true for statement typing. In particular, when typing a function block, the value parameters are added to the value environment and the block parameters are added to the block environment.

Apart from distinguishing two kinds of environment, the rules for sequencing (VAL), returning (RET), block definition (DEF) and function block calls (APP) are standard. Compared to the original version of System  $\Xi$  we have an additional rule (DO) for performing capabilities, a consequence of syntactically distinguishing them from function blocks as mentioned above. A capability has type  $\mathbf{Cap} \ \tau_1 \ \tau_2$  where  $\tau_1$  is the type of its parameter and  $\tau_2$  is the return type. Otherwise the rule is essentially the same as the one for calling function blocks. The crucial rule is TRY. The handler makes the capability available in the scope of the handled statement by adding a binding for it to the block environment for the handled statement. Since blocks cannot be returned, the capability can only be used in that scope, thus guaranteeing effect safety without any visible effect system. In the implementation statement, the continuation parameter has capability type. This is in contrast with the original version of System  $\Xi$  where it has function type. The reason for treating continuations

**Syntax of Terms:**

Statements		
$s$	$::=$ <b>val</b> $x = s$ ; $s$	sequencing
	<b>return</b> $v$	returning values
	$v[\bar{\rho}; \bar{e}](\bar{v})$	calling functions
	<b>do</b> $v[e](v)$	performing capabilities
	<b>try</b> $\{ [r; n](c) \Rightarrow s \}$ <b>with</b> $\{ (x, k) \Rightarrow s \}$	handling effects
Values		
$v$	$::=$ $x, f, k, c$	variables
	$() \mid 0 \mid 1 \mid \dots \mid \text{true} \mid \dots$	constants
	$\{ [\bar{r}; \bar{n} : \bar{\gamma}](\bar{x} : \bar{\tau}) \text{ at } \rho \Rightarrow s \}$	closures
Evidence		
$e$	$::=$ $n, \dots$	evidence variables
	$0$	reflexive evidence
	$e \oplus e$	transitive evidence

**Syntax of Types:**

Types		
$\tau$	$::=$ $\text{Unit} \mid \text{Int} \mid \text{Bool} \mid \dots$	primitives
	$\forall [\bar{r}; \bar{\gamma}](\bar{\tau}) \rightarrow \rho \tau$	functions
	<b>Cap</b> $\rho \tau \tau$	capabilities
Regions		
$\rho$	$::=$ $r$	region variable
	$\top$	oplevel region
Constraints		
$\gamma$	$::=$ $\rho \sqsubseteq \rho$	subregion

**Environments:**

$\Gamma$	$::=$ $\emptyset$	empty environment
	$\Gamma, r$	region binding
	$\Gamma, n : \gamma$	evidence binding
	$\Gamma, x : \tau$	value binding

**Names:**

Variables  $x, y, f, g, k, c \in x, \gamma, f, g, k$  Fail, Choice, ...

**Syntactic Sugar:**

**def**  $f = v \doteq$  **val**  $f =$  **return**  $v$

Fig. 6. Syntax of  $\Lambda_{\text{cap}}$ .

as capabilities is to simplify the translation to our generalized version of  $\Lambda_{\text{cap}}$ . For System  $\Xi$  this does not add expressivity as capabilities can be used in the same contexts as function blocks.

**3.1.2 Target Calculus  $\Lambda_{\text{cap}}$ .** The target calculus  $\Lambda_{\text{cap}}$  also is a calculus with lexical effect handlers in explicit capability-passing style. In contrast to System  $\Xi$ , it features first-class functions and has explicit regions and subregion evidence.

*Syntax.* The syntax is given in Figure 6. It is again in fine-grain call-by-value, but does not distinguish blocks from values. So blocks are now values, making them first class. This is also

**Value Typing**

$$\boxed{\begin{array}{c} \Gamma \vdash v : \tau \\ \uparrow \quad \uparrow \quad \downarrow \end{array}}$$

$$\frac{\Gamma(x) = \tau}{\Gamma \vdash x : \tau} [\text{VAR}] \quad \frac{}{\Gamma \vdash 3 : \text{Int}} [\text{LIT}]$$

$$\frac{\Gamma, \bar{r}, \bar{n} : \bar{y}, \bar{x} : \bar{\tau} \mid \rho \vdash s_0 : \tau_0}{\Gamma \vdash \{ [\bar{r}; \bar{n} : \bar{y}] (\bar{x} : \bar{\tau}) \text{ at } \rho \Rightarrow s_0 \} : \forall [\bar{r}; \bar{y}] (\bar{\tau}) \rightarrow \rho \tau_0} [\text{FUN}]$$

**Statement Typing**

$$\boxed{\begin{array}{c} \Gamma \mid \rho \vdash s : \tau \\ \uparrow \quad \uparrow \quad \uparrow \quad \downarrow \end{array}}$$

$$\frac{\Gamma \mid \rho \vdash s_0 : \tau_0 \quad \Gamma, x_0 : \tau_0 \mid \rho \vdash s : \tau}{\Gamma \mid \rho \vdash \text{val } x_0 = s_0; s : \tau} [\text{VAL}] \quad \frac{\Gamma \vdash v : \tau}{\Gamma \mid \rho \vdash \text{return } v : \tau} [\text{RET}]$$

$$\frac{\Gamma \vdash v_0 : \forall [\bar{r}; \bar{y}] (\bar{\tau}) \rightarrow \rho_0 \tau_0 \quad \frac{}{\Gamma \vdash e : \gamma [\bar{r} \mapsto \rho]} \quad \frac{}{\Gamma \vdash v : \tau [\bar{r} \mapsto \rho]} \quad \rho'_0 = \rho_0 [\bar{r} \mapsto \rho]}{\Gamma \mid \rho'_0 \vdash v_0 [\bar{\rho}; \bar{e}] (\bar{v}) : \tau_0 [\bar{r} \mapsto \bar{\rho}]} [\text{APP}]$$

$$\frac{\Gamma \vdash v_0 : \mathbf{Cap} \rho' \tau_1 \tau_2 \quad \Gamma \vdash e : \rho \sqsubseteq \rho' \quad \Gamma \vdash v : \tau_1}{\Gamma \mid \rho \vdash \mathbf{do} v_0 [e] (v) : \tau_2} [\text{DO}]$$

$$\frac{\Gamma, r, n : r \sqsubseteq \rho, c : \mathbf{Cap} r \tau_1 \tau_2 \mid r \vdash s_0 : \tau \quad \Gamma, x : \tau_1, k : \mathbf{Cap} \rho \tau_2 \tau \mid \rho \vdash s : \tau}{\Gamma \mid \rho \vdash \mathbf{try} \{ [r; n] (c) \Rightarrow s_0 \} \mathbf{with} \{ (x, k) \Rightarrow s \} : \tau} [\text{TRY}]$$

**Evidence Typing**

$$\boxed{\begin{array}{c} \Gamma \vdash e : \gamma \\ \uparrow \quad \uparrow \quad \downarrow \end{array}}$$

$$\frac{\Gamma(n) = \rho_1 \sqsubseteq \rho_2}{\Gamma \vdash n : \rho_1 \sqsubseteq \rho_2} [\text{EVIVAR}] \quad \frac{}{\Gamma \vdash \mathbf{0} : \rho \sqsubseteq \rho} [\text{REFLEXIVE}]$$

$$\frac{\Gamma \vdash e_1 : \rho \sqsubseteq \rho' \quad \Gamma \vdash e_2 : \rho' \sqsubseteq \rho''}{\Gamma \vdash e_1 \oplus e_2 : \rho \sqsubseteq \rho''} [\text{TRANSITIVE}]$$

Fig. 7. Type system of  $\Lambda_{\text{cap}}$ .

reflected on the level of types in that function and capability types are not in a separate syntactic category. But, to ensure effect safety, there are regions and subregion evidence. Regions are either region variables or the toplevel region. Evidence is a witness of the subregion relationship between regions. It is either an evidence variable, the trivial evidence  $\mathbf{0}$ , or the composition  $e \oplus e$  of evidence.

Functions do not distinguish value and block parameters but can now additionally abstract over a list of regions and a list of evidence. Accordingly, when calling a function, corresponding lists of regions and evidence have to be supplied. Furthermore, each function is annotated with the region it is supposed to run in. When calling a capability, it needs to be supplied with evidence but not with a region. An effect handler not only abstracts a capability but additionally a region and an evidence variable for the handled statement, whereas the implementation statement stays the same as in System  $\Xi$ . The constructs for sequencing and return are also the same. There is no construct for the definition of a local function, as it can be easily defined as syntactic sugar using sequencing now that functions are first class.

*Typing rules.* Figure 7 shows the typing rules for  $\Lambda_{\text{cap}}$ . In contrast to System  $\Xi$ , there is only one typing environment, however, in addition to value bindings it can also contain regions and

evidence bindings. Moreover, the typing judgment for statements has as an additional component the region in which the statement is typed.

As functions can abstract region and evidence parameters, these are added to the environment when typing the function body (rule FUN). Moreover, the function has to run in the region its body is typed in. This is also visible in the type of the function. When sequencing two statements, both are typed in the same region as the compound statement. Returning a result can be typed in any region. When applying a function (rule APP), its region arguments are substituted not only in the return type, but also in the types for the evidence and value arguments when typing them. Furthermore, when substituted in the region annotated in the function type, the resulting region has to coincide with the region the applied function is typed in. This allows functions to be region-polymorphic.

Rule DO defines the typing for performing capabilities. A capability can only be called in a subregion of the region annotated in its type. This is witnessed by the evidence supplied in the call. The region annotated in the type of the capability  $c$  abstracted at the handled statement  $s_0$  of an effect handler (rule TRY) is the fresh region  $r$  also abstracted there. Region  $r$  is also the region which  $s_0$  is typed in. Thus,  $c$  can only be called in a subregion of  $r$ , *i.e.*, within the handled statement. This ensures effect safety. The additionally abstracted evidence  $n$  witnesses that  $r$  is a subregion of the region  $\rho$  the overall statement is typed in. The regions and subregion evidence hence precisely reflect how handlers are nested. The implementation statement is typed in the outer region  $\rho$  which is also the region for the continuation. In contrast to the original version of  $\Lambda_{\text{cap}}$ , the continuation has capability type, not function type. This allows to call the continuation not only in region  $\rho$  but also in any subregion of  $\rho$ . Our version is thus more expressive as continuations can, for example, be called under further effect handlers, *i.e.*, we lift the restriction of scoped continuations imposed by the original version. This is important in order to fully support a lift-inference translation from System  $\Xi$  since there no such restriction for continuations exists.

The typing of evidence is rather straightforward. Evidence variables are looked up in the environment. The trivial evidence  $\emptyset$  witnesses that any region is a subregion of itself. Composition of evidence shows that subregioning is transitive.

### 3.2 Operational Semantics

We define the operational semantics of the two calculi in terms of an abstract machine. The abstract machine for  $\Lambda_{\text{cap}}$  is essentially the same as the one for System  $\Xi$ . There are only two minor differences. First, there is no stepping rule for the definition of local functions in  $\Lambda_{\text{cap}}$  as there is no separate construct for them. Second, there are regions and evidence. However, the latter are irrelevant for the machine semantics as it proceeds by searching delimiters with labels on the meta stack and does not use region and evidence information.

As the operational semantics of both languages is so similar, it is not that interesting with regards to lift inference. Nevertheless, knowing the operational semantics helps to understand lifting, in particular with respect to continuations. Therefore, we briefly sketch how the machine works, and also where it differs from the original versions of the two calculi.

A machine state  $\langle s \parallel K \rangle$  contains a statement  $s$  to be evaluated and a runtime meta stack  $K$  which is a list of delimited stacks. A stack is a list of frames ending with a delimiter  $\#_l$  containing a label  $l$ . This is a minor technical difference to the original version of the machine which treats delimiters as regular frames and does not explicitly segment the meta stack into delimited stacks. It facilitates the correctness proof of the CPS translation from our generalized version of  $\Lambda_{\text{cap}}$  to System F.

The machine implements multi-prompt delimited control [Dybvig et al. 2007]. Upon execution of a handler statement, a fresh label is generated and a new stack consisting only of a delimiter with that label is pushed onto the meta stack.

$$(try) \quad \langle \mathbf{try} \{ (c) \Rightarrow s_0 \} \mathbf{with} \{ (x, k) \Rightarrow s \} \parallel K \rangle \rightarrow \langle s_0 [c \mapsto v] \parallel \#_l :: K \rangle$$

where  $l = \text{generateFresh}()$  and  $v = \mathbf{cap}_l \{ (x, k) \Rightarrow s \}$

Execution then continues with the handled statement where the abstracted capability variable is replaced by a runtime capability which contains the just generated label  $l$  and the handler implementation. When encountering a call to a capability the machine transitions to unwinding mode and pops stacks off the meta stack until the correct label is found. These stacks are collected in a resumption which is used as continuation.

It is exactly this runtime search of correct handlers that we seek to avoid by using lifting information. To make this possible, we make sure that the evidence passed to capabilities precisely reflects the labels on the runtime meta stack and we take care that this invariant is preserved during the execution of the machine. In fact, at runtime each evidence becomes a list of appropriate labels which is adapted as the machine proceeds. The invariant ensures that evidence always contains the correct lifting information. While being irrelevant for the machine semantics, this is critically important for the CPS translation, since the latter uses evidence instead of labels (see Subsection 3.3).

The main difference to the original version of the machine is how a call of a continuation proceeds. The reason we have to treat this differently is that, compared to the original version of  $\Lambda_{\text{cap}}$ , we allow the use of continuations under further handlers, in order to fully support System  $\Xi$  in the lift-inference translation. A continuation may contain capabilities which have been provided with evidence. When calling the continuation under further handlers in the implementation statement, additional delimiters are installed on the meta stack that were not present when the continuation was created. Thus, the evidence for the capabilities inside the continuation does not precisely reflect the labels on the meta stack which violates the critical invariant described above.

To make the evidence correct again, we have to make the additional delimiters “invisible” for the continuation. This is achieved by treating continuations as capabilities as well. This way, when a continuation is called, it first captures the additional stacks with these delimiters by unwinding as described above (which is again reflected by the evidence passed to the continuation itself), and packages them into one resumption frame. Only then, the continuation is executed in the usual way by rewinding. Such a resumption frame acts a bit like an “underflow” frame [Farvardin and Reppy 2020] when returning to it, in the sense that execution then first continues with that resumption. When unwinding, however, it is treated just as another ordinary frame so that the stacks inside of it do not interfere with the unwinding. Hence, when a call to a capability inside the continuation is encountered, it only sees the labels present on the meta stack when the continuation was created so that its evidence is correct.

This difference in how continuations are treated compared to the original version does not impact the final result of the execution for all programs that can be written in the original versions of the calculi. It leads, however, to another minor difference to the original machine. As the continuation capabilities are always delimited by the next label on the meta stack at the point of their creation, execution of a closed statement  $s$  always starts with a delimiter with a special toplevel label on the otherwise empty meta stack, that is, in state  $\langle s \parallel \#_{\text{start}} :: \bullet \rangle$ . This ensures that there also is a delimiting label for continuations of  $\mathbf{try}$ -statements in the toplevel region.

### 3.3 CPS Translation to System F

For the original version of  $\Lambda_{\text{cap}}$  Schuster et al. [2022b] give a typability- and semantics-preserving CPS translation to pure System F. This CPS translation carries over almost unchanged to our version of  $\Lambda_{\text{cap}}$ . Still, it is instructive to briefly repeat the core idea, in particular, to see how evidence enables efficient compilation.



**Translation of Statements:**

$$\begin{aligned}
& \vdots \\
\mathcal{S}[\![ \text{do } v_0[e](v) ]\!] &= \mathcal{E}[\![ e ]\!] \mathcal{T}[\![ \tau_2 ]\!] (\mathcal{V}[\![ v_0 ]\!] \mathcal{V}[\![ v ]\!]) \\
\mathcal{S}[\![ \text{do } k[e](v) ]\!] &= \mathcal{E}[\![ e ]\!] \mathcal{T}[\![ \tau_2 ]\!] (\lambda k_0. \mathcal{V}[\![ k ]\!] \mathcal{V}[\![ v ]\!] k_0) \\
\mathcal{S}[\![ \text{try } \{ [r; n](c) \Rightarrow s_0 \} \text{ with } \{ (x, k) \Rightarrow s \} ]\!] &= \text{RESET} ((\Lambda r. \lambda n. \lambda c. \mathcal{S}[\![ s_0 ]\!]) \\
&\quad (\text{CPS } \mathcal{T}[\![ \rho ]\!] \mathcal{T}[\![ \tau ]\!]) (\text{LIFT}) (\lambda x. \lambda k. \mathcal{S}[\![ s ]\!]))
\end{aligned}$$

**Auxiliary Definitions:**

$$\begin{aligned}
\text{CPS } R A &= (A \rightarrow R) \rightarrow R \\
\text{RESET} &: \text{CPS } (\text{CPS } R A) A \rightarrow \text{CPS } R A & \text{LIFT} &: \forall a. \text{CPS } R a \rightarrow \text{CPS } (\text{CPS } R R') a \\
\text{RESET } m &= m (\lambda x. \lambda k. k x) & \text{LIFT} &= \Lambda a. \lambda m. \lambda k. \lambda j. m (\lambda x. k x j)
\end{aligned}$$

Fig. 8. CPS Translation from  $\Lambda_{\text{cap}}$  to System F.

The idea of the CPS translation is to use evidence information to decide how to lift a capability, *i.e.*, which parts of the runtime meta stack need to be captured. Or put in terms of the operational semantics, which delimiters have to be jumped over when unwinding. As a result, no runtime search for the correct label on the meta stack is needed anymore. To this end, the translation targets so-called iterated CPS [Schuster and Brachthäuser 2018], which uses one continuation parameter for each stack delimited by a label.

The full CPS translation is omitted here, Figure 8 only shows how handlers and calls of capabilities are translated. Note that the translation is actually defined over typing derivations, but we only write the term here. In the translation of a handler the handled statement becomes a function applied to three arguments. The region parameter  $r$  represents the polymorphic answer type that has to be instantiated appropriately ( $\tau$  and  $\rho$  are the overall type and region of the **try**-statement as in typing rule TRY), the evidence variable  $n$  is instantiated with the function LIFT and the capability parameter  $c$  is instantiated with the translated implementation statement. The whole term is then applied to an empty continuation acting as a delimiter by meta function RESET. The function LIFT increases the number of continuation parameters of its argument  $m$  by one, hence  $m$  is lifted to a different region by capturing one more delimited stack. Note that the explicitly abstracted type parameter  $a$  is the immediate return type (see also the application of a capability below), while the (answer) types  $R, R'$  are determined by the surrounding regions.

In the CPS translation of performing capabilities the capability is applied to its argument. The resulting term is then fed into the translated evidence of the capability (in the type argument  $\mathcal{T}[\![ \tau_2 ]\!]$ ,  $\tau_2$  is the return type of the capability as in typing rule Do). This evidence always eventually consists of a composition of evidence variables bound at handlers<sup>6</sup>, which means that it is a composition of LIFT-terms. Thus, the translated evidence term determines how far the capability is lifted, that is, how many stacks of the meta stack are captured.

As explained above we also treat the continuation as a capability and provide it with evidence. This ensures that the handler capabilities inside the continuation are correctly lifted since the continuation itself is lifted to the correct region. The CPS translation for continuation capabilities is almost the same as for handler capabilities, the only difference being that the applied continuation is  $\eta$ -expanded. This is necessary in order to make the following simulation theorem true.

<sup>6</sup>It may further be interspersed with trivial evidence which is, however, translated to the identity function.

**THEOREM 3.1 (SIMULATION FOR THE CPS TRANSLATION).**

If  $\vdash M$  ok and  $M \rightarrow M'$ , then  $\mathcal{M}[\![ M ]\!] \rightarrow^* \mathcal{M}[\![ M' ]\!]$ .

$M$  denotes a machine state. The operational semantics of  $\Lambda_{\text{cap}}$  hence corresponds to reduction in System F. As a corollary we obtain that evaluation is preserved by the CPS translation.

**COROLLARY 3.2 (EVALUATION FOR THE CPS TRANSLATION).**

If  $\emptyset \vdash s : \text{Int}$  and  $\langle s \parallel \#_{\text{start}} :: \bullet \rangle \rightarrow^* \langle \text{return } v \parallel \bullet \rangle$ ,  
then  $\text{RESET}(\mathcal{S}[\![ s ]\!]) \text{ done} \rightarrow^* \text{done } \mathcal{V}[\![ v ]\!]$ .

Here done is a special toplevel continuation and the RESET is necessary due to the presence of the toplevel label. As the operational semantics of our version of  $\Lambda_{\text{cap}}$  differs a bit from the original version, the proof of the simulation theorem and the necessary translation of the runtime constructs had to be adapted. In contrast, the proof of the following typability preservation stays almost unchanged.

**THEOREM 3.3 (TYPABILITY PRESERVATION FOR THE CPS TRANSLATION).**

If  $\Gamma \vdash \rho \vdash s : \tau$ , then  $\mathcal{T}[\![ \Gamma ]\!] \vdash \mathcal{S}[\![ s ]\!] : \text{CPS } \mathcal{T}[\![ \rho ]\!] \mathcal{T}[\![ \tau ]\!]$ .  
If  $\Gamma \vdash v : \tau$ , then  $\mathcal{T}[\![ \Gamma ]\!] \vdash \mathcal{V}[\![ v ]\!] : \mathcal{T}[\![ \tau ]\!]$ .

### 3.4 Lift-Inference Translation

We now present the lift-inference translation from System  $\Xi$  to  $\Lambda_{\text{cap}}$ . This is our main contribution. The translation is defined over typing derivations of System  $\Xi$  and is supposed to take well-typed terms in System  $\Xi$  to well-typed terms in  $\Lambda_{\text{cap}}$ , so we translate types and terms. In the clauses for the terms we only write the term instead of the whole typing derivation.

The translation is defined in Figure 9. As the two calculi are quite similar, the translation mainly proceeds by endowing terms in System  $\Xi$  with appropriate region and evidence abstractions and applications in the right places. To this end, we maintain a lifting environment  $E$  during the translation which is used to remember which blocks have been bound so far and in what region, while we descend recursively.

This environment is modeled as a record consisting of four components. First, it contains the current region  $\rho$  of the term to be translated. Second, it contains the block environment  $\Delta$  of the term to be translated. The types of the blocks in  $\Delta$  are not needed in the lifting environment and we usually omit them, but it eases presentation a bit to just write  $\Delta$ . The third component is a map  $m$  from the domain  $\text{dom}(\Delta)$  of the block environment to pairs of regions (Reg) and evidence (Ev). The region stands for the region of the definition-site of the entry and the evidence is supposed to witness that the current region  $\rho$  is a subregion of the definition-site region. This invariant is enabled by the second-class property of blocks which guarantees that each call-site of a block is in a subregion of the region of the definition-site. To maintain the invariant, the evidence for each block in the map has to be adapted when the current region changes during the translation. As a fourth component the lifting environment contains a typing environment  $\Gamma_E$  which consists of all the regions and evidence that are present in map  $m$ . The above invariant can now more formally be captured in the following definition.

*Definition 3.4 (Soundness of Region-and-Evidence Environment).*

We call  $E = \{ \rho, \Delta, m, \Gamma_E \}$  sound if  $\Gamma_E \vdash C[\![ f ]\!]^E : \rho \sqsubseteq \mathcal{R}[\![ f ]\!]^E$  for all  $f \in \text{dom}(\Delta)$ .

Here the functions  $\mathcal{R}[\![ \cdot ]\!]$  and  $C[\![ \cdot ]\!]$  are lookup functions for the region and evidence component of blocks in the map  $m$ , respectively (see Figure 9).

**Region-and-Evidence Environment:**

$$E = \{ \rho, \Delta, m : \text{Map}(\text{dom}(\Delta), \text{Reg} \times \text{Ev}), \Gamma_E \}$$

**Translation of Types:**

$$\mathcal{T}[(\tau) \rightarrow \tau_0]_\rho = \forall [r; r \sqsubseteq \rho] (\tau) \rightarrow r \tau_0$$

where  $r = \text{generateFresh}()$

$$\mathcal{T}[(\sigma) \rightarrow \tau_0]_\rho = \forall [r, r_f; r \sqsubseteq \rho, r \sqsubseteq r_f] (\mathcal{T}[\sigma]_{r_f}) \rightarrow r \tau_0$$

where  $r, r_f = \text{generateFresh}()$

$$\mathcal{T}[\text{Cap } \tau_1 \tau_2]_\rho = \text{Cap } \rho \tau_1 \tau_2$$

**Translation of Environments:**

$$\mathcal{T}[\emptyset]^E = \emptyset$$

$$\mathcal{T}[\Delta, f : \sigma]^E = \mathcal{T}[\Delta]^E, f : \mathcal{T}[\sigma]_{\mathcal{R}[f]^E}$$

**Translation of Blocks:**

$$B[f]^E = f$$

$$B[\{ (x : \tau) \Rightarrow s_0 \}]^E = \{ [r; n : r \sqsubseteq E.\rho](x : \tau) \text{ at } r \Rightarrow \mathcal{S}[s_0]^E \oplus n \}$$

where  $r, n = \text{generateFresh}()$

$$B[\{ (f : \sigma) \Rightarrow s_0 \}]^E =$$

$$\{ [r, r_f; n : r \sqsubseteq E.\rho, n_f : r \sqsubseteq r_f](f : \mathcal{T}[\sigma]_{r_f}) \text{ at } r \Rightarrow \mathcal{S}[s_0]^{E'} \}$$

where  $r, n, r_f, n_f = \text{generateFresh}()$  and  $E' = E \oplus n, f \mapsto (r_f \mid n_f), r_f, n_f : r \sqsubseteq r_f$

**Translation of Statements:**

$$\mathcal{S}[\text{return } v]^E = \text{return } v$$

$$\mathcal{S}[\text{val } x = s_0; s]^E = \text{val } x = \mathcal{S}[s_0]^E; \mathcal{S}[s]^E$$

$$\mathcal{S}[\text{def } f = b; s]^E = \text{def } f = B[b]^E; \mathcal{S}[s]^E, f \mapsto (\mathcal{R}[b]^E \mid C[b]^E)$$

$$\mathcal{S}[b(v)]^E = B[b]^E[E.\rho; C[b]^E](v)$$

$$\mathcal{S}[b(b_0)]^E = B[b]^E[E.\rho, \mathcal{R}[b_0]^E; C[b]^E, C[b_0]^E](B[b_0]^E)$$

$$\mathcal{S}[\text{do } c(v)]^E = \text{do } c[C[c]^E](v)$$

$$\mathcal{S}[\text{try } \{ (c) \Rightarrow s_0 \} \text{ with } \{ (x, k) \Rightarrow s \}]^E =$$

$$\text{try } \{ [r; n : r \sqsubseteq E.\rho](c) \Rightarrow \mathcal{S}[s_0]^E \oplus n, c \mapsto (r \mid 0) \} \text{ with } \{ (x, k) \Rightarrow \mathcal{S}[s]^E, k \mapsto (E.\rho \mid 0) \}$$

where  $r, n = \text{generateFresh}()$

**Lookup and Adaptions of Region-and-Evidence Environment:**

$$\mathcal{R}[f]^E = \rho \quad \text{where } (\rho, e) = E.m(f)$$

$$\mathcal{R}[w]^E = E.\rho$$

$$C[f]^E = e \quad \text{where } (\rho, e) = E.m(f)$$

$$C[w]^E = \emptyset$$

$$\emptyset \oplus n = \emptyset$$

$$(m, f \mapsto (\rho \mid n_0)) \oplus n = m \oplus n, f \mapsto (\rho \mid n \oplus n_0)$$

$$\{ \rho, \Delta, m, \Gamma_E \} \oplus n = \{ r, \Delta, m \oplus n, (\Gamma_E, r, n : r \sqsubseteq \rho) \} \quad \text{where } n : r \sqsubseteq \rho$$

$$\{ \rho, \Delta, m, \Gamma_E \}, f \mapsto (\rho' \mid e) = \{ \rho, (\Delta, f), (m, f \mapsto (\rho' \mid e)), \Gamma_E \}$$

$$\{ \rho, \Delta, m, \Gamma_E \}, r, n : \rho \sqsubseteq r = \{ \rho, \Delta, m, (\Gamma_E, r, n : \rho \sqsubseteq r) \}$$

Fig. 9. Lift-Inference Translation from System  $\Xi$  to  $\Lambda_{\text{cap}}$ .

*Translation of values.* Values in System  $\Xi$  are either variables or constants, both of which are translated trivially to the same terms in  $\Lambda_{\text{cap}}$ . Similarly, value types  $\tau$  in System  $\Xi$  are only base types and thus remain unchanged. Accordingly, value environments  $\Gamma$  need not be translated. Therefore, these parts of the calculus are omitted from the presentation in Figure 9.

*Translation of blocks.* Blocks in System  $\Xi$  are translated to values in  $\Lambda_{\text{cap}}$ . Just as value variables, block variables are translated trivially. For function blocks we only show the case of a function with one parameter and treat the cases for a value parameter and a block parameter separately to ease presentation. Multi-arity functions are translated accordingly in the obvious way.

In either case the translated function abstracts a fresh region  $r$  and fresh evidence  $n$  which witnesses that  $r$  is a subregion of the current region. Moreover, the function is annotated to run in the abstracted region  $r$ , *i.e.*, it is region-polymorphic, but the evidence enforces the constraint that the actual region in which the function will run must be a subregion of the current region of the definition-site.

In the case of a value parameter the body of the function is translated recursively but the lifting environment  $E$  is adapted to  $E \oplus n$ . This has three effects. The current region is changed to be the abstracted region  $r$ . The evidence component of all entries in the map  $m$  of  $E$  is composed with the additional evidence  $n$ . This is necessary to maintain the above-mentioned soundness invariant since the current region has changed. Moreover,  $r$  and  $n$  are added to the typing environment  $\Gamma_E$ .

In the case of a block parameter  $f : \sigma$  the translated function abstracts an additional region  $r_f$  which stands for the region of the definition-site of  $f$  and an additional evidence parameter  $n_f : r \sqsubseteq r_f$  witnessing that the region  $r$  the function will run in is a subregion of  $r_f$ . As the definition-site region of  $f$  is only known when it is actually instantiated, it is necessary to abstract over it. This region is also used to translate the type  $\sigma$  of  $f$ . The constraint that  $n_f$  imposes says that the block the parameter  $f$  later is instantiated with must be defined in a superregion of the region in which the whole function is called. The lifting environment for the translation of the body is first adapted in the same way as in the case of a value parameter, but must then be further adapted by adding an entry for  $f$ . This entry consists of the pair  $(r_f \upharpoonright n_f)$  which satisfies the soundness invariant since the current region now is  $r$ . Moreover,  $r_f$  and  $n_f$  must be added to the typing environment  $\Gamma_E$ .

Note that extending the lifting environment  $E$  with new entries for blocks and with additional region and evidence variables is both written as comma-separated concatenation, but that the two extensions affect different components of  $E$ .

*Translation of block types.* The translation of block types does not need the lifting environment as additional input but only a region standing for the region of the definition-site of the block. For functions, the additionally abstracted region and evidence parameters for the translated function itself and each of its block parameters are directly reflected in the type. For capability types the given region is simply added as the region for the capability. The translation of block environments proceeds by pointwise translation of the types of the block bindings. However, as the region input must be the definition-site region for each block, we have to look this region up in the lifting environment  $E$  using the lookup function  $\mathcal{R}[\cdot]$  for regions. The translation of block environments therefore does need  $E$  as input.

*Translation of statements.* The translation for returning values is trivial and for sequencing of statements we simply translate the substatements recursively with the same environment. The definition of a local block  $b$  is a bit more interesting. It is translated to the definition of the translated block (note that this is syntactic sugar in  $\Lambda_{\text{cap}}$ ), but for the translation of the remaining statement we have to adapt the lifting environment  $E$  by inserting an entry for this newly defined function. Now there are two cases for  $b$ . Either it is a block variable  $g$  (*i.e.*, the definition is just aliasing), then it must be in the environment and we have to look up the correct region and evidence for  $g$  in  $E$ . Or it is a block value  $w$ , then there is no binding in the environment yet. In this case, the region of the definition-site of the block is the current region and hence the correct evidence is  $\emptyset$ . The

translation for this case could thus instead be defined as

$$\mathcal{S}[\mathbf{def} f = w; s]^E = \mathbf{def} f = B[\![w]\!]^E; \mathcal{S}[s]^E + f \mapsto (E, \rho \uparrow 0)$$

However, since the lookup functions  $\mathcal{R}[\![\cdot]\!]$  and  $\mathcal{C}[\![\cdot]\!]$  are defined to yield exactly the above results for block values, we do not need to distinguish cases in the translation in Figure 9.

Calling a function block  $b$  is translated to calling the translated block. We again only show the cases for one value argument and one block argument, but the generalization to multi-arity functions is again done in the obvious way. In either case the translated block has abstracted a region and an evidence parameter. The region parameter stands for the region in which the function runs, so we instantiate it with the current region. Remember that the evidence parameter must witness that the region parameter is a subregion of the definition-site region of  $b$ . But as the region parameter was instantiated with the current region we can exploit the carefully maintained soundness and simply look the correct evidence up in the lifting environment. Note that if  $b$  is an anonymous block value, the definition-site region is the current region and so the trivial evidence yielded by the lookup function for evidence is correct. In the case of a block parameter we have to translate the block argument  $b_0$  as well, of course. But we additionally have to supply a region and evidence for the corresponding parameters that have been abstracted for the block parameter. Both of these can simply be looked up, too, since the region lookup will yield the region of the definition-site of  $b_0$  and soundness again makes sure that the corresponding evidence is correct.

Translating an applied capability is similar to calling a function with a value parameter, the needed evidence is simply looked up. The difference is just that no region needs to be supplied as capabilities are not region-polymorphic. The region of their definition-site always is the fresh region abstracted at the translation of the corresponding handler statement. This can be seen in the translation of a handler which consists of the translation of the handler statement and the translation of the implementation statement. The former is similar to how function blocks with a block parameter are translated. The lifting environment  $E$  is first adapted with the newly abstracted evidence  $n$  and then an entry for the capability  $c$  is added. As the region for this entry now is the current one, the corresponding evidence is trivial, *i.e.*,  $\emptyset$ . Note that since the fresh abstracted region is already added to the typing environment  $\Gamma_E$  by  $E \oplus n$ , it is not necessary to do this in an extra step. For the implementation statement, no region and evidence is abstracted, so we only have to add an entry for the continuation parameter  $k$  to  $E$ . The region for  $k$  is the current one, that is, the one the whole handler is defined in. The evidence for  $k$  thus again is  $\emptyset$ .

**3.4.1 Example.** To illustrate the lift-inference translation, we consider again the example from Subsection 2.2, or more precisely the definition of call.

```
def call = { (f : Int → Int) ⇒
  val x = f(1);
  try { (yield2) ⇒ f(x) }
  with { (j, k) ⇒ 42 }
};
```

Starting with the empty lifting environment  $E = \{ \top, \emptyset, \emptyset, \emptyset \}$ , we show how the environment changes as the translation proceeds. In the first step, regions and evidence for call and  $f$  are abstracted, the type of  $f$  is translated with the abstracted region and the region annotation is added. Thus, we obtain

```
def call = { [rc, rf; nc : rc ⊆ ⊤, nf : rc ⊆ rf](
  f : ∀[r; r ⊆ rf](Int) →r Int
) at rc ⇒  $\mathcal{S}[\![\dots]\!]^E$ 
};
```

The lifting environment for the recursive translation of the body is adapted since we have entered a new region and since we have to add a new entry for the block parameter. It becomes

$$\begin{aligned} E' &= E \oplus n_c, f \mapsto (r_f \mid n_f), r_f, n_f : r_c \sqsubseteq r_f \\ &= \{ r_c, (f), (f \mapsto (r_f \mid n_f)), (r_c, n_c : r_c \sqsubseteq \top, r_f, n_f : r_c \sqsubseteq r_f) \} \end{aligned}$$

In the translation of the body, the translation of the first application of  $f$  is a simple matter of looking up the current region and the evidence for  $f$  in the environment. The effect handler is endowed with a fresh region and evidence and the substatements are translated recursively. As the implementation statement is just a value, it is translated trivially, and we obtain

```

val x = f[rc; nf](1);
try { [r2; n2 : r2 ⊆ rc](yield2) ⇒ S[[ f(x) ]]E' }
with { (j, k) ⇒ 42 }

```

The lifting environment is adapted for the new region we have entered and a new entry for the capability is added. Importantly, the evidence for the existing binding for  $f$  is adapted and we find

$$\begin{aligned} E'' &= E' \oplus n_2, \text{yield}_2 \mapsto (r_2 \mid \emptyset) \\ &= \{ r_2, (f, \text{yield}_2), (f \mapsto (r_f \mid n_2 \oplus n_f), \text{yield}_2 \mapsto (r_2 \mid \emptyset)), (\Gamma_{E'}, r_2, n_2 : r_2 \sqsubseteq r_c) \} \end{aligned}$$

The translation of the second application of  $f$  is again a simple matter of looking up the current region and the evidence for  $f$  in the environment. As the latter is kept sound during the translation, the evidence is exactly right,

$$f[r_2; n_2 \oplus n_f](x)$$

### 3.5 Properties of Lift Inference

For the lift-inference translation to be sensible it should be typability- and semantics-preserving.

*Typability preservation.* For the proof of typability preservation we make heavy use of the soundness of the lifting environment. To do so, we need the following lemma stating that soundness is maintained during the translation.

LEMMA 3.5 (SOUNDNESS OF ENVIRONMENTS FOR THE LIFT-INFERENCING TRANSLATION).

*All adaptations of lifting environments made by the lift-inference translation take sound environments to sound environments.*

This enables the theorem that the translation takes well-typed terms in System  $\Xi$  to well-typed terms in  $\Lambda_{\text{cap}}$ .

THEOREM 3.6 (TYPABILITY PRESERVATION FOR THE LIFT-INFERENCING TRANSLATION).

For  $E = \{ \rho, \Delta, m, \Gamma_E \}$  sound,

if  $\Gamma \mid \Delta \vdash s : \tau$ , then  $\Gamma_E, \Gamma, \mathcal{T}[\Delta]^E \mid \rho \vdash \mathcal{S}[s]^E : \tau$ ;

if  $\Gamma \mid \Delta \vdash b : \sigma$ , then  $\Gamma_E, \Gamma, \mathcal{T}[\Delta]^E \vdash B[b]^E : \mathcal{T}[\sigma]_{\mathcal{R}[b]}^E$ ;

if  $\Gamma \vdash v : \tau$ , then  $\Gamma_E, \Gamma \vdash v : \tau$ .

Since the empty environment  $\emptyset = \{ \top, \emptyset, \emptyset, \emptyset \}$  is trivially sound, Theorem 3.6 implies typability preservation for the translation of closed terms starting with the empty environment.

*Semantics preservation.* For semantics preservation note that the operational semantics of both calculi mainly differs in the presence of regions and evidence in  $\Lambda_{\text{cap}}$ . As noted before, these are irrelevant for the operational semantics and can thus be erased. Then the only difference is that there is no rule for reduction of function definitions in  $\Lambda_{\text{cap}}$ . It is replaced by two consecutive rules. Hence, we obtain the following result.



**THEOREM 3.7 (EVALUATION FOR THE LIFT-INFERENCE TRANSLATION).**

If  $\emptyset \mid \emptyset \vdash s : \tau$  and  $\langle s \parallel \#_{start} :: \bullet \rangle \rightarrow^{n+k} \langle \mathbf{return} \ v \parallel \bullet \rangle$ ,

then  $\langle \mathcal{S} \llbracket s \rrbracket^0 \parallel \#_{start} :: \bullet \rangle \rightarrow^{n+2k} \langle \mathbf{return} \ v \parallel \bullet \rangle$ ,

where  $k$  is the number of steps for function definitions and  $n$  the number of other steps in System  $\Xi$ .

Note that together with the corresponding results of the papers we build on (see the overview in Figure 1), the above theorems guarantee typability and semantics preservation along the whole compilation pipeline down to System F, fully proven. Hence, a well-typed term in Effekt is guaranteed to have the same semantics after translation to System F. In particular, type safety implies effect safety and thus guarantees that no effect goes unhandled.

## 4 EVALUATION

To evaluate our approach, we have implemented lift inference for the Effekt language. This way, we could close the gap illustrated in Section 1 and write benchmark programs directly in Effekt. We first describe the implementation and some limitations, before we discuss the benchmark results.

### 4.1 Implementation

Effekt is a functional language with support for lexical effect handlers, effect inference, data types, type polymorphism, interface types that generalize functions, backtrackable local state, and many more features. The implementation of the compiler amounts to around 23k lines of code in Scala. For this paper, we have extended the Effekt compiler in two ways. We have implemented the lift-inference translation presented in Subsection 3.4, and we have implemented a new backend targeting SML in continuation-passing style.

**4.1.1 Lift Inference.** The overview in Figure 1 doubles as an overview over the compiler phases in the Effekt compiler. A translation of the source language to explicit capability-passing System  $\Xi$ , which is called *Core* in the implementation, was already implemented by Brachthäuser et al. [2020]. To evaluate the feasibility of the translation described in this paper, we have added a new intermediate representation that corresponds to  $\Lambda_{cap}$ , which is called *Lifted* in the implementation. Lifted is typed and includes explicit subregion evidence, but regions are erased from the type level.

Both Core and Lifted differ from the presentation in this paper in that they also support various other features of the language, such as data types, pattern matching, local mutable state, and more. Like Leijen [2017], the Effekt compiler also distinguishes potentially effectful expressions from pure expressions, in order to generate more efficient code. The implementation of lift inference itself is a straightforward translation of the algorithm presented in Section 3.4 to Scala. The lifting environment  $E$  is implemented as a `Map[Symbol, List[Lift]]` mapping block variables (*i.e.*, `Symbols`) to a chain of evidence variables (*i.e.*, `Lifts`) that witness the subregion relationship.

**4.1.2 SML Backend in CPS.** The translation from  $\Lambda_{cap}$  to System F in iterated continuation-passing style is conceptually described by Schuster et al. [2022b]. However, they do not present an implementation. As a second implementation step for this paper, we have thus implemented a translation from Lifted to Standard ML [Milner et al. 1997] in continuation-passing style. Specifically, we target MLton since it is a whole program optimizing compiler. We conjectured that MLton could discover many of the static abstractions identified by Schuster et al. [2020] at compile-time and thus heavily optimize the generated programs in CPS. We can, however, imagine that a setting with separate compilation could profit from lift inference as well. For example, our approach might be applied within a compilation unit and one might moreover rely on link-time optimizations or just-in-time compilation to obtain further improvements from lift inference at runtime. We leave closer investigation of this to future work.

While the translation conceptually translates effectful programs to pure System F, in the toplevel region our implementation supports all native effects present in the target language, like native mutable references, file IO, etc.

Our CPS translation employs standard techniques to avoid administrative  $\beta$ - and  $\eta$ -redexes [Danvy and Filinski 1992; Schuster and Brachthäuser 2018], is curried [Hillerström et al. 2017; Schuster et al. 2022b], but not fully iterated [Schuster and Brachthäuser 2018; Schuster et al. 2020], which is to say that we do not abstract more than one continuation in function definitions. Rather, additional continuation parameters are added by instantiating the answer type with another layer of CPS as required.

Effekt has higher-rank polymorphism originating from function parameters as well as from polymorphic effect signatures (e.g., `effect Exc { def raise[A](): A }`). Since SML is a language with a Hindley-Milner-style [Hindley 1969; Milner 1978] type system, it does not support higher-rank polymorphism. Due to this limitation, our SML backend currently does neither support higher-rank function types nor polymorphic effect signatures.

Another problem is that the translation from  $\Lambda_{\text{cap}}$  to System F presented by Schuster et al. [2022b] makes heavy use of higher-rank types. Region abstractions are translated to type abstractions which makes region-polymorphic function parameters have a higher-rank type. Moreover, the type of subregion evidence  $\rho_1 \sqsubseteq \rho_2$  is translated to  $\forall a. \text{CPS } \mathcal{T}[\rho_2] a \rightarrow \text{CPS } \mathcal{T}[\rho_1] a$ . Functions that take evidence parameters, which are virtually all translated functions, have rank-2 type. While this might sound like a severe limitation, we observed that in many functions evidence is not actually used but only passed on and thus can often be treated parametrically. When evidence is actually used, its type argument is often inferred monomorphically. However, using the same evidence at two different types in HM will lead to a type mismatch, a problem not present in System F.

In order to admit more Effekt programs to be translated to SML, we perform evidence monomorphization, where we effectively partially evaluate programs with respect to evidence and specialize effect handler implementations to the region they are called in. We consider both, the limitation of SML not supporting higher-rank types, as well as the implemented evidence monomorphization as non-essential aspects of the present paper. We could have chosen an arbitrary different target platform that does support higher-rank types.

## 4.2 Benchmarks

One of our goals was to reproduce the performance results of Schuster et al. [2020] in a realistic source-level language, in particular, their conjecture that specialized optimizations or special reduction theories are not needed to remove abstraction overhead; rather, existing optimizing compilers can do the job. In Table 1 we present the results of measuring the running time of programs written in Effekt and compiled to our SML backend against the running time of the same programs written in other languages with effect handlers. The benchmark programs are taken from a community benchmark suite that has been designed specifically for effect handler implementations [Hillerström et al. 2023]. The repository contains detailed explanations for each of the benchmark programs. Benchmarks were conducted on a 12th Gen Intel(R) Core(TM) i7-1255U running Ubuntu 22.04.

**4.2.1 Systems.** We compare Eff, Multicore OCaml, and Koka with our own implementation in Effekt. Eff [Karachalias et al. 2021; Saleh et al. 2018] is a language with effect handlers and specific optimizations for those. After optimization it generates OCaml code which is then compiled with `ocamlopt 4.14.1`. Unfortunately, investigation of the generated code reveals that in many programs the specialization of functions to handlers is not triggered. Multicore OCaml [Dolan et al. 2014] 4.12.0 is a language with effect handlers and a very fast runtime. While it does not officially support

Table 1. Benchmark results comparing Eff, OCaml, Koka, our implementation of lift inference in Effekt, and a hand-optimized baseline. Fastest mean for each benchmark is highlighted in gray.

Benchmark	Mean time in ms (standard deviation)				
	Eff	OCaml	Koka	Effekt	Baseline
Countdown (200M)	72.0 ( $\pm 13.2$ )	1976.1 ( $\pm 26.6$ )	1598.0 ( $\pm 24.0$ )	44.5 ( $\pm 1.0$ )	44.5 ( $\pm 0.9$ )
Fibonacci (42)	1093.6 ( $\pm 5.5$ )	1161.5 ( $\pm 12.0$ )	1222.6 ( $\pm 27.0$ )	1335.4 ( $\pm 12.0$ )	1335.4 ( $\pm 24.5$ )
Product Early (100k)	535.7 ( $\pm 71.9$ )	113.0 ( $\pm 0.4$ )	1506.6 ( $\pm 20.0$ )	238.2 ( $\pm 33.4$ )	113.0 ( $\pm 0.9$ )
Iterator (40M)	516.3 ( $\pm 17.6$ )	195.4 ( $\pm 1.3$ )	1082.0 ( $\pm 9.6$ )	92.5 ( $\pm 10.7$ )	13.7 ( $\pm 0.5$ )
Queens (12)	262.2 ( $\pm 6.1$ )	635.6 ( $\pm 1.8$ )	2643.5 ( $\pm 26.6$ )	117.2 ( $\pm 0.3$ )	96.6 ( $\pm 1.0$ )
Tree Explore (16)	161.1 ( $\pm 3.8$ )	142.9 ( $\pm 2.2$ )	278.4 ( $\pm 5.1$ )	187.1 ( $\pm 4.4$ )	179.1 ( $\pm 2.0$ )
Triples (300)	125.0 ( $\pm 4.4$ )	315.5 ( $\pm 3.3$ )	2635.8 ( $\pm 11.4$ )	30.0 ( $\pm 0.5$ )	25.1 ( $\pm 0.4$ )
Resume Non-tail (10k)	182.4 ( $\pm 15.9$ )	190.4 ( $\pm 1.0$ )	1601.5 ( $\pm 16.6$ )	85.9 ( $\pm 3.5$ )	62.5 ( $\pm 3.0$ )
Parsing (20k)	2061.7 ( $\pm 177.3$ )	1443.5 ( $\pm 14.6$ )	3220.4 ( $\pm 253.6$ )	88.6 ( $\pm 0.8$ )	88.1 ( $\pm 1.0$ )

multiple resumptions, which some benchmarks use, it has limited support for those which is sufficient to run these benchmarks. Koka [Xie and Leijen 2021] is a language with effect handlers, a fast runtime, and an optimizing compiler. The Koka compiler generates C code which then is further compiled with gcc. Our own Effekt compiler produces code in Standard ML and then uses MLton 20210117 to compile it. We instruct MLton to choose numbers to be 64bit integers to match the behavior of the other languages. Finally, as a baseline, we have taken the programs produced by our Effekt compiler and minimized and hand-optimized them using native effects where possible.

**4.2.2 Results.** Our findings (presented in Table 1) are generally positive. Effekt outperforms the other languages in most benchmarks, sometimes by an order of magnitude. Speedups range from around 1.6x–16.3x to the next best system for each particular benchmark. On the other side, we only observe slowdowns of 1.2x–2.1x compared to the best system for the specific benchmark. Our benchmarks are available as an artifact (see Section 7).

The Countdown benchmark uses the state effect to tail-recursively count down from a given number. Some implementations (OCaml and Koka) use references to implement the state effect, others (such as Eff and ours after evidence monomorphization) modify the answer type to be a function taking the state. In OCaml and Koka getting and setting the state goes through performing an effect operation, while Eff and Effekt are able to optimize this indirection away.

The Fibonacci benchmark does not actually use effect handlers. Eff, Koka, and Effekt generate special code for pure functions and the performance is competitive. The code generated by Eff and Effekt (after MLton optimizations) is very similar to the handwritten direct-style OCaml code and runtime differences amount to the different language runtimes used to execute the code.

The Product Early benchmark pushes 1,000 frames onto the stack and then discards all of them by throwing an exception. We can see a slowdown compared to native exceptions in OCaml and in the MLton baseline. This is due to the fact that the implementations of exceptions in both runtimes are very efficient and indeed faster than our approach of translating to CPS.

The Iterator benchmark models push streams by using an effect to emit values. The handler uses state to add all values and calls the continuation in tail position. We can see speedups compared to OCaml of around 2.1x. The optimizations performed by Eff seem to be blocked and thus the handler is not optimized away. We expect this problem to be technical and not fundamental in their approach.

The Queens benchmark searches for a solution to place  $n$  queens on a chessboard. It heavily uses continuations in a non-trivial way to perform backtracking search. We can observe a speedup of 2.2x over Eff, but note again that their optimizations seem to be blocked.

The Tree Explore benchmark constructs a tree and then traverses it to collect all leaves. It uses a choice effect to simulate a non-deterministic traversal. We can see a slowdown of 1.3x compared to OCaml. The reason for this is again the difference between OCaml and MLton. Indeed, we have translated the code we generate from Standard ML to OCaml and observed that it runs faster than the OCaml variant using native effects.

The Triples benchmark makes heavy use of continuations to perform a backtracking search. We see speedups of around 4.2x compared to Eff, but yet again note that the rewrite rules of Karachalias et al. [2021] seem not to be applied fully.

The Resume Non-tail benchmark calls an effect operation in a loop. The handler resumes in non-tail position and thus aggregates  $N$  stack frames. After the loop returns, the stack frames are popped one-after-another. Again, we can see speedups of 2.1x over Eff, where the rewrites are not fully applied.

Finally, the Parsing benchmark defines a streaming parser which uses three effects:

```
def parse(a: Int): Unit / {Read, Emit, Stop} = ...
```

The `Read` effect reads a character, the `Emit` effect emits the result of parsing a line, and the `Stop` effect stops when an unrecognized character is found. The function `parse` is used under three handlers, one for each of the effects:

```
sum { catch { feed(n) { parse(0) } } }
```

The program has non-trivial control flow, which is abstracted away by the use of effect handlers. Moreover we could use the same function `parse` with a different source of characters and a different target of emitted values. Our Effekt implementation significantly outperforms the other languages by a factor of 16.3x–36.3x. It is the only benchmark that relies on evidence monomorphism in order to compile. Our implementation specializes this function to the handlers surrounding it, which after optimization results in a single tight loop, which is exactly our original goal: to remove all abstraction overhead introduced by using effect handlers.

In general, our approach works better in the cases where effects and resumptions are used extensively. In these cases we observe large speedups over the other implementations. That said, the optimizations for Eff were often blocked in these benchmarks. We would expect the results for Eff to be much closer to ours, if the optimizations kick in. Our results are often quite close to the hand-optimized baseline. In these cases our implementation, of which lift-inference is an integral part, is able to remove all abstraction introduced by the use of effect handlers to structure the program. In the other cases, more investigation is needed in order to remove the gap between compiled code using effect handlers and hand-optimized code using native effects.

## 5 RELATED WORK

We have presented a translation from System  $\Xi$ , a calculus with second-class capabilities to  $\Lambda_{\text{cap}}$ , a calculus with region-based effects. In combination with a translation to iterated CPS this enables efficient compilation of effect handlers. In this section, we compare our approach to existing work.

### 5.1 Efficient Compilation of Effect Handlers

Closely related is the work on explicit effect subtyping for algebraic effect handlers in Eff [Saleh et al. 2018]. Their main motivation is to use this explicit information in the optimization of programs using effect handlers [Karachalias et al. 2021]. In particular, they define source-to-source rewrite

rules on an intermediate representation called  $\text{ExEff}$ . The rewrite rules are designed to propagate handlers down until they reach an effect operation in which case the effect operation can be statically reduced. Intermediate frames are aggregated in the return clause of the handler. While our motivation is ultimately the same, our work is in the context of *lexical* effect handlers as they appear in Effekt. Also, we do not define a reduction theory on a language with effects and handlers, but instead (via  $\Lambda_{\text{cap}}$ ) define a translation to System F in CPS. This way, existing optimizing compilers for functional languages (such as SML) can readily be used. Karachalias et al. [2021] support first-class functions which makes their language more general than ours. While their explicitly typed language applies subtyping coercions to arbitrary computations, we pass evidence values along and only use them at effect operations where they are needed.

Also highly related in this regard is the work on evidence passing [Xie et al. 2020; Xie and Leijen 2021], which provides the basis for the implementation of effect handlers in the Koka language. The idea is to pass evidence vectors down to effect operations. These evidence vectors consist of pairs of labels and handler implementations so that handlers can be looked up in place. In contrast, evidence in  $\Lambda_{\text{cap}}$  is just a list of labels and the handler implementations are passed as capabilities. Evidence passing is defined in the context of *dynamically* scoped handlers and hence does not reflect the *lexical* nesting of handlers as regions and subregion evidence in  $\Lambda_{\text{cap}}$  do. Xie et al. [2020] define an evidence passing translation from an algebraic effect calculus to an evidence calculus, thus determining the evidence vectors statically. This translation is facilitated by the effect system of the algebraic effect calculus based on rows of effect labels. In contrast, our source calculus System  $\Xi$  does not feature a visible effect system and instead relies on second-class capabilities. Xie and Leijen [2021] do not define a translation but achieve evidence passing by defining appropriate evaluation rules for the algebraic effect calculus, hence their evidence vectors are created dynamically. This allows them to lift the restriction of scoped resumptions, imposed by Xie et al. [2020]. Both approaches support first-class functions. Moreover, both papers define a translation to System  $F^v$ , a polymorphic lambda calculus with support for multi-prompt monads. Instead, by building on Schuster et al. [2022b], we directly compile to System F in CPS.

## 5.2 Languages with Second-Class Values

Our lift inference consumes programs in System  $\Xi$  [Brachthäuser et al. 2020], a language with second-class functions and capabilities. It is inspired by the work of Oswald et al. [2016] who present  $\lambda^{1/2}$  a lambda calculus that features both first-class and second-class functions, but no control effects. Their work in turn builds on type-based escape analysis [Hannan 1998]. In  $\lambda^{1/2}$ , second-class functions cannot be returned, nor closed over by first-class functions. In contrast, System  $\Xi$  does not support first-class functions and in consequence our translation does not have to handle them. We do not expect any complications in extending System  $\Xi$  and our translation to first-class functions in the style of  $\lambda^{1/2}$ —that is, to first-class functions that cannot close over second-class functions and capabilities. As they do not contain capabilities that need to satisfy some subregion constraint, they can always run in the toplevel region, so we could just pretend that one to be their definition-site region.

Brachthäuser et al. [2022] present System  $C$  as an extension of System  $\Xi$  to support a fine-grained notion of second-class values. Their calculus introduces explicit box and unbox constructs, inspired by modal logics. They also extend the type system to track which capabilities are used by a statement or block. Boxing takes a second-class block and turns it into a first-class value, where the type of the boxed block specifies the necessary capabilities (e.g.,  $\text{Int} \Rightarrow \text{Int at } \{\text{yield}\}$ ). To call a boxed function it needs to be unboxed first. When unboxing, the type system ensures that the necessary capabilities are still available, preventing functions from closing over capabilities



and then leaving the scope of a handler. This system is more expressive than  $\lambda^{1/2}$  and it is less clear to us how to translate the sets of capabilities (e.g., `{yield}`) to the corresponding region. We leave studying the translation of System C to  $\Lambda_{\text{cap}}$  to future work.

Xhebraj et al. [2022] present another variant of  $\lambda^{1/2}$ , called  $\lambda_{\leftrightarrow}^{1/2}$ , which supports returning second-class functions and is designed to be used for stack allocating memory. Safety is achieved by modifying the runtime semantics: When a second-class value is returned, the returning frame is simply not removed. While this is an elegant solution, our goal is to target standard runtime systems like System F.

### 5.3 Languages with Regions and Subregioning

Our lift inference produces programs in  $\Lambda_{\text{cap}}$  [Schuster et al. 2022b] featuring explicit regions and subregion evidence. We use a generalization which allows for non-scoped continuations. While we follow Schuster et al. [2022b] and use regions and evidence to track the lexical nesting of handlers on the stack, the original usage of regions is in memory management [Tofte and Talpin 1997] and more generally resource management. Our notion of regions could in principle also be used for resource management. A recent calculus in this regard which is close to  $\Lambda_{\text{cap}}$  is presented by Schuster et al. [2022a]. They use regions and subregion evidence to deal with the management of resources in the presence of exception handlers. In contrast to this work, they do not deal with general effect handlers and do not consider inference of regions and evidence. The work of Schuster et al. [2022a] is based on Kiselyov and Shan [2008], who also perform region inference. Their approach, however, is very different from ours, as they encode regions using monad transformers and hence rely on type inference to infer regions. Likewise, other algorithms [Tofte et al. 2001; Tofte and Talpin 1997] for region inference in the context of memory management are different from ours, often creating fresh regions for each variable and subsequently analysing which of them can be unified. In contrast to this prior work, we infer regions by establishing a connection between the lexical scoping of second-class values and regions.

It might be especially interesting to consider our CPS translation with a target language which already has a notion of regions, like the ML Kit [Tofte et al. 2001], and potentially try to preserve region annotations. We leave exploration of this to future work.

### 5.4 Dictionary Passing and Monad Polymorphism in Haskell

In Haskell it is typical to use stacks of monad transformers [Liang et al. 1995] and type classes to compose different programs using different effects into one [Jones 1995]. When type classes are implemented by dictionary passing, this is not unlike our passing of capabilities, but implicit. When effect operations corresponding to a lower layer in the monad-transformer stack are used, they have to be lifted through all layers above, just like in our work. Finding the correct composition of lifts is automatic and works by type class resolution guided by the type of computations. In contrast to this, we assume that capabilities are passed explicitly and find the correct composition of evidence by a transformation guided by program terms. This has the advantage that different instances of the same effect are easily disambiguated by passing different capabilities.

This problem was observed by Figueroa et al. [2015] and named effect interference. For them, the interference between different instances of the same effect is also a security concern. As a solution they propose the explicit passing of capabilities. However, their notion of capability is different from ours in that they use them to ensure certain security guarantees on top of monad transformers while we use them as an implementation technique for effect handlers. Consequently, our capabilities contain the operation to execute when they are used while theirs do not.



Schrijvers et al. [2019] compare and contrast monad transformers and the traditional implementation of dynamic effect handlers in terms of a free monad in Haskell. There, the problem of lift inference manifests in a different way. Effectful programs are written against an open union of signatures [Kiselyov and Ishii 2015]. The challenge is for a given effect operation to find the correct injection into this open union. Again, Haskell type classes can be used for this to some extent.

However, using type class resolution to find the correct handler is problematic when there are multiple instances of the same effect in the same program. As a solution, Devriese [2019] propose explicit passing of type class dictionaries, which essentially are what we call capabilities. Moreover, in order to nest handlers, they propose the explicit use of liftings, essentially what we call evidence. They also present a case study that speaks for the feasibility of their explicit approach. We, however, infer the correct use of evidence, so programmers do not have to do so explicitly. In their setting liftings are general monad morphisms, while in our setting we specialize them to the continuation- and state monads. Another difference is that they apply these liftings to capabilities, while we pass evidence to the places where capabilities are used.

## 6 CONCLUSION

In this paper, we have presented a way to infer lifting information for lexical effects and handlers, by giving a typed translation from a calculus System  $\Xi$  with second-class capabilities to a calculus  $\Lambda_{\text{cap}}$  with explicit regions and subregion evidence. Our translation preserves typability and semantics. It makes use of the second-class property to provide a clear connection for the definition-site and each call-site of a function. This establishes a precise relation between reasoning based on the second-class property and region-based regioning.

Moreover, we have evaluated the implications of lift inference practically, by implementing it as a compiler phase for a source-level language. To this end, we have further implemented the CPS translation for  $\Lambda_{\text{cap}}$  described by Schuster et al. [2022b], which makes heavy use of the information provided by lift inference to enable efficient compilation of effect handlers. Our benchmarks indicate that our approach is competitive with other state-of-the-art implementations of effect handlers and often outperforms them.

While the second-class property of our source calculus is particularly helpful for lift inference, it can sometimes be a restriction in programming. In the future, it would be interesting to investigate whether it is possible to extend our approach to a language which has a controlled way of using effectful first-class functions, as described, e.g., by Brachthäuser et al. [2022]. Furthermore, it would be interesting to see how lift inference can be performed for languages with traditional effect handlers. Also, while System  $\Xi$  uses types to enforce the second-class property for capabilities, we believe that any mechanism to enforce this would do. But we leave a fully precise exploration of this for future work.

## 7 DATA-AVAILABILITY STATEMENT

The benchmarks from Section 4 are available to be run in a Docker container in the accompanying artifact [Müller et al. 2023a].

## ACKNOWLEDGMENTS

We thank Matthew Fluet for his helpful assistance with MLton. The work on this project was supported by the Deutsche Forschungsgemeinschaft (DFG – German Research Foundation) – project number DFG-448316946.

## REFERENCES

- Andrej Bauer and Matija Pretnar. 2015. Programming with algebraic effects and handlers. *Journal of Logical and Algebraic Methods in Programming* 84, 1 (2015), 108–123. <https://doi.org/10.1016/j.jlamp.2014.02.001>
- Dariusz Biernacki, Maciej Piróg, Piotr Polesiuk, and Filip Sieczkowski. 2017. Handle with Care: Relational Interpretation of Algebraic Effects and Handlers. *Proc. ACM Program. Lang.* 2, POPL, Article 8 (Dec. 2017), 30 pages. <https://doi.org/10.1145/3158096>
- Dariusz Biernacki, Maciej Piróg, Piotr Polesiuk, and Filip Sieczkowski. 2019. Binders by Day, Labels by Night: Effect Instances via Lexically Scoped Handlers. *Proc. ACM Program. Lang.* 4, POPL, Article 48 (Dec. 2019), 29 pages. <https://doi.org/10.1145/3371116>
- Jonathan Immanuel Brachthäuser, Philipp Schuster, Edward Lee, and Aleksander Boruch-Gruszecki. 2022. Effects, Capabilities, and Boxes: From Scope-Based Reasoning to Type-Based Reasoning and Back. *Proc. ACM Program. Lang.* 6, OOPSLA, Article 76 (apr 2022), 30 pages. <https://doi.org/10.1145/3527320>
- Jonathan Immanuel Brachthäuser, Philipp Schuster, and Klaus Ostermann. 2020. Effects as Capabilities: Effect Handlers and Lightweight Effect Polymorphism. *Proc. ACM Program. Lang.* 4, OOPSLA, Article 126 (Nov. 2020). <https://doi.org/10.1145/3428194>
- Oliver Danvy and Andrzej Filinski. 1992. Representing control: A study of the CPS transformation. *Mathematical Structures in Computer Science* 2, 4 (1992), 361–391. <https://doi.org/10.1017/S0960129500001535>
- Dominique Devriese. 2019. Modular Effects in Haskell through Effect Polymorphism and Explicit Dictionary Applications: A New Approach and the  $\mu$ VeriFast Verifier as a Case Study. In *Proceedings of the 12th ACM SIGPLAN International Symposium on Haskell (Berlin, Germany) (Haskell 2019)*. Association for Computing Machinery, New York, NY, USA, 1–14. <https://doi.org/10.1145/3331545.3342589>
- Stephen Dolan, Leo White, and Anil Madhavapeddy. 2014. Multicore OCaml. In *OCaml Workshop*.
- R. Kent Dybvig, Simon Peyton Jones, and Amr Sabry. 2007. A Monadic Framework for Delimited Continuations. *Journal of Functional Programming* 17, 6 (Nov. 2007), 687–730. <https://doi.org/10.1017/S0956796807006259>
- Kavon Farvardin and John Reppy. 2020. From Folklore to Fact: Comparing Implementations of Stacks and Continuations. In *Proceedings of the 41st ACM SIGPLAN Conference on Programming Language Design and Implementation (London, UK) (PLDI 2020)*. Association for Computing Machinery, New York, NY, USA, 75–90. <https://doi.org/10.1145/3385412.3385994>
- Ismael Figueroa, Nicolas Tabareau, and Éric Tanter. 2015. Effect capabilities for Haskell: Taming effect interference in monadic programming. *Science of Computer Programming* 119 (Nov 2015). <https://doi.org/10.1016/j.scico.2015.11.010>
- John Hannan. 1998. A Type-based Escape Analysis for Functional Languages. *Journal of Functional Programming* 8, 3 (May 1998), 239–273. <https://doi.org/10.1017/S0956796898003025>
- Daniel Hillerström, Sam Lindley, Bob Atkey, and KC Sivaramakrishnan. 2017. Continuation Passing Style for Effect Handlers. In *Formal Structures for Computation and Deduction (LIPIcs, Vol. 84)*. Schloss Dagstuhl–Leibniz-Zentrum für Informatik.
- Daniel Hillerström, Filip Koprivec, and Philipp Schuster (benchmarking chairs). 2023. Effect handlers benchmarks suite. (2023). <https://github.com/effect-handlers/effect-handlers-bench>
- Daniel Hillerström, Sam Lindley, and Robert Atkey. 2020. Effect handlers via generalised continuations. *Journal of Functional Programming* 30 (2020), e5. <https://doi.org/10.1017/S0956796820000040>
- J.R. Hindley. 1969. The principal type scheme of an object in combinatory logic. *Trans. of the American Mathematical Society* 146 (Dec. 1969), 29–60. <https://doi.org/10.2307/1995158>
- Mark P. Jones. 1995. Functional programming with overloading and higher-order polymorphism. In *Advanced Functional Programming*, Johan Jeuring and Erik Meijer (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 97–136. [https://doi.org/10.1007/3-540-59451-5\\_4](https://doi.org/10.1007/3-540-59451-5_4)
- Georgios Karachalias, Filip Koprivec, Matija Pretnar, and Tom Schrijvers. 2021. Efficient Compilation of Algebraic Effect Handlers. *Proc. ACM Program. Lang.* 5, OOPSLA, Article 102 (oct 2021), 28 pages. <https://doi.org/10.1145/3485479>
- Oleg Kiselyov and Hiromi Ishii. 2015. Freer Monads, More Extensible Effects. In *Proceedings of the Haskell Symposium (Vancouver, BC, Canada)*. ACM, New York, NY, USA, 94–105. <https://doi.org/10.1145/2887747.2804319>
- Oleg Kiselyov and Chung-chieh Shan. 2008. Lightweight Monadic Regions. In *Proceedings of the Haskell Symposium (Victoria, BC, Canada) (Haskell '08)*. ACM, New York, NY, USA. <https://doi.org/10.1145/1411286.1411288>
- Daan Leijen. 2017. Type directed compilation of row-typed algebraic effects. In *Proceedings of the Symposium on Principles of Programming Languages*. ACM, New York, NY, USA, 486–499. <https://doi.org/10.1145/3093333.3009872>
- Paul Blain Levy, John Power, and Hayo Thielecke. 2003. Modelling environments in call-by-value programming languages. *Information and Computation* 185, 2 (2003), 182–210. [https://doi.org/10.1016/S0890-5401\(03\)00088-9](https://doi.org/10.1016/S0890-5401(03)00088-9)
- Sheng Liang, Paul Hudak, and Mark Jones. 1995. Monad Transformers and Modular Interpreters. In *Proceedings of the Symposium on Principles of Programming Languages (San Francisco, California, USA)*. ACM, New York, NY, USA, 333–343. <https://doi.org/10.1145/199448.199528>
- Sam Lindley, Conor McBride, and Craig McLaughlin. 2017. Do Be Do Be Do. In *Proceedings of the Symposium on Principles of Programming Languages (Paris, France)*. ACM, New York, NY, USA, 500–514. <https://doi.org/10.1145/3009837.3009897>

- Robin Milner. 1978. A theory of type polymorphism in programming. *J. Comput. System Sci.* 17 (1978), 248–375. [https://doi.org/10.1016/0022-0000\(78\)90014-4](https://doi.org/10.1016/0022-0000(78)90014-4)
- Robin Milner, Mads Tofte, Robert Harper, and David MacQueen. 1997. *The definition of standard ML: revised*. MIT press. <https://doi.org/10.7551/mitpress/2319.001.0001>
- Marius Müller, Philipp Schuster, Jonathan Lindegaard Starup, Klaus Ostermann, and Jonathan Immanuel Brachthäuser. 2023a. *Artifact of the paper 'From Capabilities to Regions: Enabling Efficient Compilation of Lexical Effect Handlers'*. <https://doi.org/10.5281/zenodo.8315298>
- Marius Müller, Philipp Schuster, Jonathan Lindegaard Starup, Klaus Ostermann, and Jonathan Immanuel Brachthäuser. 2023b. *From Capabilities to Regions: Enabling Efficient Compilation of Lexical Effect Handlers*. Extended Technical Report. University of Tübingen, Germany. <https://se.informatik.uni-tuebingen.de/publications/mueller23lift>.
- Leo Osvald, Grégory Essertel, Xilun Wu, Lilliam I González Alayón, and Tiark Rompf. 2016. Gentrification gone too far? affordable 2nd-class values for fun and (co-) effect. In *Proceedings of the Conference on Object-Oriented Programming, Systems, Languages and Applications*. ACM, New York, NY, USA, 234–251. <https://doi.org/10.1145/3022671.2984009>
- Gordon Plotkin and Matija Pretnar. 2009. Handlers of algebraic effects. In *European Symposium on Programming*. Springer-Verlag, 80–94. [https://doi.org/10.1007/978-3-642-00590-9\\_7](https://doi.org/10.1007/978-3-642-00590-9_7)
- Gordon D. Plotkin and Matija Pretnar. 2013. Handling Algebraic Effects. *Logical Methods in Computer Science* 9, 4 (2013). [https://doi.org/10.2168/LMCS-9\(4:23\)2013](https://doi.org/10.2168/LMCS-9(4:23)2013)
- Matija Pretnar, Amr Hany Shehata Saleh, Axel Faes, and Tom Schrijvers. 2017. *Efficient compilation of algebraic effects and handlers*. Technical Report. Department of Computer Science, KU Leuven; Leuven, Belgium.
- Amr Hany Saleh, Georgios Karachalias, Matija Pretnar, and Tom Schrijvers. 2018. Explicit Effect Subtyping. In *Programming Languages and Systems*, Amal Ahmed (Ed.). Springer International Publishing, Cham, Switzerland, 327–354. [https://doi.org/10.1007/978-3-319-89884-1\\_12](https://doi.org/10.1007/978-3-319-89884-1_12)
- Tom Schrijvers, Maciej Piróg, Nicolas Wu, and Mauro Jaskieloff. 2019. Monad Transformers and Modular Algebraic Effects: What Binds Them Together. In *Proceedings of the 12th ACM SIGPLAN International Symposium on Haskell (Berlin, Germany) (Haskell 2019)*. Association for Computing Machinery, New York, NY, USA, 98–113. <https://doi.org/10.1145/3331545.3342595>
- Philipp Schuster and Jonathan Immanuel Brachthäuser. 2018. Typing, Representing, and Abstracting Control. In *Proceedings of the Workshop on Type-Driven Development (St. Louis, Missouri, USA)*. ACM, New York, NY, USA, 14–24. <https://doi.org/10.1145/3240719.3241788>
- Philipp Schuster, Jonathan Immanuel Brachthäuser, Marius Müller, and Klaus Ostermann. 2022b. A Typed Continuation-Passing Translation for Lexical Effect Handlers. In *Proceedings of the 43rd ACM SIGPLAN International Conference on Programming Language Design and Implementation (San Diego, CA, USA) (PLDI 2022)*. Association for Computing Machinery, New York, NY, USA, 566–579. <https://doi.org/10.1145/3519939.3523710>
- Philipp Schuster, Jonathan Immanuel Brachthäuser, and Klaus Ostermann. 2020. Compiling Effect Handlers in Capability-Passing Style. *Proc. ACM Program. Lang.* 4, ICFP, Article 93 (Aug. 2020), 28 pages. <https://doi.org/10.1145/3408975>
- Philipp Schuster, Jonathan Immanuel Brachthäuser, and Klaus Ostermann. 2022a. Region-based Resource Management and Lexical Exception Handlers in Continuation-Passing Style. In *Programming Languages and Systems*, Ilya Sergey (Ed.). Springer International Publishing, Cham, 492–519. [https://doi.org/10.1007/978-3-030-99336-8\\_18](https://doi.org/10.1007/978-3-030-99336-8_18)
- KC Sivaramakrishnan, Stephen Dolan, Leo White, Tom Kelly, Sadiq Jaffer, and Anil Madhavapeddy. 2021. Retrofitting Effect Handlers onto OCaml. In *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation (PLDI 2021)*. Association for Computing Machinery, New York, NY, USA, 206–221. <https://doi.org/10.1145/3453483.3454039>
- Hayo Thielecke. 2003. From Control Effects to Typed Continuation Passing. In *Proceedings of the 30th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (New Orleans, Louisiana, USA) (POPL '03)*. Association for Computing Machinery, New York, NY, USA, 139–149. <https://doi.org/10.1145/604131.604144>
- Mads Tofte, Lars Birkedal, Martin Elsman, Niels Hallenberg, and Peter Sestoft. 2001. Programming with Regions in the ML Kit (for Version 4). (10 2001).
- Mads Tofte and Jean-Pierre Talpin. 1997. Region-Based Memory Management. *Inf. Comput.* 132, 2 (Feb. 1997), 109–176. <https://doi.org/10.1006/inco.1996.2613>
- Anxhelo Xhebraj, Oliver Bračevac, Guannan Wei, and Tiark Rompf. 2022. What If We Don't Pop the Stack? The Return of 2nd-Class Values. In *36th European Conference on Object-Oriented Programming (ECOOP 2022) (Leibniz International Proceedings in Informatics (LIPIcs), Vol. 222)*, Karim Ali and Jan Vitek (Eds.). Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl, Germany, 15:1–15:29. <https://doi.org/10.4230/LIPIcs.ECOOP.2022.15>
- Ningning Xie, Jonathan Immanuel Brachthäuser, Daniel Hillerström, Philipp Schuster, and Daan Leijen. 2020. Effect Handlers, Evidently. *Proc. ACM Program. Lang.* 4, ICFP, Article 99 (Aug. 2020), 29 pages. <https://doi.org/10.1145/3408981>
- Ningning Xie and Daan Leijen. 2021. Generalized Evidence Passing for Effect Handlers: Efficient Compilation of Effect Handlers to C. *Proc. ACM Program. Lang.* 5, ICFP, Article 71 (aug 2021), 30 pages. <https://doi.org/10.1145/3473576>

Yizhou Zhang and Andrew C. Myers. 2019. Abstraction-safe Effect Handlers via Tunneling. *Proc. ACM Program. Lang.* 3, POPL, Article 5 (Jan. 2019), 29 pages. <https://doi.org/10.1145/3290318>

Received 2023-04-14; accepted 2023-08-27