

MPhil Thesis

Michael Jing Long Lee

May 5, 2025

Contents

1	Introduction	4
1.1	Contributions	6
2	Background	7
2.1	Metaprogramming	7
2.1.1	Metaprogramming for Fast and Maintainable Code . .	7
2.1.2	The Design Space of Metalanguages	12
2.2	Effect Handlers	13
2.2.1	Composable and Customisable Effects	13
2.2.2	λ_{op} : A Calculus for Effect Handlers	16
2.2.3	The Design Space of Effect Handlers	24
2.3	Scope Extrusion	25
2.3.1	Existing Solutions to the Scope Extrusion Problem . . .	27

Abstract

There are many different ways to manage the interaction between compile-time metaprogramming and effect handlers, but there has been no way to evaluate them against each other. This dissertation introduces such a mechanism, focusing on correctness, expressiveness, and efficiency. Second, we evaluate existing approaches by the aforementioned mechanism, illustrating the trade-offs in the design space. Finally, we introduce a novel approach, a dynamic Best-Effort check, that we show is correct, maximally expressive, and likely more efficient in the common case, as compared to existing checks.

1 Introduction

The promise of an optimising compiler is that the programmer can focus on writing maintainable code, entrusting the responsibility of *optimising* said code to the compiler.

Entrusting optimisation to the compiler is sufficient for many cases, but not all. For reasons that are both theoretical [16] and practical [17], compilers will never be able to identify *all* opportunities for optimisation. What can be done when compiler optimisation is not enough?

Metaprogramming (for example, C++ templates) [1] is one possible solution. In languages that support metaprogramming, programmers can identify what should be executed at compile-time. Importantly, programmers may write **maintainable** code, which when executed by the compiler, generates **efficient** code. For example, assume two implementations of the same function, one that runs faster on Arm, and one on Intel.

The maintainable way to switch between these functions is via a conditional:

```
1 let f () = if (cpu_type == Arm) then
2             (* return this program *)
3             else
4             (* return this program *)
5 f ()
```

OCaml

However, it might be expensive to determine processor type at run-time. If efficiency is paramount, programmers might have to duplicate their code, doubling maintenance costs. With metaprogramming, however, we may write

```
1 (* macro and $ can be read as annotations which
2    indicate that the code should be run at compile-time *)
3 macro m () = if (cpu_type == Arm) then
4             (* generate this program, to be run later *)
5             else
6             (* gen this program, to be run later *)
7 $(m ())
```

MacroCaml

Importantly, the macro `m` is executed by the compiler, *at compile time*. The Arm and Intel variants are *generated from a single specification*, not manually duplicated.

Importantly, to support metaprogramming, languages have to provide the

ability to build programs *to be run later*. These “suspended programs” are best thought of as abstract syntax trees. The following code constructs the program $(\lambda x.x)1$.

```

1  if (cpu_type == Arm) then
2    let build_app (f: (int -> int) expr) (v: int expr) = App(f, v)
3    let id_ast: (int -> int) expr = Lam(Var(x), Var(x))
4    build_app(id_ast, Int(1)) (* int expr *)
5  else
6    (*...*)

```

OCaml

Effect handlers are a powerful language construct that can simulate many other language features (state, I/O, greenthreading) [15], and have recently been added to OCaml [19]. It is thus strategic, and timely, to investigate the interaction between metaprogramming and effect handlers.

Unfortunately, metaprogramming is known to interact unexpectedly, and poorly, with effect handlers, a problem known as *scope extrusion* [10]. Scope extrusion occurs when the programmer accidentally directs the compiler to generate code with unbound variables. In the following program, we store the variable `Var(x)` into a heap cell (1), and retrieve it outside its scope. The program thus evaluates to `Var(x)`, which is unbound.

```

1  let l: int expr ref = new(Int(1)) in
2  let id_ast : (int -> int) expr = Lam(Var(x), l := Var(x); Int(1)) in
3  !l

```

OCaml

The problem of scope extrusion has been widely studied, resulting in multitudinous mechanisms for managing the interaction between metaprogramming and effect handlers. Some solutions adapt the type system (Refined Environment Classifiers [12, 7], Closed Types [2]), others insert dynamic checks into the generated code [10]. However, there are two issues.

First, we are a picky lot. Type-based approaches require unfeasible modification of the OCaml type-checker, and tend to limit expressiveness, disallowing a wide range of programs that do not lead to scope extrusion. Dynamic solutions are inefficient, reporting scope extrusion long after the error occurred, or unpredictable, allowing some programs, but disallowing morally equivalent programs. Further, many of the dynamic solutions have not been proved correct.

Second, and more importantly, we lacked a way to evaluate solutions fairly and holistically. Our evaluation criteria is three-pronged: correctness, efficiency, and expressiveness. Each solution is described in its own calculus, which made solutions difficult to compare. It is non-trivial to show that the definitions of scope extrusion in differing calculi agree. Hence, one solution may be correct with respect to its definition, but wrong with respect to an-

other. Further, expressiveness and efficiency are inherently comparative criteria.

This thesis will solve both these problems, in decreasing order of priority.

1.1 Contributions

Concretely, our contributions are:

1. A novel “universal” calculus that allows for effect handlers both in the generating and generated code.
2. In the calculus, three precise (but different) definitions of scope extrusion.
3. In the calculus, encodings of the following existing solutions:
 - a) Lazy dynamic check (with a proof of correctness)
 - b) Eager dynamic check (with a refutation of correctness)
 - c) Refined Environment Classifiers
4. A new solution, a Best-Effort dynamic check, that we justify is correct, expressive, and efficient.
5. Implementations of the three dynamic checks in Macocaml.

The scope of this dissertation thus extends to effect systems with **deep** (but **unnamed**) handlers and **multi-shot continuations**.

2 Background

The aim of this dissertation is to evaluate, and propose, policies for mediating the interaction between Macocaml-style **metaprogramming** and **effect handlers**, by addressing the issue of **scope extrusion**.

In this chapter, I provide technical overviews of each of the key concepts: metaprogramming ([Section 2.1](#)), effect handlers ([Section 2.2](#)), and scope extrusion ([Section 2.3](#)).

2.1 Metaprogramming

What is MacoCaml-style metaprogramming? I will provide an answer in two steps. First, I motivate metaprogramming, by illustrating the challenge of writing code that is both fast and maintainable ([Section 2.1.1](#)). Second, I will consider the design space of metaprogramming ([Section 2.1.2](#)), highlighting decisions made by MacoCaml.

2.1.1 Metaprogramming for Fast and Maintainable Code

Metaprogramming helps programmers write fast and maintainable code. How does one write fast and maintainable code? A naïve answer is “by being a skilled programmer”¹. Programmer skill is insufficient, because maintainability and efficiency are in constant tension.

I illustrate this tension by considering a concrete problem. Consider computing the gradient of a differentiable function as part of backpropagation over a neural network. More precisely, assume f of type differentiable

```
1 type differentiable = Sin | Tanh | Sigmoid | ...
2                   | Polynomial of float list
3                   | Compose of differentiable * differentiable
```

OCaml

For example, the following expression represents $\sin \circ \tanh$.

```
1 Compose(Tanh, Sin)
```

OCaml

We wish to write a function `grad` such that $\text{grad } f = f'$. For simplicity, assume the existence of a helper function, `grad_of`, that returns the gradient

¹Though not the worst answer: “use ChatGPT”

2 Background

of basic functions. For example, $\text{grad_of } \text{Sin } 0.0 = \cos 0.0 = 1.0$. `grad_main` (Listing 1) illustrates one way to compute gradients maintainably:

```
1 let rec grad_main f x = match f with
2   | Sin
3   | Tanh
4   | Sigmoid
5   | ...
6   | Polynomial(_) -> grad_of f x
7   | Compose(f, g) -> (grad_main f x) * (grad_main g (app f x))
```

OCaml

Listing 1: A maintainable implementation of grad

However, `grad_main` may not be the most efficient implementation. Performing a `match` on every recursive call might result in expensive branches. If x is a vector, and the weights of a polynomial are vectors, then `grad_main` could hide opportunities for cache prefetching.

If f is known in advance, for example, $f = \text{Compose}(\text{Tanh}, \text{Sin})$, we could implement a more efficient `grad_fast` function (Listing 2), whose body is simply a hardcoded equation:

```
1 let grad_fast x = (cos x) /. (cosh (sin x) ** 2)
```

OCaml

Listing 2: A fast implementation of grad, assuming $f = \text{Compose}(\text{Tanh}, \text{Sin})$

Although `grad_fast` only works for a single f , it has eliminated the branching overhead, and enabled opportunities for prefetching. It is thus likely to be faster.

The grad example illustrates the trade-off between maintainability and efficiency. `grad_main` is maintainable in part because it parameterises over f . More generally, abstraction centralises implementations, thus reducing maintenance costs. However, `grad_fast` is more efficient because it is more specialised to a specific f (`Compose(Tanh, Sin)`). More generally, many compiler optimisations, like monomorphisation, eliminate abstraction, simplifying functions by applying known arguments in advance.

The tension between maintainability (abstraction) and efficiency (specialisation) has also been observed in regex matching [20], parsing [25], linking [18], statistical modelling [22], and hardware design [21].

A more informed answer might therefore be “by letting the compiler generate optimised versions of my maintainable code”. Not quite: for reasons both theoretical and practical, compiler optimisations can be insufficient. In theory, we proposed an optimisation that assumed we would always know f at, or before, compile-time. Is this a reasonable assumption? It is: we assumed that grad would perform backpropagation over neural networks. The network over which backpropagation is performed is known at compile-time.

However, notice that this justification appeals to domain-specific knowledge regarding how `grad` will be used. In the general case, `grad` could be applied to a function not known until runtime. It is not feasible to expect a compiler to spot all opportunities for optimisation [16]. In practice, while compiler engineers might have an economic incentive to write optimisations for the machine learning community, this may not be true for less lucrative domains [17]. Even in machine learning, many libraries are built on top of existing languages, like Python, which might not perform the desired optimisations.

How does one write maintainable and efficient code, **when one cannot trust the compiler to optimise one's code**?

One answer is metaprogramming, which gives users the ability to perform code-generation. Programmers may thus take matters into their own hands: manually generating optimised code when the compiler may not automatically do so for them. The `grad` function in JAX, a Python-based machine learning framework, uses metaprogramming for precisely this purpose [8].

Speeding up exponentiation with Metaprogramming

Metaprogramming allows for code-generation via the creation and manipulation of abstract syntax trees (ASTs). I will now illustrate how metaprogramming works with reference to `MacoCaml`, which implements *compile-time* metaprogramming. While the `grad` example motivated metaprogramming, for pedagogical reasons, I switch to a morally equivalent, but simpler example: raising an integer x to an exponent n . One maintainable implementation is the `pow` function (Listing 3):

```
1 let rec pow (n: int) (x: int) =
2   if n == 0 then 1
3   else x * pow (n-1) x
```

OCaml

Listing 3: A maintainable implementation of an exponentiation function

`pow`, which can be applied to any exponent n , is analogous to `grad_main` (Listing 1), which could be applied to *any* differentiable function f .

However, should we know the exponent in advance, for example $n = 2$, then a more efficient, but less maintainable implementation, is the `square` function (Listing 4)

```
1 let square x = x * x
```

OCaml

Listing 4: An efficient implementation of exponentiation, assuming $n = 2$

`square` is analogous to `grad_fast` (Listing 2).

Metaprogramming can be utilised to write a function, `pow_gen`, which resembles `pow` (inheriting its maintainability), but that generates a program

2 Background

which resembles square (inheriting its efficiency). Listing 5 presents the meta-programmed `pow_gen` function. Compilation generates the body of square `y` (line 4), `y * y * 1`. I will now explain the mechanics of generation.

```
1 macro rec pow_gen (n: int) (x: int expr) =
2   if n == 0 then <<1>>
3   else <<$x * $(pow_gen (n-1) x)>>
4   let square y = $(pow_gen 2 <<y>>) (*after compile-time: y * y * 1 *)
5   square 3 (*at runtime: 9*)
```

OCaml

Listing 5: A meta-programmed `pow_gen` function, which resembles `pow` but generates square

Recall that (compile-time) metaprogramming gives the programmer the ability to generate programs at compile-time, for use at run-time. We may build this in two steps, by:

1. Deciding on a representation for code values, such that code can be created, and manipulated by programs. Once we have a representation for code values, it is possible to write expressions that return code values. These expressions serve as program generators.
2. Building a mechanism for executing expressions *at compile-time*. We can constrain this mechanism, using types, so only generators can be executed at compile-time.

First, we represent code values as ASTs. Generated programs are ASTs, and program generators are expressions that evaluate to ASTs. We can assume, for clarity, that the language offers, for each program construct, a corresponding AST node: for example, the integer 1 has AST node `Int(1)`. If a program has type 'a, then its AST node has type 'a expr. One can now write program generators, that evaluate to ASTs, for example:

```
1 let rec pow_gen (n: int) (x: int expr) =
2   if n == 0 then Int(1)
3   else Mul(x, (pow_gen (n-1) x))
4   pow_gen 2 Int(3) (*Mul(Int(3), Mul(Int(3), 1))**)
5   pow_gen 2 Var(y) (*Mul(Var(y), Mul(Var(y), 1))**)
6   pow_gen 3 Var(y) (*Mul(Var(y), Mul(Var(y), Mul(Var(y), 1)))**)
```

MacoCaml

Second, we need a mechanism to execute expressions at compile-time. In MacoCaml, this is the “top-level splice”, a splice (\$) annotation not surrounded by quotes (<<>>). For example, in Listing 5, there is only one top-level splice, on line 4: `$(pow_gen 2 <<y>>)`. We may now shift program generators (and only program generators) under top-level splices, to perform generation at

compile time. Note that to access `pow_gen` at compile-time, we must also move it under the top-level splice.

```

1  let square y = $(let rec pow_gen (n: int) (x: int expr) = ...
2                      in pow_gen 2 Var(y))
3                      (*Mul(Var(y), Mul(Var(y), 1))**)
4  let cube y   = $(let rec pow_gen (n: int) (x: int expr) = ...
5                      in pow_gen 3 Var(y))
6                      (*Mul(Var(y), Mul(Var(y), Mul(Var(y), 1))))*)

```

MacoCam1

To allow compile-time functions, like `pow_gen`, to be re-used across multiple top-level splices, MacoCam1 introduces the **macro** (Listing 6)

```

1  macro rec pow_gen (n: int) (x: int expr) =
2      if n == 0 then Int(1)
3      else Mul(x, (pow_gen (n-1) x))
4  let square y = $(pow_gen 2 Var(y)) (*Mul(Var(y), Mul(Var(y), 1))**)
5  let cube y   = $(pow_gen 3 Var(y)) (*Mul(Var(y), Mul(Var(y), Mul(Var(y), 1))))*)

```

OCaml

Listing 6: In MacoCam1, **macro** allows for definitions to be shared across top-level splices

Further, rather than explicit AST constructors, ASTs are created by the `<<>>` (“quote”) and `$` annotations. Quotation creates ASTs, by converting a program into its AST representation. For example,

$$<< \$x + 0 >> = \text{Plus}(\text{Var}(x), \text{Int}(0))$$

Under a quotation, the `$` annotation stops this conversion, allowing for programs that *manipulate* ASTs.

$$<< \$x + 0 >> = \text{Plus}(x, \text{Int}(0))$$

In MacoCam1, the programmer interleaves quotes and splices to perform code generation

$$<< $(\text{add_zero } <<1>>) + 0 >> = \text{Plus}(\text{add_zero } \text{Int}(1), \text{Int}(0))$$

Notice that `$` is overloaded. We must be careful to disambiguate between “top-level splices”, which execute programs at compile-time, and splices under quotations, which stop conversion to AST.

Re-writing Listing 6 in this style (being careful about non-top-level splices), we obtain exactly Listing 5.

Applying this technique to the grad example, we obtain

```

1  macro rec grad_gen f x = match f with
2    | Sin
3    | Tanh
4    | Sigmoid
5    | ...
6    | Polynomial(_) -> grad_of f x
7    | Compose(f, g) -> <<$ (grad_gen f x) * $(grad_gen g (app f x))>>
8
9  let grad_fast x = $(grad_gen Compose(Tanh, Sin) <<x>>)

```

MacoCam1

whathere `grad_of` and `app` are appropriately modified.

2.1.2 The Design Space of Metalanguages

Different metalanguages provide slightly different variants of metaprogramming to the user. In this section, I broadly taxonomise these languages by considering three key design decisions:

1. Homogenous or Heterogenous

Do the generated (“object”) and generating (“meta”) languages agree or differ?

If the object and meta languages are the same, this is known as homogenous metaprogramming. Otherwise, it is heterogenous [11].

MacoCam1 allows for homogenous metaprogramming, where OCaml code generates OCaml code. In contrast, MetaHaskell [13] programs generate C code, allowing for heterogenous metaprogramming.

2. Run-time or Compile-Time

When does the generation take place?

Code generation could take place at compile-time (as with MacoCam1 programs or C macros), or at run-time (as with MetaOCaml [10]).

Run-time and compile-time metaprogramming differ non-trivially. Run-time metaprogramming requires a language construct (!, or “run”) for explicit invocation of the compiler. Further, in run-time metaprogramming, generated and generating programs may share a heap.

MacoCam1 supports compile-time metaprogramming, and we will pay no further attention to run-time metaprogramming.

3. Two-stage or Multi-stage

How many stages of code generation are allowed?

When introducing MacoCam1, I illustrated how one uses top-level splices to perform shift computation from run-time (“level 0”) to compile-time (“level -1”). Might it be possible to shift computation from compile-time to a pre-compile-time (“level -2” phase), for example, via a nested splice?

```
1  $( $ pow_gen 2 Var(y))
```

MacoCam1

In a two-stage system, one is restricted to operating between two levels, so this is disallowed. In contrast, in a multi-stage system, one can operate between any number of levels. Multi-stage metaprogramming is thus strictly more general than two-stage metaprogramming.

Although nested splices are disallowed in MacoCam1, it is a multi-stage system, since entire modules may be imported at a decremented level [24].

MacoCam1 offers homogenous, compile-time, multi-stage metaprogramming. The scope of this dissertation is slightly more restrictive: I focus on two-stage, not multi-stage metaprogramming. This restriction was motivated by a cost-benefit analysis:

1. **Cost:** Since in MacoCam1, the module system is the only mechanism for achieving multi-stage programming, investigating multi-stage metaprogramming would require the investigation of module systems, effects, and metaprogramming. The interaction between module systems and metaprogramming is still an ongoing area of research [3].
2. **Benefit:** In practice, “almost all uses” of multi-stage metaprogramming only use two stages [6]. Further, scope extrusion can be observed, and is often studied, in two-stage systems [7, 12].

2.2 Effect Handlers

What is an effect handler? I will first motivate effect handlers by considering the problem of adding resumable exceptions to OCaml ([Section 2.2.1](#)). Second, I will introduce a fine-grained call-by-value calculus for studying the operational behaviour of effect handlers, à la Pretnar [15] ([Section 2.2.2](#)). This calculus will be useful both for precise description of effect handlers, and as a basis for investigating the interaction between metaprogramming and effect handlers (once the calculus has been extended with metaprogramming facilities). Finally, since different design decisions for effect handlers could affect the nature of their interaction with metaprogramming, I will consider the design space of effect handlers ([Section 2.2.3](#)).

2.2.1 Composable and Customisable Effects

Effects are a mechanism by which a program interacts with its environment. Examples of effects include state, (resumable) exceptions, non-determinism, and I/O. Effects are typically defined and understood separately, meaning

2 Background

they are not easily composable. They are also implemented by compiler engineers rather than programmers, meaning they are not customisable. Effect handlers provide a programmable, unifying framework that may be instantiated into different effects. This allows for composable and customisable treatment of effects.

To illustrate the need for effect handlers, consider the following problem, by Kiselyov [9]. Assume a binary search tree of (key, value) pairs. The following code provides two functions. The first finds a value v associated with key k , raising a `NotFound` exception if k is not in the tree. The second updates the dictionary with a fresh key value pair, overwriting old values.

```
1  type ('a, 'b) tree = Lf | Br of 'a * 'b * tree * tree
2
3  let rec find (t: tree) (k: 'a) = match t with
4  | Lf -> raise NotFound()
5  | Br(k', v, l, r) -> if k == k' then v
6                        else if k < k' then find l k
7                        else find r k
8
9  let rec update (t: tree) (k: 'a) (v: 'b) = match t with
10 | Lf -> Br(k, v, Lf, Lf)
11 | Br(k', v', l, r) -> if k == k' then Br(k, v, l, r)
12                       else if k < k' then Br(k', v', update l k v, r)
13                       else Br(k', v', l, update r k v)
```

Assume the task is to build a `findOrInsert` function that either finds the value associated with a key, *or* inserts a default value. A naïve approach to writing this function would be

```
1  let rec findOrInsert (t: tree) (k: 'a) (default: 'b) =
2  try find t k with NotFound -> insert t k default
```

This function is **inefficient**. If a `NotFound` exception is raised, then the `find` function will have raised the exception at the point where the default value should be inserted. The function could be twice as efficient if the exception could be resumed at the point where the exception was raised, in the following style.

```
1  let rec findOrInsert (t: tree) (k: 'a) (default: 'b) =
2  try find t k with NotFound(p) -> continue p Br(k, default, Lf, Lf)
```

p represents the suspended program to be resumed, and is known as a *delimited continuation*.

The aforementioned problem motivates the need for resumable exceptions. To understand the need for effect handlers, consider how one might go about

implementing resumable exceptions. One approach might be to fork the implementation of handlers and tweak it ever-so-slightly. This solution does not scale well. First, the solution may not be **composable**. The intended informal semantics for resumable exceptions is “effectively equivalent to exceptions, with the additional power to resume programs”. Resumable exceptions should thus interact with other exceptions in a predictable way, but this is difficult to guarantee, and *continually* guarantee, especially as implementations evolve, and more variants of exceptions are demanded. Second, the solution is not **customisable**. To add resumable exceptions requires a compiler engineer to modify the compiler. With the exception of raising an issue, there is nothing the programmer may do, in the moment, to meet their need.

Effect handlers resolve both composability and customisability issues. Much like how exception handlers allow users to create custom exceptions with custom semantics, effect handlers provide a general framework for creating custom effects with custom semantics. The interaction between effect handlers is described abstractly, parameterising over the exact semantics of the effect. Hence, implementing effects (in the earlier example, resumable exceptions, but more generally, state, I/O, greenthreading, non-determinism, and more) as effect handlers ensures composability by design.

With effect handlers, we can re-write the previous example to obtain the behaviour of resumable exceptions, even if the OCaml compiler does not support it, with the guarantee that `NotFound` will interact predictably with other defined effects.

```

1  type _ Effect.t += NotFound: unit -> tree t
2
3  let rec find (t: tree) (k: 'a) = match t with
4  | Lf -> NotFound()
5  | Br(k', v, l, r) -> if k == k' then v
6                        else if k < k' then find l k
7                        else find r k
8
9  let rec findOrInsert (t: tree) (k: 'a) (default: 'b) =
10    match find t k with NotFound(p) with
11    | v -> v
12    | effect NotFound k -> continue p Br(k, default, Lf, Lf)

```

OCaml

Since effect handlers may be instantiated into a range of different effects, considering the interaction of metaprogramming with effect handlers is an exercise in killing many birds with a single stone. Additionally, effect handlers were recently added to OCaml [19], making their interaction a timely problem.

Values	$v ::= x \mid n \mid \lambda x.c \mid \kappa x.c \mid h$	λ_{op}
Computations	$c ::= v_1 v_2 \mid \text{return } v \mid \text{do } x \leftarrow c_1 \text{ in } c_2 \mid \text{op}(v) \mid \text{handle } c \text{ with } \{h\} \mid \text{continue } v_1 v_2$	
Handlers	$h ::= \text{return}(x) \mapsto c \mid h; \text{op}(x, k) \mapsto c$	

Figure 2.1: The syntax of λ_{op} . Terms are syntactically divided into values v and computations c

2.2.2 λ_{op} : A Calculus for Effect Handlers

Having motivated effect handlers, I will now describe a calculus, which I call λ_{op} , for reasoning about their operational behaviour. λ_{op} is a slight variant of the calculus described by Pretnar [15]. Understanding λ_{op} will be useful for two reasons. First, a precise description of the operational behaviour of effects will aid reasoning about their interaction with metaprogramming. Second, my universal calculus will be described by extending λ_{op} . Throughout this section, I will use the λ_{op} program in [Listing 7](#) as a running example.

```

handle
  do  $x \leftarrow \text{print}(1); \text{return } 1$  in do  $y \leftarrow \text{print}(2); \text{return } 2$  in  $x + y$ 
with
  {return( $x$ )  $\mapsto$  return ( $x, ""$ );
   print( $x, k$ )  $\mapsto$  do ( $v, s$ )  $\leftarrow$  continue  $k()$  in return ( $v, f"\{x\};" ^ s$ )}

return ( $3, "1;2"$ )

```

Listing 7: An λ_{op} program that returns $(3, "1;2")$. It will be used as a running example throughout this section.

[Figure 2.1](#) collates the base syntax of λ_{op} . In addition to this base syntax, in this section, I will assume λ_{op} is extended with the following language extensions: a unit value $()$, pairs $(1, 2)$ which can be destructured **do** $(x, y) \leftarrow \text{return } (1, 2)$ **in** $x + y$, strings "Hello", format strings $f"\{1\}"$, and string concatenation $^$. For example, the following code evaluates to "Revolution 9".

```

do ( $x, y$ )  $\leftarrow$  return ("Revolution",  $f"\{9\}"$ ) in  $x ^ y$ 

return "Revolution 9"

```

Further, I use $c_1; c_2$ as syntactic sugar for **do** $_ \leftarrow c_1$ **in** c_2 . I will explain key language constructs in turn.

Sequencing computations: do and return

Effects force us to carefully consider the order of evaluation. For example, consider the following OCaml programs

```
1 let pure      = (1+0) + (2+0)
2 let effectful = let l = new 0 in (l := 1; 1) + (l := 2; 2)
```

OCaml

The result of `pure`, which has no effects, is independent of the evaluation order. In contrast, the result of `effectful` is dependent on the evaluation order. If terms are evaluated left-to-right, the value of `!l` is 2, otherwise, it is 1.

In order to be precise about the order of evaluation, λ_{op} terms are stratified into two syntactic categories, “inert values” v and “potentially effectful computations” c [15]. **return** v lifts values into computations, and is also the result of fully evaluating a computation. **do** $x \leftarrow c_1$ **in** c_2 sequences computations, forcing programmers to be explicit the order of evaluation. First, c_1 is fully evaluated to obtain some **return** v . The value v is then bound to x , and finally c_2 is evaluated.

For example, extending λ_{op} with a plus function, what is the order of evaluation of `plus 1 2`? Do we evaluate both arguments before applying them, or interleave evaluation and application? The syntax forces programmers to choose explicitly. We can either fully evaluate both arguments before applying them in turn,

```
do x ← return 1 in (do y ← return 2 in (do f ← plus x in fy))
```

 λ_{op}

or alternatively, evaluate 1, apply it, then evaluate 2

```
do x ← return 1 in (do f ← plus x in do y ← return 2 in fy)
```

 λ_{op}

Both choices are valid, but the programmer must choose. For clarity, where the ordering cannot affect the result (both of the aforementioned choices evaluate to **return** 3), I will abuse notation and write (for instance) `1 + 2`.

Performing effects: op, handle, and continue

Having made explicit the order of operation, we may now add effect handlers. Recall that effect handlers allow users to register custom effects with custom semantics. I will now illustrate how this is supported by λ_{op} .

For simplicity, λ_{op} assumes that the effects have been registered in advance, parameterising over them with the placeholder `op(v)`. Assume that the user has declared the effects `print` and `read_int` in advance. This would allow the user to write programs like

```
do  $x \leftarrow \text{print}(1)$ ; return 1 in do  $y \leftarrow \text{print}(2)$ ; return 2 in  $x + y$ 
```



In the program fragment above, we know that **print** is an effect, but we do not know its semantics. Effect handlers, which comprise a **return handler** and zero or more **operation handlers**, specify how effects interact with their environment, and thus may be used to give effects meaning. I will define an effect handler that accumulates print statements in a string (some “stdout”). For example, the aforementioned program should return (3, “1; 2”).

We begin by considering how to handle the case where there are no calls to **print**. For example, in the program **return 3**. We may wish to return both the value, and the empty string (empty stdout) to the environment: in this case, (3, “”). We can achieve this by specifying a *return handler*.

$$\text{return}(x) \mapsto c$$

In this case, we set c to **return** (x , “”). All effect handlers must specify a return handler. In many cases, the return handler is simply the identity (c is set to **return** x): for brevity and clarity, if the return handler is the identity, I may drop it.

Next, we consider how to handle a call to **print**. We use an operation handler of the form

$$\text{print}(x, k) \mapsto c$$

Where c is the user-defined semantics for **print**. Concretely, one instance of c is

$$\text{print}(x, k) \mapsto \text{do } (v, s) \leftarrow \text{continue } k() \text{ in return } (v, f" \{x\}; " \wedge s)$$

In the definition of c , the programmer may refer to x and k , which I will now explain. x allows programs to send values (for example, values to be printed) to their environment. k is a delimited continuation representing a suspended program, awaiting a value from the environment. Effects also allow programs to receive data from their environment, as in

$$1 + \text{get_int_from_user}()$$

Note that the program is suspended until the value is received. We may write the suspended program as $1 + [-]$, where $[-]$ indicates an as-yet-unknown value. This suspended program is represented by the continuation k . The expression

$$\text{continue } k v$$

is used to resume the suspended program with value v .

We are now able to interpret the concrete operation handler c : we resume the suspended program, supplying a unit value, since **print** effects do not receive values from their environment. This returns a value v and some partially accumulated stdout s . We prepend the printed value, x , onto s .

Having defined the semantics for **print**, the user may now interpret the earlier example with their semantics, using the **handle** e **with** $\{h\}$ construct. Doing so results in the program in [Listing 7](#).

Notice that multiple effects may be handled by the same handler, and the same effect might be handled by multiple handlers, potentially with different semantics.

Operational Semantics

Having described informally the desired semantics of λ_{op} , we may now make our intuitions precise, by means of an operational semantics. The operational semantics is collated in [Figure 2.2](#).

Auxiliary Definitions		λ_{op}
Evaluation Frame	$F ::= \text{do } x \leftarrow - \text{ in } c_2 \mid \text{handle } - \text{ with } \{h\}$	
Evaluation Context	$E ::= [] \mid E :: F$	
Domain of Handler	$\text{dom}(h) \triangleq \begin{aligned} &\text{dom}(\text{return}(x) \mapsto c) = \emptyset, \\ &\text{dom}(h; \text{op}(x, k) \mapsto c) = \text{dom}(h) \cup \{\text{op}\} \end{aligned}$	
Handled Effects	$\text{handled}(E) \triangleq \begin{aligned} &\text{handled}([]) = \emptyset, \\ &\text{handled}(E :: \text{do } x \leftarrow - \text{ in } c_2) = \text{handled}(E), \\ &\text{handled}(E :: \text{handle } - \text{ with } \{h\}) = \text{handled}(E) \cup \text{dom}(h), \end{aligned}$	
Operational Semantics		
(RED-APP)	$(\lambda x. c)v; E \rightsquigarrow c[v/x]; E$	
(RED-SEQ)	$\text{do } x \leftarrow \text{return } v \text{ in } c; E \rightsquigarrow c[v/x]; E$	
(RED-HDL)	$\text{handle return } v \text{ with } \{h\}; E \rightsquigarrow c[v/x]; E \quad (\text{where } \text{return}(x) \mapsto c \in h)$	
(CNG-PSH)	$F[c]; E \rightsquigarrow c; E :: F$	
(CNG-POP)	$\text{return } v; E :: F \rightsquigarrow F[v]; E$	
(EFF-OP)	$\text{op}(v); E_1 @ [h] @ E_2 \rightsquigarrow c[v/x, \kappa x. \text{handle } E_2[\text{return } x] \text{ with } \{h\}/k]; E_1$ (where $\text{op}(x, k) \mapsto c \in h$ and $\text{op} \notin \text{handled}(E_2)$)	
(EFF-CNT)	$\text{continue } E_2 v; E_1 \rightsquigarrow \text{return } v; E_1 @ E_2$	

Figure 2.2: The operational semantics of λ_{op} . The semantics is given on configurations of the form $\langle c, E \rangle$, with the brackets dropped for clarity. Rules are divided into three classes: reduction rules RED- X , which perform computation, congruence rules CNG- Y which manipulate the evaluation context, and effect rules EFF- Z that are special to λ_{op}

The operational semantics is given on configurations of the form $\langle c, E \rangle$, where c is a term and E is an evaluation context, in the style of Felleisen and Friedman [4]. Evaluation contexts are represented as a stack of evaluation frames F , à la Kiselyov [9]. For clarity, I will compress multiple stack frames using nesting. For example, instead of $\text{do } x \leftarrow - \text{ in } c_2 :: \text{do } y \leftarrow - \text{ in } c_1$, I will write $\text{do } x \leftarrow (\text{do } y \leftarrow - \text{ in } c_1) \text{ in } c_2$. Most of the rules are standard. We will

2 Background

focus on two rules: **EFF-OP**, the mechanism for giving effects custom semantics, and **EFF-CNT**, the mechanism for resuming programs.

To illustrate the operation of **EFF-OP** and **EFF-CNT**, consider the evaluation of the running example in [Listing 7](#), beginning with an empty context. Let h be the handler body

```
{return( $x$ )  $\mapsto$  return ( $x$ , "");
 print( $x$ ,  $k$ )  $\mapsto$  do ( $v$ ,  $s$ )  $\leftarrow$  continue  $k$  () in return ( $v$ , f"{" $x$ "}; " ^  $s$ )}
```

After several applications of **CNG-POP**, we obtain the configuration

```
<print(1) ; handle - with { $h$ } ::
  do  $x \leftarrow -$ ; return 1 in do  $y \leftarrow$  print(2); return 2 in  $x + y$ >
```

Let $E = \text{do } x \leftarrow \text{return } u; \text{return } 1 \text{ in do } y \leftarrow \text{print}(2); \text{return } 2 \text{ in } x + y$. Applying **EFF-OP**, we can suspend the program, find the handler h with the user's semantics for **print**, and give the **print** effect the desired semantics

```
<do ( $v$ ,  $s$ )  $\leftarrow$  continue ( $\kappa u.E$ ) () in return ( $v$ , f"{"1"}; " ^  $s$ ) ; handle - with { $h$ }>
```

Applying **CNG-POP**,

```
<continue ( $\kappa u.E$ ) () ; handle - with { $h$ } ::
  do ( $v$ ,  $s$ )  $\leftarrow -$  in return ( $v$ , f"{"1"}; " ^  $s$ )>
```

Applying **EFF-CNT**, we can resume the program that was suspended

```
<return () ; handle - with { $h$ } ::
  do ( $v$ ,  $s$ )  $\leftarrow -$  in return ( $v$ , f"{"1"}; " ^  $s$ )
  do  $x \leftarrow -$ ; return 1 in do  $y \leftarrow$  print(2); return 2 in  $x + y$ >
```

The side-condition on **EFF-OP** is needed because the user may define multiple handlers with different semantics for the same effect. The side-condition resolves any ambiguity by using the *latest* handler. For example, the following program has a **read** effect that is given two definitions: it could read either 1 or 2. The ambiguity is resolved by choosing the latest handler: in this case, 1.

```
handle
  handle read() with {return( $y$ )  $\mapsto$  return  $y$ ; read( $x$ ,  $k$ )  $\mapsto$  continue  $k$  1}
with
  {return( $y$ )  $\mapsto$  return  $y$ ; read( $x$ ,  $k$ )  $\mapsto$  continue  $k$  2}
return 1
```

λ_{op}

Effects row	Δ	$::= \emptyset \mid \Delta \cup \{\text{op}_i\}$	λ_{op}
Value type	T	$::= \mathbb{N}$ $\mid T_1 \xrightarrow{\Delta} T_2$ functions $\mid T_1 ! \Delta \Longrightarrow T_2 ! \Delta'$ handlers $\mid T_1 \xrightarrow{\Delta} T_2$ continuations	
Computation type	$T ! \Delta$		

Figure 2.3: λ_{op} types. Notice that, just as terms are divided into values and computations, types are divided into value types (T) and computation types ($T ! \Delta$)

Type-and-Effect System

We now give a type-and-effect system to λ_{op} . Figure 2.3 collates the syntax of λ_{op} types, which I will now briefly describe.

Syntactically, just as terms are divided into values and computations, types are divided into value types (for example, \mathbb{N}) and computation types (for example, $\mathbb{N} ! \{\text{print}\}$). Since computations may have effects, computation types track unhandled effects using an effects row (typically labelled Δ), which in this system is simply a set. This type-and-effect system allows us to distinguish between values, computations that return values, and computations that return values and additionally have some unhandled side effects.

Term	Type
3	\mathbb{N}
$\text{do } x \leftarrow \text{return } 1 \text{ in } \text{do } y \leftarrow \text{return } 2 \text{ in } x + y$	$\mathbb{N} ! \emptyset$
$\text{do } x \leftarrow \text{print}(1); \text{return } 1 \text{ in } \text{do } y \leftarrow \text{return } 2 \text{ in } x + y$	$\mathbb{N} ! \{\text{print}\}$

Functions are values, and are applied to other values, but produce computations on application. For example, the function

$$\lambda x : \mathbb{N}. \text{print}(x); \text{return } x$$

is a value that accepts a value of type \mathbb{N} and returns a computation of type $\mathbb{N} ! \{\text{print}\}$. We thus say functions have suspended effects, which we write $T_1 \xrightarrow{\Delta} T_2$. In this case, the function has type $\mathbb{N} \xrightarrow{\{\text{print}\}} \mathbb{N}$. For technical reasons, continuations and functions need to be distinguished, but in most cases they may be treated equivalently.

Handlers transform computations of one type to computations of another type. This happens in two ways: first, by handling effects, and thus removing them from the effects row (which recall represents unhandled effects). Second, by modifying the return type of computations. To reflect both abilities,

λ_{op}

<p>(NAT)</p> $\frac{}{\Gamma \vdash n : \mathbb{N}}$	<p>(VAR)</p> $\frac{\Gamma(x) = T}{\Gamma \vdash x : T}$	<p>(LAMBDA)</p> $\frac{\Gamma, x : T_1 \vdash c : T_2 ! \Delta}{\Gamma \vdash \lambda x. c : T_1 \xrightarrow{\Delta} T_2}$	<p>(CONTINUATION)</p> $\frac{\Gamma, x : T_1 \vdash c : T_2 ! \Delta}{\Gamma \vdash \kappa x. c : T_1 \xrightarrow{\Delta} T_2}$
<p>(APP)</p> $\frac{\Gamma \vdash v_1 : T_1 \xrightarrow{\Delta} T_2 \quad \Gamma \vdash v_2 : T_1}{\Gamma \vdash v_1 v_2 : T_2 ! \Delta}$		<p>(CONTINUE)</p> $\frac{\Gamma \vdash v_1 : T_1 \xrightarrow{\Delta} T_2 \quad \Gamma \vdash v_2 : T_1}{\Gamma \vdash \mathbf{continue} \ v_1 v_2 : T_2 ! \Delta}$	
<p>(RETURN)</p> $\frac{\Gamma \vdash v : T}{\Gamma \vdash \mathbf{return} \ v : T ! \Delta}$	<p>(DO)</p> $\frac{\Gamma \vdash c_1 : T_1 ! \Delta \quad \Gamma, x : T_1 \vdash c_2 : T_2 ! \Delta}{\Gamma \vdash \mathbf{do} \ x \leftarrow c_1 \ \mathbf{in} \ c_2 : T_2 ! \Delta}$		
<p>(OP)</p> $\frac{\Gamma \vdash v : T_1 \quad \text{op} : T_1 \rightarrow T_2 \in \Sigma \quad \text{op} \in \Delta}{\Gamma \vdash \mathbf{op}(v) : T_2 ! \Delta}$	<p>(HANDLE)</p> $\frac{\Gamma \vdash c : T_1 ! \Delta \quad \Gamma \vdash h : T_1 ! \Delta \Longrightarrow T_2 ! \Delta' \quad \forall \text{op} \in \Delta \setminus \Delta'. \text{op} \in \text{dom}(h)}{\Gamma \vdash \mathbf{handle} \ c \ \mathbf{with} \ \{h\} : T_2 ! \Delta'}$		
<p>(RET-HANDLER)</p> $\frac{\Gamma, x : T_1 \vdash c : T_2 ! \Delta'}{\Gamma \vdash \mathbf{return}(x) \mapsto c : T_1 ! \Delta \Longrightarrow T_2 ! \Delta'}$			
<p>(OP-HANDLER)</p> $\frac{\text{op} : A \rightarrow B \in \Sigma \quad \Gamma \vdash h : T_1 ! \Delta \Longrightarrow T_2 ! \Delta' \quad \Gamma, x : A, k : B \xrightarrow{\Delta'} T_2 \vdash c : T_2 ! \Delta' \quad \Delta' \subseteq \Delta \setminus \{\text{op}\} \quad \mathbf{op}(x', k') \mapsto c' \notin h}{\Gamma \vdash h; \mathbf{op}(x, k) \mapsto c : T_1 ! \Delta \Longrightarrow T_2 ! \Delta'}$			

Figure 2.4: Typing rules for λ_{op} terms

handlers are given a type of the form $T_1 ! \Delta \Longrightarrow T_2 ! \Delta'$. For example, a handler of the form

```
{return(x) ↦ return (x, "");
  print(x, k) ↦ do (v, s) ← continue k () in return (v, f" {x}; " ^ s)}
```

may be given type $\mathbb{N} ! \{\mathbf{print}\} \Longrightarrow (\mathbb{N} \times \text{String}) ! \emptyset$, reflecting both the handling of the **print** effect and the transformation of the return type to include the collated print statements.

I now consider the typing rules for terms, which are collated in [Figure 2.4](#). Most rules are standard, but a few are worth paying attention to.

First, the RETURN and DO rules. In the RETURN rule, we are allowed to assign

the term **return** v any set of effects. For example, we could write:

$$\frac{}{\Gamma \vdash \mathbf{return} \ 0 : \mathbb{N}! \{\mathbf{print}\}}$$

This flexibility is important, because to type $\mathbf{do} \ x \leftarrow c_1 \ \mathbf{in} \ c_2$, the Do rule requires both c_1 and c_2 to have the same effects. For example, without this flexibility, we would not be able to complete the following typing derivation

$$\frac{\vdots}{\Gamma \vdash \mathbf{do} \ x \leftarrow \mathbf{print}(0) \ \mathbf{in} \ \mathbf{return} \ 0 : \mathbb{N}! \{\mathbf{print}\}}$$

A valid alternative would be to forbid this flexibility and add explicit subtyping. However, such an approach would no longer be syntax directed.

Second, the OP rule. Previously, we assumed that the user declared their effects in advance. We also assume that they declare the types of their effects in advance, and that we store the mapping from effects to types in Σ . For example, we might assume $\Sigma = \{\mathbf{print} : \mathbb{N} \rightarrow 1\}$. In OCaml, this would correspond to writing:

```
1 type _ Effect.t += Print: nat -> unit OCaml
```

Note further the $\text{op} \in \Delta$ restriction – flexibility allows us to over-approximate the effects in a term, but never underapproximate them.

Third, the RET-HANDLER and OP-HANDLER rules, which are used to type handlers, which I will explain by means of an example. Assume we are trying to type the handler

$$\begin{aligned} \{\mathbf{return}(x) &\mapsto \mathbf{return} \ (x, ""); \\ \mathbf{print}(x, k) &\mapsto \mathbf{do} \ (v, s) \leftarrow \mathbf{continue} \ k \ () \ \mathbf{in} \ \mathbf{return} \ (v, f''\{x\}; '' ^ s)\} \end{aligned}$$

with the type $\mathbb{N}! \{\mathbf{print}\} \Longrightarrow (\mathbb{N} \times \text{String})! \emptyset$. We apply the OP-HANDLER rule, which is transcribed below. Preconditions are numbered for reference.

$$\begin{array}{c} \text{(OP-HANDLER)} \\ \begin{array}{l} (1) \text{op} : A \rightarrow B \in \Sigma \\ (2) \Gamma \vdash h : T_1! \Delta \Longrightarrow T_2! \Delta' \quad (3) \Gamma, x : A, k : B \xrightarrow{\Delta'} T_2 \vdash c : T_2! \Delta' \\ (4) \Delta' \subseteq \Delta \setminus \{\text{op}\} \quad (5) \mathbf{op}(x', k') \mapsto c' \notin h \end{array} \\ \hline \Gamma \vdash h; \mathbf{op}(x, k) \mapsto c : T_1! \Delta \Longrightarrow T_2! \Delta' \end{array}$$

The preconditions of the OP-HANDLER rule direct us to check, in turn:

- (1) $\mathbf{print} : \mathbb{N} \rightarrow 1 \in \Sigma$, which is true by assumption
- (2) Recursively check the rest of the handler $h = \mathbf{return}(x) \mapsto \mathbf{return} \ (x, "")$, ensuring it has type $\mathbb{N}! \{\mathbf{print}\} \Longrightarrow (\mathbb{N} \times \text{String})! \emptyset$. This follows from a trivial application of the RET-HANDLER rule.

2 Background

- (3) Assuming x has type \mathbb{N} and k has type $1 \xrightarrow{\emptyset} (\mathbb{N} \times \text{String})$, the body

do $(v, s) \leftarrow \text{continue } k ()$ **in return** $(v, f''\{x\}; " ^ s)$

has type $(\mathbb{N} \times \text{String}) ! \emptyset$. This is easy to show.

- (4) That the handler *only* removes **print** from the effects row, and no other effects. This check passes, but would fail if we tried to type the handler with, for example, $\mathbb{N} ! \{\text{print}, \text{get}\} \Longrightarrow (\mathbb{N} \times \text{String}) ! \emptyset$
- (5) That there are no other handlers for **print** in h .

A full typing derivation may be found in the appendix.

2.2.3 The Design Space of Effect Handlers

The design space of effect handlers is large. I consider three key design decisions made by different systems.

1. Deep, Shallow, or Sheep

Are multiple instances of the same effect handled by the same handler? In λ_{op} , continuations reinstate handlers (EFF-CNT) and thus multiple instances of the same effect are handled by the same handler. For example, in the following example, the effect **addn** is handled by the same handler, adding one each time. We say these handlers are **deep**.

```

handle
  addn(1) + addn(2)
with
  {return( $x$ )  $\mapsto$  return  $x$ ;
   addn( $y, k$ )  $\mapsto$  continue  $k (y + 1)$ }

return 5

```

λ_{op}

We may also choose *not* to reinstate the handler, in an approach known as **shallow** handlers [5]. The example above would be stuck, since the second **addn** would not be handled.

Finally, we could choose to modify the interface for **continue** such that it accepts a handler

continue $k \ v \ h$

This would allow multiple effects to be handled by different handlers. That is, we could add one the first time **addn** is performed, and two the second time. These handlers behave as a hybrid of shallow and deep handlers, and are thus termed **sheep** handlers [14].

OCaml allows the programmer to choose between shallow and deep handlers. Since most prior work on scope extrusion focuses on deep handlers [7], we focus on those.

2. One-Shot or Multi-Shot

How many times can one resume the same continuation?

In λ_{op} , continuations may be resumed multiple times. For example, we can write

```

handle
  performTwice(1)
with
  {return( $x$ )  $\mapsto$  return  $x$ ;
   performTwice( $y, k$ )  $\mapsto$  (continue  $k$   $y$ ) + (continue  $k$   $y$ )}

return 2

```

We say the effect system permits **multi-shot continuations**. Multi-shot continuations are useful for simulating certain effects, like non-determinism [14].

In other systems, like OCaml and WasmFX [14], this is not allowed: continuations are only allowed to be resumed once. These systems permit **one-shot continuations**.

Although continuations in OCaml are one-shot, due to the utility of multi-shot continuations, I believe it is worthwhile to study effect systems with multi-shot continuations.

3. Named or Unnamed Handlers

Can I invoke a specific handler for an operation?

In λ_{op} , when there are multiple handlers for the same effect, we invoke the “nearest” or “most recent” handler for that effect. An alternative approach is **named handlers** [23], where, by associating each handler with a *name*, we can more easily specify which handler should be invoked.

While named handlers can be more ergonomic, they do not provide greater expressiveness than using distinct effects [23]. I do not consider named handlers in this thesis.

2.3 Scope Extrusion

I now turn my attention to scope extrusion, which arises from the unexpected interaction of effects and metaprogramming. To illustrate scope extrusion, I

2 Background

will first extend λ_{op} (Page 16) with AST constructors $\text{Var}(x_T)$, $\text{Nat}(n)$, Lam , and Plus . For example, we may generate the AST of $\lambda x : \mathbb{N}. x + 0$ as follows:

```
return Lam(Var( $x_{\mathbb{N}}$ ), Plus(Var( $x_{\mathbb{N}}$ ), Nat(0)))
```

λ_{op}

Listing 8 illustrates the problem of scope extrusion. The program constructs the AST of $\lambda x : \mathbb{N}. x$, but additionally performs an effect, **extrude**, with type $\mathbb{N} \text{expr} \rightarrow \mathbb{N} \text{expr}$. The handler for **extrude** discards the continuation, simply returning the value it was given: $\text{Var}(x_{\mathbb{N}})$. The entire program evaluates to $\text{Var}(x_{\mathbb{N}})$, and the generated AST is ill-scoped. We say that the result of evaluation demonstrates scope extrusion.

```
handle
  do body  $\leftarrow$  extrude(Var( $x_{\mathbb{N}}$ )) in return Lam(Var( $x_{\mathbb{N}}$ ), body)
with
  {return( $u$ )  $\mapsto$  return Nat(0);
   extrude( $y, k$ )  $\mapsto$  return  $y$ }

return Var( $x_{\mathbb{N}}$ )
```

λ_{op}

Listing 8: A λ_{op} program that evaluates to the $\text{Var}(x_{\mathbb{N}})$. The AST is ill-scoped, and thus exhibits scope extrusion. It will be used as a running example.

It is difficult to give a precise definition to scope extrusion, because there are multiple competing definitions [10, 12], and many are given informally. For example, is scope extrusion a property of the *result* of evaluation [12], as in **Listing 8**, or is it a property of *intermediate* configurations [10]? We can, for example, build ASTs with extruded variables, that are bound at some future point. In **Listing 11**, we produce the intermediate AST $\text{Plus}(\text{Nat}(0), \text{Var}(x_{\mathbb{N}}))$, which is not well scoped. However, the result of evaluation is well scoped: $\text{Lam}(\text{Var}(x_{\mathbb{N}}), \text{Plus}(\text{Nat}(0), \text{Var}(x_{\mathbb{N}})))$. Does **Listing 11** exhibit scope extrusion?

Nevertheless, all definitions agree on the example in **Listing 8**. Making precise the competing definitions of scope extrusion these competing definitions, and their relation to one another, is a contribution of this dissertation.

```
handle
  do body  $\leftarrow$  extrude(Var( $x_{\mathbb{N}}$ )); return Var( $x_{\mathbb{N}}$ ) in return Lam(Var( $x_{\mathbb{N}}$ ), body)
with
  {return( $u$ )  $\mapsto$  return  $u$ ;
   extrude( $y, k$ )  $\mapsto$  do  $z \leftarrow$  return  $y$  in continue  $k z$ }

return Lam(Var( $x_{\mathbb{N}}$ ), Plus(Nat(0), Var( $x_{\mathbb{N}}$ )))
```

λ_{op}

Listing 9: A λ_{op} program that is a slight variation of Listing 11, but that (unlike Listing 11) passes the Eager Dynamic Check.

2.3.1 Existing Solutions to the Scope Extrusion Problem

There are multiple solutions to the problem of scope extrusion. The solution space can be broadly divided into two types of approaches: static (type-based) and dynamic. I will now survey two dynamic approaches, which I term the lazy and eager checks, and one static approach, the method of refined environment classifiers [12, 7].

Lazy Dynamic Check

Scope extrusion, at least, of the kind in Listing 12, may seem trivial to resolve: evaluate the program to completion, and check that the resulting AST is well-scoped [10]. I term this the Lazy Dynamic Check. This approach, while clearly correct and maximally expressive, is not ideal for efficiency and error reporting reasons.

To illustrate the inefficiency of this approach, consider a slight variation of Listing 8, Listing 10. In Listing 10, we can, in theory, report a warning as soon scope extrusion is detected. However, waiting for the result of the program can be much more inefficient.

In terms of error reporting, note that, in waiting for the result of execution, we lose information about *which program fragment* was responsible for scope extrusion, reducing the informativeness of reported errors [10].

```

do  $x \leftarrow$  handle
  do body  $\leftarrow$  extrude(Var( $x_{\mathbb{N}}$ )) in return Lam(Var( $x_{\mathbb{N}}$ ), body)
  with
    {return( $u$ )  $\mapsto$  return Nat(0);
     extrude( $y, k$ )  $\mapsto$  return  $y$ }
  in some very long program; return  $x$ 

return Var( $x_{\mathbb{N}}$ )

```

λ_{op}

Listing 10: A λ_{op} program that evaluates to the Var($x_{\mathbb{N}}$). Executing the entire program to determine if it exhibits scope extrusion is inefficient.

Eager Dynamic Check

A second dynamic check, motivated by the problems with the Lazy Dynamic Check, adopts a stricter definition of scope extrusion. During the code generation process, one inserts checks into the running program, reporting errors

2 Background

when one encounters unbound free variables in intermediate ASTs. Hence, the Eager Dynamic Check would classify the program in [Listing 11](#) as exhibiting scope extrusion. The Eager Dynamic Check has been adopted by BER MetaOCaml, and offers better efficiency and error reporting guarantees over the Lazy Dynamic Check [10].

However, the Eager Dynamic Check is not without issue. The problem relates to the manner in which checks are inserted, which we will make precise later. To illustrate the problem, consider [Listing 9](#), a slight variation of [Listing 11](#) in which we replace the program fragment `Plus(Nat(0), y)` with `y`.

λ_{op}

```

handle
  do body ← extrude(Var( $x_N$ )); return Var( $x_N$ ) in return Lam(Var( $x_N$ ), body)
with
  { return( $u$ ) ↦ return  $u$ ;
    extrude( $y, k$ ) ↦ do  $z$  ← return Plus(Nat(0),  $y$ ) in continue  $k z$  }

return Lam(Var( $x_N$ ), Plus(Nat(0), Var( $x_N$ )))

```

Listing 11: A λ_{op} program that may, or may not demonstrate scope extrusion, depending on one's definition. The final result of the program is well-scoped, but not all intermediate results are well-scoped. If scope extrusion is a property of the resulting AST, then this does not display scope extrusion. If, instead, it is a property of intermediate ASTs, then this does display scope extrusion.

While the program in [Listing 9](#) produces an intermediate AST with extruded variables, because of the mechanism for inserting checks, it passes the Eager Dynamic Check. I assert that this behaviour is unintuitive, and exposes too much of the internals to the programmer.

Refined Environment Classifiers

Refined Environment Classifiers are a static check that uses the type system to prevent scope extrusion. Recall that metaprogramming involves the *creation* and *manipulation* of ASTs. Refined environment classifiers prevent scope extrusion by checking:

1. *Created* ASTs are well-scoped
2. *Manipulating* ASTs preserves well-scopedness

We shall first ignore the ability to manipulate ASTs, and consider how to ensure *created* ASTs are well-scoped. What does it mean to be well-scoped? Consider [Figure 2.5](#), the AST of

$$(\lambda f. \lambda x. f x)(\lambda y. y)$$

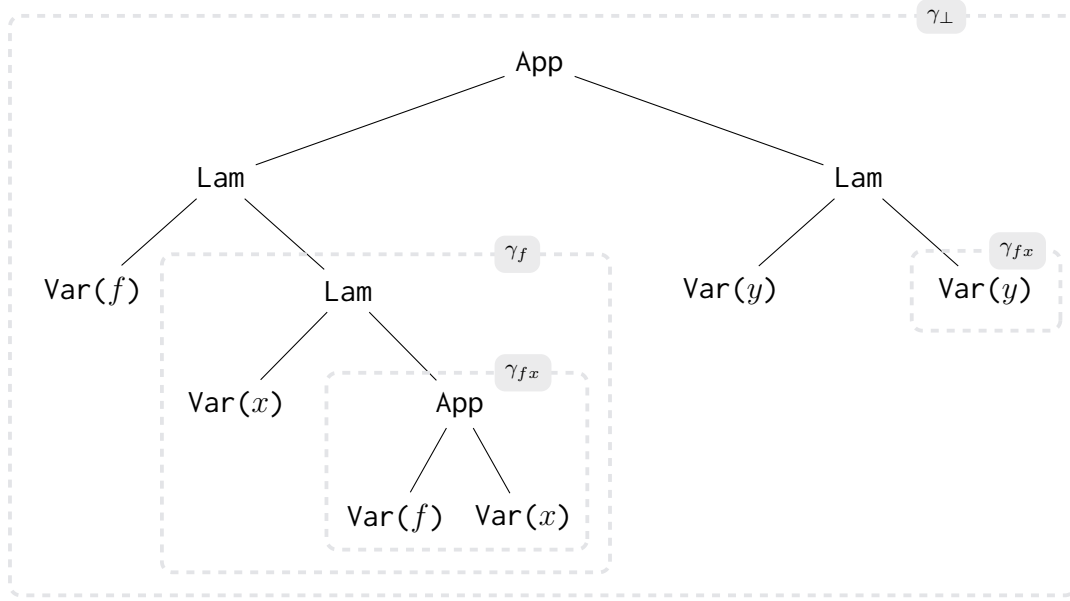


Figure 2.5: The AST of $(\lambda f.\lambda x.fx)(\lambda y.y)$, where each scope is labelled with the corresponding environment classifier.

Informally, a scope represents a set of variables that are permitted to be free. In the example, there are four scopes: one where no variables are free, one where only f is free, one where f and x are free, and one where y is free. An AST is well-scoped *at a scope* if it is well-typed, and where the only free variables are those permitted by the scope.

Refined environment classifiers make this notion precise. Each classifier represents a scope. The AST has four scopes, corresponding to four classifiers:

γ_{\perp} The top-level, where no free variables are permitted

γ_f Only $\text{Var}(f)$ permitted to be free

γ_{fx} Only $\text{Var}(f)$ and $\text{Var}(x)$ permitted to be free

γ_y Only $\text{Var}(y)$ permitted to be free

As every variable binder creates a scope, we may refer to the classifier created by $\text{Var}(\alpha)$, $\text{classifier}(\text{Var}(\alpha))$. For example, $\gamma_{fx} = \text{classifier}(\text{Var}(x))$.

With classifiers, we may be precise about “the variables permitted to be free (within the scope)”. As illustrated by the nesting in [Figure 2.5](#), scopes are related to other scopes. For example, since the scope γ_{fx} is created within scope γ_f , any variable tagged with γ_f may be safely used in the scope γ_{fx} . We say γ_f is compatible with γ_{fx} , and write

$$\gamma_f \sqsubseteq \gamma_{fx}$$

2 Background

The compatibility relation (\sqsubseteq) is a partial order, meaning \sqsubseteq is reflexive, anti-symmetric, and transitive, and identifies a smallest classifier. Reflexivity expresses that $\text{Var}(\alpha)$ may be used within the scope it creates

$$\forall \gamma. \gamma \sqsubseteq \gamma$$

anti-symmetric, since nesting only proceeds in one direction

$$\forall \gamma_1, \gamma_2. \gamma_1 \sqsubseteq \gamma_2 \wedge \gamma_2 \sqsubseteq \gamma_1 \implies \gamma_1 = \gamma_2$$

and transitive, since we should count nestings within nestings

$$\forall \gamma_1, \gamma_2, \gamma_3. \gamma_1 \sqsubseteq \gamma_2 \wedge \gamma_2 \sqsubseteq \gamma_3 \implies \gamma_1 \sqsubseteq \gamma_3$$

γ_\perp acts as the least element of this partial order

$$\forall \gamma. \gamma_\perp \sqsubseteq \gamma$$

Given a classifier γ , we may now define the variables permitted to be free in γ , written $\text{permitted}(\gamma)$

$$\text{permitted}(\gamma) \triangleq \{\text{Var}(\alpha) \mid \text{classifier}(\text{Var}(\alpha)) \sqsubseteq \gamma\}$$

For example, $\text{permitted}(\gamma_{fx}) = \{\text{Var}(f), \text{Var}(x)\}$

We now say that an AST n is well-scoped at type T and scope γ if it is well-typed at T , and all free variables in n are in $\text{permitted}(\gamma)$. We write

$$\Gamma \vdash^\gamma n : T$$

Ensuring that created ASTs are well-scoped is the responsibility of the type system. The key rule is the C-Abs rule, which, **assuming we know that we are creating an AST**, has roughly the following shape²:

$$\begin{array}{c} \text{(C-Abs)} \\ \frac{\gamma \in \Gamma \quad (2) \gamma' \text{fresh} \quad (3) \Gamma, \gamma', \gamma \sqsubseteq \gamma', (x : T_1)^{\gamma'} \vdash^{\gamma'} n : T_2}{(1) \Gamma \vdash^\gamma \lambda x. n : T_1 \longrightarrow T_2} \end{array}$$

The premises and conclusions have been numbered for reference, and many technical details have been simplified for clarity. [Figure 2.6](#) visually depicts the typing rule.

- (1) The goal of the typing rule is to ensure the function is well-scoped at type $T_1 \rightarrow T_2$ and scope γ .
- (2) Since the function introduces a new variable binder, $\text{Var}(x)$, one has to create a new scope. This is achieved by picking a fresh classifier γ' .

²I will revisit this assumption when introducing the type system for my calculus

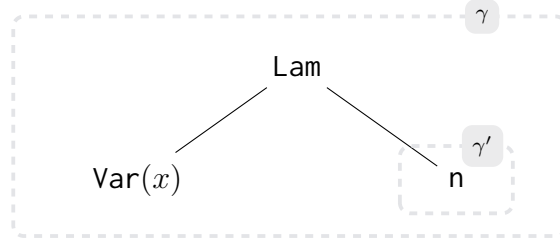
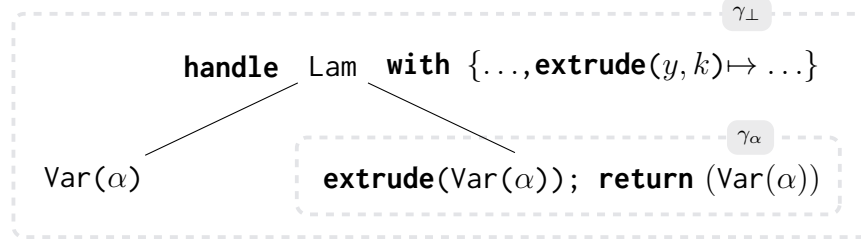


Figure 2.6: Visual depiction of the (C-Abs) typing rule.

Figure 2.7: The “AST” of the scope extrusion example, Listing 8. Notice that in place of AST nodes, we may now have compile-time executable code that *evaluate* to AST nodes.

(3) We record the following:

- (a) Since γ' is created within the scope of γ , $\gamma \sqsubseteq \gamma'$.
- (b) $\text{classifier}(\text{Var}(x)) = \gamma'$ (as a shorthand $(x : T_1)^{\gamma'}$)

With this added knowledge, we ensure that the function body is well-scoped at type T_2 and γ' .

The above example focused on **creating** ASTs, and had no compile-time executable code. We now consider how to maintain well-scopedness while **manipulating** ASTs. We consider the “AST” of the scope extrusion example, Listing 8 (Figure 2.7). Notice that in place of AST nodes, we may now have compile-time executable code that *evaluate* to AST nodes. Thus, both code and AST nodes reside within scopes. We have two classifiers: γ_\perp and γ_α , with $\text{classifier}(\text{Var}(\alpha)) = \gamma_\alpha$.

Key to the prevention of scope extrusion is the typing of handlers and operations, like **extrude**, that manipulate ASTs. Since these rules are complex, we describe them informally. The handle expression

handle e **with** $\{h\}$

is in scope γ_\perp . Therefore, for each operation handled by h , such as **extrude**, the argument to the operation must either not be an AST, or be an AST that is well-scoped at some $\gamma \sqsubseteq \gamma_\perp$. However, $\text{Var}(\alpha)$ is typed at γ_α , and clearly, $\gamma_\alpha \not\sqsubseteq \gamma_\perp$. There is thus no way to type the scope extrusion example in Listing 8.

2 Background

Note that the analysis was independent of the *body* of the handler. Therefore, the examples in Listings 9 and 11 are *also* not well-typed. Perhaps somewhat surprisingly, so too is Listing 12 (which would pass both the Eager and Lazy Dynamic Checks). Refined environment classifiers statically prevent variables ($\text{Var}(\alpha)$) from becoming *available* in program fragments ($\text{op}(y, k) \mapsto \dots$) where, *if misused, might* result in scope extrusion. This is, of course, an over-approximation. It means that the refined environment classifiers check prevents not only both types of scope extrusion, but even more benign examples, such as that in Listing 12.

λ_{op}

```
handle
  do body  $\leftarrow$  extrude(Var( $x_{\mathbb{N}}$ )); return Var( $x_{\mathbb{N}}$ ) in return Lam(Var( $x_{\mathbb{N}}$ ), body)
with
  {return( $u$ )  $\mapsto$  return App( $u$ , Nat(1));
   extrude( $y, k$ )  $\mapsto$  return Nat(0)}

return Nat(0)
```

Listing 12: A λ_{op} program that passes the Eager and Lazy Dynamic Checks, but is not well-typed under the Refined Environment Classifiers type system.

Listing 12 thus illustrates one of the key drawbacks of Refined Environment Classifiers: the check is too stringent, and restricts expressiveness.

Bibliography

- [1] D. Abrahams and A. Gurtovoy. *C++ Template Metaprogramming: Concepts, Tools, and Techniques from Boost and Beyond (C++ in Depth Series)*. Addison-Wesley Professional, 2004. ISBN 0321227255.
- [2] C. Calcagno, E. Moggi, and W. Taha. Closed types as a simple approach to safe imperative multi-stage programming. In U. Montanari, J. D. P. Rolim, and E. Welzl, editors, *Automata, Languages and Programming*, pages 25–36, Berlin, Heidelberg, 2000. Springer Berlin Heidelberg. ISBN 978-3-540-45022-1.
- [3] T.-J. Chiang, J. Yallop, L. White, and N. Xie. Staged compilation with module functors. *Proc. ACM Program. Lang.*, 8(ICFP), Aug. 2024. doi: 10.1145/3674649. URL <https://doi.org/10.1145/3674649>.
- [4] M. Felleisen and D. P. Friedman. Control operators, the secd-machine, and the λ -calculus. In M. Wirsing, editor, *Formal Description of Programming Concepts - III: Proceedings of the IFIP TC 2/WG 2.2 Working Conference on Formal Description of Programming Concepts - III, Ebberup, Denmark, 25-28 August 1986*, pages 193–222. North-Holland, 1987.
- [5] D. Hillerström and S. Lindley. Shallow effect handlers. In S. Ryu, editor, *Programming Languages and Systems*, pages 415–435, Cham, 2018. Springer International Publishing. ISBN 978-3-030-02768-1.
- [6] J. Inoue and W. Taha. Reasoning about multi-stage programs. In H. Seidl, editor, *Programming Languages and Systems*, pages 357–376, Berlin, Heidelberg, 2012. Springer Berlin Heidelberg. ISBN 978-3-642-28869-2.
- [7] K. Isoda, A. Yokoyama, and Y. Kameyama. Type-safe code generation with algebraic effects and handlers. In *Proceedings of the 23rd ACM SIGPLAN International Conference on Generative Programming: Concepts and Experiences, GPCE '24*, page 53–65, New York, NY, USA, 2024. Association for Computing Machinery. ISBN 9798400712111. doi: 10.1145/3689484.3690731. URL <https://doi.org/10.1145/3689484.3690731>.
- [8] Jax-ML. Using grad on vmap on map on function containing sinc results in error. URL <https://github.com/jax-ml/jax/issues/10750>.

- [9] O. Kiselyov. Delimited control in ocaml, abstractly and concretely. *Theoretical Computer Science*, 435:56–76, 2012. ISSN 0304-3975. doi: <https://doi.org/10.1016/j.tcs.2012.02.025>. URL <https://www.sciencedirect.com/science/article/pii/S0304397512001661>. Functional and Logic Programming.
- [10] O. Kiselyov. The design and implementation of ber metaocaml. In M. Codish and E. Sumii, editors, *Functional and Logic Programming*, pages 86–102, Cham, 2014. Springer International Publishing. ISBN 978-3-319-07151-0.
- [11] O. Kiselyov. Generating c: Heterogeneous metaprogramming system description. *Science of Computer Programming*, 231:103015, 2024. ISSN 0167-6423. doi: <https://doi.org/10.1016/j.scico.2023.103015>. URL <https://www.sciencedirect.com/science/article/pii/S0167642323000977>.
- [12] O. Kiselyov, Y. Kameyama, and Y. Sudo. Refined environment classifiers. In A. Igarashi, editor, *Programming Languages and Systems*, pages 271–291, Cham, 2016. Springer International Publishing. ISBN 978-3-319-47958-3.
- [13] G. Mainland. Explicitly heterogeneous metaprogramming with meta-haskell. *SIGPLAN Not.*, 47(9):311–322, Sept. 2012. ISSN 0362-1340. doi: [10.1145/2398856.2364572](https://doi.org/10.1145/2398856.2364572). URL <https://doi.org/10.1145/2398856.2364572>.
- [14] L. Phipps-Costin, A. Rossberg, A. Guha, D. Leijen, D. Hillerström, K. Sivaramakrishnan, M. Pretnar, and S. Lindley. Continuing webassembly with effect handlers. *Proc. ACM Program. Lang.*, 7(OOPSLA2), Oct. 2023. doi: [10.1145/3622814](https://doi.org/10.1145/3622814). URL <https://doi.org/10.1145/3622814>.
- [15] M. Pretnar. An introduction to algebraic effects and handlers. invited tutorial paper. *Electronic Notes in Theoretical Computer Science*, 319:19–35, 2015. ISSN 1571-0661. doi: <https://doi.org/10.1016/j.entcs.2015.12.003>. URL <https://www.sciencedirect.com/science/article/pii/S1571066115000705>. The 31st Conference on the Mathematical Foundations of Programming Semantics (MFPS XXXI).
- [16] H. G. Rice. Classes of recursively enumerable sets and their decision problems. *Transactions of the American Mathematical Society*, 74(2):358–366, 1953. ISSN 00029947, 10886850. URL <http://www.jstor.org/stable/1990888>.
- [17] A. D. Robison. Impact of economics on compiler optimization. In *Proceedings of the 2001 Joint ACM-ISCOPE Conference on Java Grande, JGI '01*,

- page 1–10, New York, NY, USA, 2001. Association for Computing Machinery. ISBN 1581133596. doi: 10.1145/376656.376751. URL <https://doi.org/10.1145/376656.376751>.
- [18] M. Servetto and E. Zucca. A meta-circular language for active libraries. In *Proceedings of the ACM SIGPLAN 2013 Workshop on Partial Evaluation and Program Manipulation, PEPM '13*, page 117–126, New York, NY, USA, 2013. Association for Computing Machinery. ISBN 9781450318426. doi: 10.1145/2426890.2426913. URL <https://doi.org/10.1145/2426890.2426913>.
 - [19] K. Sivaramakrishnan, S. Dolan, L. White, T. Kelly, S. Jaffer, and A. Madhavapeddy. Retrofitting effect handlers onto ocaml. In *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation, PLDI 2021*, page 206–221, New York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450383912. doi: 10.1145/3453483.3454039. URL <https://doi.org/10.1145/3453483.3454039>.
 - [20] L. Tratt. Domain specific language implementation via compile-time meta-programming. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 30(6):1–40, 2008.
 - [21] J. Vandebon, J. G. F. Coutinho, W. Luk, and E. Nurvitadhi. Enhancing high-level synthesis using a meta-programming approach. *IEEE Transactions on Computers*, 70(12):2043–2055, 2021. doi: 10.1109/TC.2021.3096429.
 - [22] H. Wickham. *Advanced R*. Chapman and Hall/CRC, 2019.
 - [23] N. Xie, Y. Cong, K. Ikemori, and D. Leijen. First-class names for effect handlers. *Proc. ACM Program. Lang.*, 6(OOPSLA2), Oct. 2022. doi: 10.1145/3563289. URL <https://doi.org/10.1145/3563289>.
 - [24] N. Xie, L. White, O. Nicole, and J. Yallop. Macocaml: Staging composable and compilable macros. *Proc. ACM Program. Lang.*, 7(ICFP), Aug. 2023. doi: 10.1145/3607851. URL <https://doi.org/10.1145/3607851>.
 - [25] J. Yallop, N. Xie, and N. Krishnaswami. flap: A deterministic parser with fused lexing. *Proc. ACM Program. Lang.*, 7(PLDI), June 2023. doi: 10.1145/3591269. URL <https://doi.org/10.1145/3591269>.