

# Implementation of type inference for a programming language with algebraic effects

(Inferencja typów dla języka programowania z efektami algebraicznymi)

Dominik Gulczyński

Praca licencjacka

**Promotor:** dr Filip Sieczkowski

Uniwersytet Wrocławski  
Wydział Matematyki i Informatyki  
Instytut Informatyki

9 września 2020



## Abstract

Algebraic effects and handlers are a novel and powerful programming construct that mediates between manageable side effects and meaningful type safety.

Based on the recent works in this field, we proposed a calculus and a type-and-effect system that enabled us to implement ML-style polymorphism for decidable type inference. The inference algorithm we present has proven to be a non-trivial refinement of a classical, syntax-oriented algorithm  $\mathbb{W}$ , due to the introduced subtyping relation, additional calculus constructs and extended typing rules. Theoretical description of our algorithm is accompanied by practical implementation.

---

Handlery i efektami algebraicznymi to nowoczesna i skuteczna konstrukcja programistyczna, pośrednicząca pomiędzy wykorzystaniem efektów ubocznych a bezpiecznym systemem typów.

Na podstawie ostatnich prac w tej dziedzinie, zaproponowaliśmy rachunek oraz system typów i efektów, który umożliwił nam wdrożenie polimorfizmu w stylu ML dla rozstrzygalnej inferencji typów. Przedstawiony przez nas algorytm inferencji jest nietrywialnym rozszerzeniem klasycznego algorytmu syntaktycznej inferencji typów  $\mathbb{W}$ , ze względu na wprowadzoną relację podtypowania, dodatkowe konstrukcje języka i rozszerzone reguły typowania. Teoretycznemu opisowi naszego algorytmu towarzyszy praktyczna implementacja.



# Contents

<b>1</b>	<b>Introduction</b>	<b>7</b>
<b>2</b>	<b>Calculus</b>	<b>9</b>
<b>3</b>	<b>Type system</b>	<b>11</b>
3.1	Subtyping . . . . .	13
3.2	Principal type . . . . .	14
<b>4</b>	<b>Inference algorithm</b>	<b>17</b>
4.1	Remarks on effect unification . . . . .	18
4.2	Generating constraints . . . . .	18
4.3	Solving constraints . . . . .	21
4.3.1	Simple constraints . . . . .	21
4.3.2	Interesting constraints . . . . .	22
4.4	Example . . . . .	24
<b>5</b>	<b>Implementation</b>	<b>25</b>
5.1	Project structure . . . . .	25
5.2	Tutorial . . . . .	27
<b>6</b>	<b>Conclusions and future work</b>	<b>29</b>
	<b>Bibliography</b>	<b>31</b>



# Chapter 1

## Introduction

Programming languages that enjoy *static type safety* require that each variable and expression is given a *type*. Assigning values to variables of different types usually results in a compile error.

Type inference method was designed to aid the process of *type-checking* and to ease the lives of programmers. It is a process of *deducing* the type of an expression, by the analysis of its structure and the environment in which we carry out the inference process. Many programming languages use it to free programmers from the burden of writing explicit type annotations, which makes code easier to read and modify. For example, in programming language C# programmer may use the `var` keyword to omit the type in variable declaration.

A *complete* type inference system can deduce a type for every expression without the need for any type annotations. Most popular of such type systems, called *Hindley-Milner*, is based on the works of J. Roger Hindley[1], Robin Milner, who presented *the* type inference algorithm W[2], and Luis Damas[3]. Their work is a corner stone of type theory for programming languages and the basis for type systems and type inference present in most of the statically typed functional programming languages.

Type systems like System-F or Hindley-Milner's, based on the *lambda calculus* are great mathematical models of describing pure computations and types of pure terms concisely. Yet, a lot of real-life programs heavily use side effects, be it writing to memory, performing I/O actions, or mutating state, and those systems cannot express them. There are different ways of describing side effects in both scientific literature and real world implementations of programming languages. Some languages do not restrict side effects at all, like OCaml, while others have their unique ways of expressing side effects, like monadic actions in Haskell. Finally there are so called type-and-effect systems. One of the earliest works by Talpin and Jouvelot[4] in this research area describe studying effects as an additional way of reasoning about programs: while the type of an expression tells us about the *results* of computation,

the effects describe *how* it computes.

The particular system and calculus (which will from now on be referred to as *the original calculus*) presented by Biernacki et al[5] uses *algebraic effects* to describe programming with computational effects. An effect is defined by a set of available operations (defined by its *signature*) and a handler that provides its implementation. Operators are meaningless on their own, acting much like ordinary functions calling the handler's body, while the handler itself defines the way operation calls execute.

One could think of effects as a generalization of exceptions: calling an operator corresponds to throwing an exception, and surround an expression with a handler, corresponds to *try {...} catch {...}* construct known in some form in many mainstream languages. But, in most of those languages exceptions are not *resumable*, meaning that once we *leave* the inner expression of a handler, by calling one of its operators, we cannot continue the computation of handled expression. In *the original calculus*, the programmer is free to use the *continuation* however they like when programming handlers.

We transcribe *the original calculus* to the world of ML programming by presenting modified type-and-effect system and corresponding calculus, together with type inference algorithm for it.



## Chapter 2

# Calculus

The calculus we present and analyze in this work is a subset of the one described by Biernacki et al[5], adjusted to match the style of *ML-the-calculus* and *ML-the-type-system*[6]. It is an extension of standard *call-by-value lambda calculus with let* by effect handlers and operators.

$x, \dots$	(term variables)
$\alpha, \dots$	(quantified variables)
$t, \dots$	(type variables)
$\epsilon, \dots$	(effect variables)
$a, \dots$	(instance variables)
$\tau ::= \alpha \mid t \mid \mathbf{Unit} \mid \mathbf{Int} \mid \tau \rightarrow_{\epsilon} \tau \mid \forall(a : \sigma). \tau$	(types)
$\pi ::= \forall \alpha. \pi \mid \tau$	(type schemes)
$\varepsilon ::= \alpha \mid \epsilon \mid a \mid \varepsilon \cdot \varepsilon$	(effects)
$\sigma ::= \mathbf{Error} \mid \mathbf{State} \tau$	(signatures)
$e ::= x \mid () \mid n \mid \lambda x. e \mid \mathbf{fun} \ f, x. e \mid e \ e \mid \mathbf{let} \ x = e \ \mathbf{in} \ e$ $\quad \mid \lambda(a : \sigma). e \mid e \ a \mid op_a \ e \mid \mathbf{handle}_a \ e \{h; \mathbf{return} \ x.e\}$	(terms)
$op ::= \mathbf{Raise} \mid \mathbf{Get} \mid \mathbf{Put}$	(operators)
$h ::= [(\mathbf{Raise}, x, k.e)] \mid [(\mathbf{Get}, x, k.e); (\mathbf{Put}, x, k.e)]$	(handlers)

Terms (or expressions) are given by:

- variables, bound by abstractions, let-expressions, handlers, or environment  $\Gamma$ ,
- constants:  $()$  and  $n \in \mathbb{C}$ ,
- abstractions: anonymous functions  $\lambda x. e$  with argument  $x$  and body  $e$  and recursive functions denoted  $\mathbf{fun} \ f, x. e$ ,
- instance abstraction:  $\lambda(a : \sigma). e$ , that lets programmers write code unspecific to certain effect *instance*, but capable of working with any instance of specified

signature,

- applications:  $e_1 e_2$  and  $e a$ , for applying arguments to respective abstractions,
- let-construct: **let**  $x = e_1$  **in**  $e_2$ , which first evaluates body of  $e_1$ , and bounds its value to variable  $x$  in expression  $e_2$ ,
- operation calls:  $op_a e$  calling the  $op$  operator handler of instance  $a$  with value of  $e$ ,
- handlers: **handle** <sub>$a$</sub>   $e \{h; \text{return } x.e_r\}$  of instance  $a$ , which provides meaning to operators: calling  $op_a e_x$ , executes body  $e_{op}$  of a construct  $(op, x, k.e_{op})$  defined in handler  $h$ .

Instead of allowing effects of arbitrary signatures, we limited them to instances of either **Error** or **State**  $\tau$ . Arbitrary signatures could be dealt with in similar fashion as *ADTs* (algebraic data types) and are *orthogonal* to type inference.

For the *complete type inference* we need to limit the system to *prenex* polymorphism and thus we differentiate between *monomorphic* types and *polymorphic* type schemes (which are types proceeded by  $\forall$ -quantifier followed by  $\alpha$  variables). Accordingly we have dropped the type-lambdas ( $\Lambda$ ) from calculus in favor of the polymorphic let-construct.

Notice how each expression, type and effect of our calculus can be expressed in terms of *original calculus*. The semantics of our calculus strictly follow the rules defined in the original paper *Binders by Day, Labels by Night* by Biernacki et al[5].

## Chapter 3

# Type system

We have tweaked the type-and-effect system of *the original calculus* to match the restrictions required by the classical methods of type inference. Thus, we introduce *unification variables* for both types and effects (written  $t$  and  $\epsilon$ , respectively) to denote “yet to be determined” types and effects. During the type inference we will *generate* these unification variables, and deduce rigid types for them, as described in Chapter 4.

Judgement  $\Gamma, \Theta \vdash e : \tau/\epsilon$  states that in environments  $\Gamma$  (assigning types to variables), and  $\Theta$  (assigning signatures to instances), term  $e$  *inhabits* syntactical type and effect  $\tau/\epsilon$  (which means that computing  $e$  would yield a value of type  $\tau$  and possibly cause effect  $\epsilon$ ).

$$\frac{}{\Gamma; \Theta \vdash () : \mathbf{Unit}/\iota} \quad \frac{}{\Gamma; \Theta \vdash n : \mathbf{Int}/\iota}$$

There are two *base types* **Unit** and **Int** for constants.

$$\frac{(x : \pi) \in \Gamma \quad \pi[\vec{\tau}_\alpha/\vec{\alpha}] = \tau}{\Gamma; \Theta \vdash x : \tau/\iota}$$

Judgement for variables follows the usual let-polymorphism typing, where variables bound by let clauses are generalized and need to be instantiated. Variables do not cause effects as only the value is assigned to them, while the effects caused by their computation (if any occur) are bound to the term which introduced that variable.

$$\frac{\Gamma, (x : \tau_1); \Theta \vdash e : \tau_2/\epsilon}{\Gamma; \Theta \vdash \lambda x. e : \tau_1 \rightarrow_\epsilon \tau_2/\iota} \quad \frac{\Gamma, (f : \tau_1 \rightarrow_\epsilon t_2), (x : \tau_1); \Theta \vdash e : \tau_2/\epsilon}{\Gamma; \Theta \vdash \mathbf{fun} \ f, x. e : \tau_1 \rightarrow_\epsilon \tau_2/\iota}$$

The *type constructor*  $\rightarrow$  is used to type abstractions, where type  $\tau_1 \rightarrow_\epsilon \tau_2$  is given to functions that when applied with input of type  $\tau_1$ , computing the function’s body  $e$  produces some output of type  $\tau_2$ , possible causing effect  $\epsilon$

$$\frac{\Gamma; \Theta \vdash e_1 : \tau_2 \rightarrow_\varepsilon \tau/\varepsilon_1 \quad \Gamma; \Theta \vdash e_2 : \tau_2/\varepsilon_2}{\Gamma; \Theta \vdash e_1 e_2 : \tau/\varepsilon_1 \cdot \varepsilon \cdot \varepsilon_2}$$

Accordingly, the effect of application combines effects of: computing the left-hand side term, effect that “hangs” on it’s arrow type, and the effect of computing the right-hand side term.

$$\frac{\Gamma; \Theta \vdash e_1 : \tau_1/\iota \quad \text{gen}(\Gamma, \tau_1) = \pi \quad \Gamma, (x : \pi); \Theta \vdash e_2 : \tau/\varepsilon}{\Gamma; \Theta \vdash \text{let } x = e_1 \text{ in } e_2 : \tau/\varepsilon}$$

$$\frac{\Gamma; \Theta \vdash e_1 : \tau_1/\varepsilon_1 \quad \varepsilon_1 \neq \iota \quad \Gamma, (x : \tau_1); \Theta \vdash e_2 : \tau/\varepsilon}{\Gamma; \Theta \vdash \text{let } x = e_1 \text{ in } e_2 : \tau/\varepsilon}$$

As usually in let-polymorphism schemes, we *generalize* the type derived for  $e_1$  before we add it to the environment in which we derive type for  $e_2$ . This is achieved by the function  $\text{gen}(\Gamma, \tau) = \forall \vec{\alpha}. \tau$  where  $\vec{\alpha} = \text{freevars}(\tau) \setminus \text{freevars}(\Gamma)$ . We restrict generalization to only *pure* terms, i.e. such that their computation would cause no effects.

$$\frac{\Gamma; \Theta, (a : \sigma) \vdash e : \tau/\iota}{\Gamma; \Theta \vdash \lambda(a : \sigma). e : \forall(a : \sigma). \tau/\iota} \quad \frac{\Gamma; \Theta \vdash e : \forall(a : \sigma). \tau/\iota \quad (b : \sigma) \in \Theta}{\Gamma; \Theta \vdash e b : \tau[b/a]/\iota}$$

For handling operators of distinct occurrences of the same signature, i.e many cells of memory via **State**  $\tau$ , we differentiate them by assigning them different instances. *Instance abstractions* allow us to write code that works not only with one instance of particular signature, but with any that is supplied using *instance application* of instances bound by handlers or other instance lambdas.

$$\frac{\Theta \vdash op_a : \tau_e \Rightarrow \tau \quad \Gamma; \Theta \vdash e : \tau_e/\varepsilon}{\Gamma; \Theta \vdash op_a e : \tau/a \cdot \varepsilon}$$

Operators of type  $\tau_1 \Rightarrow \tau_2$  expect arguments of type  $\tau_1$  and return value of type  $\tau_2$ . Type of  $e$  must match the operator type. Effect of computing  $e$  is extended by  $\{a\}$ .

$$\frac{\Gamma; \Theta \vdash h : \sigma \triangleright \tau/\varepsilon \quad \Gamma; \Theta, (a : \sigma) \vdash e : \tau'/a \cdot \varepsilon \quad \Gamma, (x : \tau'); \Theta \vdash e_r : \tau/\varepsilon}{\Gamma; \Theta \vdash \text{handle}_a e \{h; \text{return } x.e_r\} : \tau/\varepsilon}$$

For type-checking handler body, there’s a separate relation  $\vdash : \triangleright$ , explained below. Types of all *execution paths* of the handler must match, whether the value returns from within handler body or via return expression  $e_r$  (where value of evaluating  $e$  gets bound to  $x$ ).

Finally, we allow every type and effect to *grow* as needed:

$$\frac{\Gamma; \Theta \vdash e : \tau'/\varepsilon' \quad \tau' <: \tau \quad \varepsilon' <: \varepsilon}{\Gamma; \Theta \vdash e : \tau/\varepsilon}$$

$$\begin{array}{c}
\frac{\Gamma, (x : \mathbf{Unit}), (k : \tau' \rightarrow_{\varepsilon} \tau); \Theta \vdash e' : \tau/\varepsilon}{\Gamma; \Theta \vdash [(\mathbf{Raise}, x, k.e')] : \mathbf{Error} \triangleright \tau/\varepsilon} \\
\\
\frac{\Gamma, (x : \mathbf{Unit}), (k : \tau' \rightarrow_{\varepsilon} \tau); \Theta \vdash e' : \tau/\varepsilon \quad \Gamma, (x : \tau'), (k : \mathbf{Unit} \rightarrow_{\varepsilon} \tau'); \Theta \vdash e'' : \tau/\varepsilon}{\Gamma; \Theta \vdash [(\mathbf{Get}, x, k.e'); (\mathbf{Put}, x, k.e'')] : \mathbf{State} \tau' \triangleright \tau/\varepsilon}
\end{array}$$

Supplying the *continuation*  $k$  with some value  $v$  continues evaluation of the expression surrounded by the handler, with  $v$  substituted in place of operation call, thus in the body of operator handler of type  $\tau_1 \Rightarrow \tau_2$ ,  $k$  is given type  $\tau_2 \rightarrow_{\varepsilon} \tau$  and  $x$  is given  $\tau_2$ . In a sense  $k$  acts just like an ordinary function, and programmer may use it in many different ways or not even use it all, returning a value straight from the handler code. However they choose to do so, the type-and-effect of handlers expression must match the type of the whole expression in which the handler was used.

$$\frac{(a : \mathbf{Error}) \in \Theta}{\Theta \vdash \mathbf{Raise}_a : \mathbf{Unit} \Rightarrow \tau} \quad \frac{(a : \mathbf{State} \tau) \in \Theta}{\Theta \vdash \mathbf{Put}_a : \mathbf{Unit} \Rightarrow \tau} \quad \frac{(a : \mathbf{State} \tau) \in \Theta}{\Theta \vdash \mathbf{Get}_a : \tau \Rightarrow \mathbf{Unit}}$$

Relation  $\Theta \vdash op_a : \tau_1 \Rightarrow \tau_2$  simply finds the signature  $\sigma$  bound to instance  $a$  in environment  $\Theta$  and returns the operator's type.

### 3.1 Subtyping

$$\frac{}{\tau <: \tau} \quad \frac{\tau'_1 <: \tau_1 \quad \varepsilon <: \varepsilon' \quad \tau_2 <: \tau'_2}{\tau_1 \rightarrow_{\varepsilon} \tau_2 <: \tau'_1 \rightarrow_{\varepsilon'} \tau'_2} \quad \frac{}{\varepsilon <: \varepsilon} \quad \frac{\varepsilon <: \varepsilon'}{\varepsilon <: \varepsilon' \cdot \varepsilon''}$$

The subtyping rule we propose is *structural*, meaning that only types of *matching shape* are related. However, while leaves containing types must be equal, we allow effects to differ as long as they are related. Notice that the  $\rightarrow$  is contravariant to subtyping relation (in the premise of rule for  $\rightarrow$ , the order of relation on  $\tau_1$  and  $\tau'_1$  is reversed).

Although the effect constructor  $\cdot$  we used to define effects in Chapter 1 forms binary trees, we implicitly use an equivalence relation on effects that interprets these trees as finite sets. To be precise, we treat effects  $\varepsilon_1$  and  $\varepsilon_2$  as equivalent as long as each leaf in the tree of  $\varepsilon_1$  also appears somewhere in the tree of  $\varepsilon_2$ . For example, we have  $a \equiv a \cdot a$  and  $b \cdot a \cdot b \equiv a \cdot b$ .

The subtyping plays a vital role in usability of the calculus. Consider a term  $f$ , some higher order function that does some calculation when applied with some arithmetical function.

Consider a term  $f$ , that takes as an arithmetical function as argument, which is allowed to cause no effects other than failing to compute via effect  $e$ .

$$\emptyset, (e : \mathbf{Error}) \vdash f : (\mathbf{Int} \rightarrow_e \mathbf{Int}) \rightarrow_e \mathbf{Int}/\iota$$

Thanks to the subtyping relation, we could apply some pure function to  $f$  and the application would still type-check. On the other hand, it would be undesirable if we could supply a term that expects a pure function with an effectful one. Clearly this property gives us flexibility, while keeping effects under control.

## 3.2 Principal type

In ML type system, the *principal type* property states that there exists a *most general* type for any correct program[3]. A type scheme  $\pi$  is called *principal* if any other type that could be given to  $e$  is an *instantiation* of it.

For example, in ML calculus, term  $e = \lambda x. x$  could be given type  $\mathbf{Unit} \rightarrow \mathbf{Unit}$  or  $(\mathbf{Int} \rightarrow \mathbf{Int}) \rightarrow \mathbf{Int} \rightarrow \mathbf{Int}$ , but clearly any correct type we could think of would not be more general than  $\forall \alpha. \alpha \rightarrow \alpha$ , which is in fact, the principal type of  $e$ .

### Definition 3.2.1. Principal Type

In our type system,  $\pi$  is a principal type of  $e$  in environments  $\Gamma, \Theta$  iff

$$\forall \vec{\tau}. \Gamma; \Theta \vdash e : \pi[\vec{\tau}/\vec{\alpha}]/\iota \wedge \forall \pi'. \exists \vec{\tau}'. \Gamma; \Theta \vdash e : \pi'[\vec{\tau}'/\vec{\alpha}']/\iota \wedge \pi[\vec{\tau}/\vec{\alpha}] <: \pi'[\vec{\tau}'/\vec{\alpha}']$$

We talk about this property more in Chapter 6. For now we will give an example of how to approximate a principal type. Consider the two type schemes that could be assigned to a term  $compose = \lambda f. \lambda g. \lambda x. f(g x)$ , representing function composition:

$$\begin{aligned} \pi_1 &::= \forall \alpha, \beta. (\tau_b \rightarrow_\alpha \tau_c) \rightarrow (\tau_a \rightarrow_\beta \tau_b) \rightarrow \tau_a \rightarrow_{\alpha \cdot \beta} \tau_c \\ \pi_2 &::= \forall \gamma. (\tau_b \rightarrow_\gamma \tau_c) \rightarrow (\tau_a \rightarrow_\gamma \tau_b) \rightarrow \tau_a \rightarrow_\gamma \tau_c \end{aligned}$$

At first glance, it may look like  $\pi_1$  is the “correct” type for  $compose$ , as it seems more natural, i.e. given functions  $f$  causing effect  $\varepsilon_f$ , and  $g$  causing  $\varepsilon_g$ ,  $compose f g$  is a function that would apply them both, so it clearly must be causing effect  $\varepsilon_f \cdot \varepsilon_g$ .

While this reasoning is sound, we are interested in deriving the most concise type possible. Let’s see if there exists  $\varepsilon_h$  such that  $\pi_1$  instantiated with arbitrary effects  $\varepsilon_f, \varepsilon_g$  subtypes  $\pi_2$  instantiated with  $\varepsilon_h$ :

$$\begin{aligned} \pi_1[\varepsilon_f, \varepsilon_g/\alpha, \beta] &<: \pi_2[\varepsilon_h/\gamma] \\ &\iff \\ (\tau_b \rightarrow_{\varepsilon_f} \tau_c) \rightarrow (\tau_a \rightarrow_{\varepsilon_g} \tau_b) \rightarrow \tau_a \rightarrow_{\varepsilon_f \cdot \varepsilon_g} \tau_c &<: (\tau_b \rightarrow_{\varepsilon_h} \tau_c) \rightarrow (\tau_a \rightarrow_{\varepsilon_h} \tau_b) \rightarrow \tau_a \rightarrow_{\varepsilon_h} \tau_c \\ &\iff \\ \varepsilon_h &<: \varepsilon_f \wedge \varepsilon_h <: \varepsilon_g \wedge \varepsilon_f \cdot \varepsilon_g <: \varepsilon_h \end{aligned}$$

Clearly,  $\pi_1[\varepsilon_f, \varepsilon_g/\alpha, \beta] <: \pi_2[\varepsilon_h/\gamma]$  does not hold for  $\varepsilon_f$  and  $\varepsilon_g$  other than  $\iota$ , thus  $\pi_1$  cannot be a principal type of  $e$ . On the other hand, if we were to check if for

arbitrary  $\varepsilon_h$  there exist  $\varepsilon_f$  and  $\varepsilon_g$  such that  $\pi_2[\varepsilon_h/\gamma] <: \pi_1[\varepsilon_f, \varepsilon_g/\alpha, \beta]$ , we would need to find witnesses for such formula:

$$\varepsilon_f <: \varepsilon_h \wedge \varepsilon_g <: \varepsilon_h \wedge \varepsilon_h <: \varepsilon_f \cdot \varepsilon_g$$

Clearly if we choose  $\varepsilon_f = \varepsilon_g = \varepsilon_h$ , it is satisfied, which means that  $\pi_2$  is indeed more general  $\pi_1$ . We designed our inference algorithm with this intuition in mind and  $\pi_2$  is the result that our implementation actually infers, as we show in Chapter 4.





## Chapter 4

# Inference algorithm

The algorithm we present loosely follows original algorithm  $\mathbb{W}$  and executes in two distinct phases:

1. Constraint gathering: the algorithm traverses the structure of given expression, generating sub-expressions' types, effects and constraints,
2. Constraint solving: the algorithm builds a substitution that resolves gathered constraints, in respect to the laws of type checking.

In practice, these phases are often interleaved, as even in pure *let-polymorphism* type inference when handling expression **let**  $x = e_1$  **in**  $e_2$  we need to solve constraints regarding the inferred type of  $e_1$ , so we can *generalize* it before adding it to the environment for inferring  $e_2$ .

$$C ::= \emptyset \mid \{\tau <: \tau\} \mid \{\varepsilon <: \varepsilon\} \mid C \cup C$$

$$https : //www.overleaf.com/project/5f3c5361269a720001d5cf2dS ::= id \mid [t \mapsto \tau] \mid [\epsilon \mapsto \varepsilon] \mid S \circ S$$

Finally, the key difference in our algorithm arises. In the standard HM algorithm and its derivatives that do not support subtyping, the constraints are rather of form  $\{x = y\}$ . Such constraints are *simple*, as they undoubtedly require  $x$  and  $y$  to be unified. Instead, we have constraints of type  $\{x <: y\}$ . Presence of subtyping relation in type inference process proves to be problematic, as described by francois in his work. well well well

For base types, it means the same, but for effects is hella more complicated and ya know.

are a way of describing necessary conditions that need to be satisfied for the expression and inferred type to type-check.

A substitution is a mapping from type and effect *unification variables* to inferred

types or effects, respectively. We will write  $S[t \mapsto \tau]$  for  $[t \mapsto \tau] \circ S$ . Substitution  $id$  is simply an identity function.

## 4.1 Remarks on effect unification

With combining effect *unification variables* and constraint solving, a problem arises. Consider constraint  $a <: \epsilon_1 \cdot \epsilon_2$ , for some instance variable  $a$  and some effect unification variables  $\epsilon_1, \epsilon_2$ . To resolve such constraint we have a few viable options:

1. Expand  $\epsilon_1$  with  $a$ .
2. Expand  $\epsilon_2$  with  $a$ .
3. Expand both  $\epsilon_1$  and  $\epsilon_2$  with  $a$ .

But how would we choose one over the other? Maybe one of those makes the program ill-typed, while the others do not? What if there's more than two unification variables? Clearly such constraints are undesirable.

To tackle this problem, in our algorithm we permit no more than one effect *unification variable* or *quantified variable* in effects. Thus the effects are defined differently than in Chapter 1:

$$\begin{aligned} \varepsilon &::= I \mid I * \alpha \mid I * \epsilon && \text{(effects)} \\ I &::= \iota \mid \{a\} \mid I \cup I \mid I \setminus I && \text{(sets of instances)} \end{aligned}$$

So an effect is either just a finite set of instances, or a union of one with either effect *unification variable* or *generalized* effect variable. If we think about effects as sets, then the subtyping relation of effects simply boils down to set inclusion. Notice how every effect defined by the new grammar is expressible in the previous one as well.

## 4.2 Generating constraints

As in the *syntax-directed* algorithm  $\mathbb{W}$ , to infer the type for a given term, we build it by working *bottom-up* from the leaves through the whole expression tree. To this end, we have defined judgement  $\Gamma; \Theta \vdash_{\text{Gen}} e : \tau/\varepsilon \rightsquigarrow C; S$  that states in the the premise what conditions need to be satisfied for deducing type, and under what constraints  $C$  and substitution  $S$  shall we interpret it.

We chose to present  $\vdash_{\text{Gen}}$  typing rules as close as possible to the practical inference algorithm. Hence, we heavily use *unification variables* (denoted  $t$  for types and  $\epsilon$  for effects), and the judgement should be thought of as a semi-formal description of type and constraints generating process (working bottom-up) rather than a type checker ( $\vdash$  working top-down).

It is important that we “return” not only constraints, but also a substitution, because some constraints may have been already resolved in one sub-tree of the term, effectively changing the environment, and we need to take it into account while inferring the other sub-trees. We abstract solving constraints  $C$  in environment  $\Gamma$  (under substitution  $S$ ) to a high-level functions *solve* and *solve'* which return a reduced set of constraints  $C'$  and a new substitution  $S'$ .

**Definition 4.2.1.** Syntactic type soundness

$$\begin{aligned} \Gamma; \Theta \vdash_{\mathbf{Gen}} e : \tau/\varepsilon \rightsquigarrow C; S &\implies \\ \text{solve}(\Gamma, \tau/\varepsilon, C, S) = \emptyset; S' &\implies \\ S'\Gamma; S'\Theta \vdash e : S'\tau/S'\varepsilon \end{aligned}$$

We would say that an inference algorithm enjoys *syntactic type soundness* if the type generated by it (after we have solved all the constraints and built the substitution) for the given term checks by the syntactic rules we have defined in chapter 3. Proof of *soundness* our algorithm is left for future work.

$$\begin{array}{c} \frac{}{\Gamma; \Theta \vdash_{\mathbf{Gen}} () : \mathbf{Unit}/\iota \rightsquigarrow \emptyset; id} \quad \frac{}{\Gamma; \Theta \vdash_{\mathbf{Gen}} n : \mathbf{Int}/\iota \rightsquigarrow \emptyset; id} \\[1em] \frac{(x : \pi) \in \Gamma \quad \text{inst}(\pi) = \tau}{\Gamma; \Theta \vdash_{\mathbf{Gen}} x : \tau/\iota \rightsquigarrow \emptyset; id} \end{array}$$

As environment  $\Gamma$  contains type schemes rather than types, before we return the type of  $x$ , we switch out every variable quantified in  $\pi$  to a fresh unification variable via function *inst*.

$$\begin{array}{c} \frac{\Gamma, (x : t); \Theta \vdash_{\mathbf{Gen}} e : \tau/\varepsilon \rightsquigarrow C; S \quad \text{fresh}(t)}{\Gamma; \Theta \vdash_{\mathbf{Gen}} \lambda x. e : t \rightarrow_{\varepsilon} \tau/\iota \rightsquigarrow C; S} \\[1em] \frac{\begin{array}{ccc} \text{fresh}(t_1) & \text{fresh}(\epsilon) & \text{fresh}(t_2) \\ \Gamma, (f : t_1 \rightarrow_{\epsilon} t_2), (x : t_1); \Theta \vdash_{\mathbf{Gen}} e : \tau/\varepsilon \rightsquigarrow C; S \end{array}}{\Gamma; \Theta \vdash_{\mathbf{Gen}} \mathbf{fun} f, x. e : t_1 \rightarrow_{\epsilon} t_2/\iota \rightsquigarrow C \cup \{\varepsilon <: \epsilon, \tau <: t_2\}; S} \end{array}$$

For recursive functions, we must ensure that the actual type-and-effect of function's body  $e$  subtypes one we have generated for  $f$ .

$$\frac{\begin{array}{ccc} \Gamma; \Theta \vdash_{\mathbf{Gen}} e_1 : \tau_1/\varepsilon_1 \rightsquigarrow C_1; S_1 & S_1\Gamma; S_1\Theta \vdash_{\mathbf{Gen}} e_2 : \tau_2/\varepsilon_2 \rightsquigarrow C_2; S_2 \\ \text{fresh}(t) & \text{fresh}(\epsilon) & C = C_1 \cup C_2 \cup \{\varepsilon_1 <: \epsilon, \varepsilon_2 <: \epsilon, \tau_1 <: \tau_2 \rightarrow_{\epsilon} t\} \end{array}}{\Gamma; \Theta \vdash_{\mathbf{Gen}} e_1 e_2 : t/\epsilon \rightsquigarrow C; S_2 \circ S_1}$$

As mentioned earlier, we use the generated substitution  $S_1$  to update the environment for handling  $e_2$ . Additional constraints ensure that the fresh effect variable we return is sybtyped by the effects of computing  $e_1$ ,  $e_2$ , and the effect hung on the  $\rightarrow$  type of  $e_1$ .

$$\begin{array}{c}
\frac{\Gamma; \Theta \vdash_{\mathbf{Gen}} e_1 : \tau_1/\varepsilon_1 \rightsquigarrow C_1; S_1 \quad \text{solve}(\Gamma, \tau_1/\varepsilon_1, C_1, S_1) = C; S \quad S\varepsilon_1 = \iota \quad \text{gen}(S\Gamma, S\tau_1) = \pi \quad S\Gamma, (x : \pi); S\Theta \vdash_{\mathbf{Gen}} e_2 : \tau_2/\varepsilon_2 \rightsquigarrow C_2; S_2}{\Gamma; \Theta \vdash_{\mathbf{Gen}} \text{let } x = e_1 \text{ in } e_2 : \tau_2/\varepsilon_2 \rightsquigarrow C \cup C_2; S_2 \circ S} \\
\\
\frac{\Gamma; \Theta \vdash_{\mathbf{Gen}} e_1 : \tau_1/\varepsilon_1 \rightsquigarrow C_1; S_1 \quad \text{solve}(\Gamma, \tau_1/\varepsilon_1, C_1, S_1) = C; S \quad S\varepsilon_1 \neq \iota \quad S\tau_1 = \tau \quad S\Gamma, (x : \tau); S\Theta \vdash_{\mathbf{Gen}} e_2 : \tau_2/\varepsilon_2 \rightsquigarrow C_2; S_2 \quad \text{fresh}(\epsilon)}{\Gamma; \Theta \vdash_{\mathbf{Gen}} \text{let } x = e_1 \text{ in } e_2 : \tau_2/\epsilon \rightsquigarrow C \cup C_2 \cup \{\varepsilon_1 <: \epsilon, \varepsilon_2 <: \epsilon\}; S_2 \circ S}
\end{array}$$

Before we can handle  $e_2$ , we need to check if type generated for  $e_1$  can be generalized. To this end, we deploy the *solve* function, and if the substitution returned ensured us that computing of  $e_1$  would cause no effects, we can safely generalize it via *gen* function, otherwise we take its effect into account, adding constraints regarding fresh variable  $\epsilon$ . Again,  $S$  is applied to environments for handling  $e_2$ .

$$\begin{array}{c}
\frac{\Gamma; \Theta, (a : \sigma) \vdash_{\mathbf{Gen}} e : \tau/\varepsilon \rightsquigarrow C; S}{\Gamma; \Theta \vdash_{\mathbf{Gen}} \lambda(a : \sigma). e : \forall(a : \sigma). \tau/\iota \rightsquigarrow C \cup \{\varepsilon <: \iota\}; S} \\
\\
\frac{\Gamma; \Theta \vdash_{\mathbf{Gen}} e : \tau'/\varepsilon \rightsquigarrow C'; S' \quad \text{solve}'(\Gamma, C', S') = C; S \quad S\tau' = \forall(a : \sigma). \tau \quad S\varepsilon = \iota \quad (b : \sigma) \in S\Theta}{\Gamma; \Theta \vdash_{\mathbf{Gen}} e b : \tau[b/a]/\iota \rightsquigarrow C; S}
\end{array}$$

For instance abstraction, we have to make sure that the expression underneath it is pure. With instance application, we need to *know* that the type for  $e$  is quantified by some instance  $a$ , so we use the weaker *solve'* function to find out its type, before substituting  $a$  for  $b$  in it.

$$\frac{\Theta \vdash op_a : \tau_1 \Rightarrow \tau_2 \quad \Gamma; \Theta \vdash_{\mathbf{Gen}} e : \tau_e/\varepsilon \rightsquigarrow C; S}{\Gamma; \Theta \vdash_{\mathbf{Gen}} op_a e : \tau_2/a * \varepsilon \rightsquigarrow C \cup \{\tau_e <: \tau_1\}; S}$$

Operators resemble functions, so we make sure that the type of  $e$  subtypes the type for arguments of  $op_a$ . Effect returned is the consequence of computing  $e$  combined with  $\{a\}$ .

$$\frac{\begin{array}{c} \Gamma; \Theta \vdash_{\mathbf{Gen}} h : \sigma \triangleright t/\epsilon \rightsquigarrow C_h; S_h \\ S_h \Gamma; S_h \Theta, (a : S_h \sigma) \vdash_{\mathbf{Gen}} e : \tau/\varepsilon \rightsquigarrow C_e; S_e \\ S_e S_h \Gamma, (x : \tau); S_e S_h \Theta \vdash_{\mathbf{Gen}} e_r : \tau_r/\varepsilon_r \rightsquigarrow C_r; S_r \end{array} \quad \text{fresh}(t) \quad \text{fresh}(\epsilon) \quad C = C_h \cup C_e \cup C_r \cup \{\varepsilon <: a * \epsilon, \tau_r <: t, \varepsilon_r <: \epsilon\}}{\Gamma; \Theta \vdash_{\mathbf{Gen}} \text{handle}_a e \{h; \text{return } x. e_r\} : t/\epsilon \rightsquigarrow C; S_r \circ S_e \circ S_h}$$

Unsurprisingly, rules for typing handlers looks rather complicated, but they are nothing more than the rules from  $\vdash$ , extended with subtyping and propagation of

constraints and substitution.

$$\begin{array}{c}
\frac{\Gamma, (x : \mathbf{Unit}), (k : t \rightarrow_{\varepsilon} \tau); \Theta \vdash_{\mathbf{Gen}} e : \tau' / \varepsilon' \rightsquigarrow C; S \quad \text{fresh}(t)}{\Gamma; \Theta \vdash_{\mathbf{Gen}} [(\mathbf{Raise}, x, k. e')] : \mathbf{Error} \triangleright \tau / \varepsilon \rightsquigarrow C \cup \{\tau' <: \tau, \varepsilon' <: \varepsilon\}; S} \\
\\
\frac{\begin{array}{c} \Gamma, (x : \mathbf{Unit}), (k : t \rightarrow_{\varepsilon} \tau); \Theta \vdash_{\mathbf{Gen}} e' : \tau' / \varepsilon' \rightsquigarrow C'; S' \\ S' \Gamma, (x : t), (k : \mathbf{Unit} \rightarrow_{\varepsilon} \tau); S' \Theta \vdash_{\mathbf{Gen}} e'' : \tau'' / \varepsilon'' \rightsquigarrow C''; S'' \\ C = C' \cup C'' \cup \{\tau' <: \tau, \varepsilon' <: \varepsilon, \tau'' <: \tau, \varepsilon'' <: \varepsilon\} \quad S = S'' \circ S' \quad \text{fresh}(t) \end{array}}{\Gamma; \Theta \vdash_{\mathbf{Gen}} [(\mathbf{Get}, x, k.e'); (\mathbf{Put}, x, k.e'')] : \mathbf{State} t \triangleright \tau / \varepsilon \rightsquigarrow C; S}
\end{array}$$

### 4.3 Solving constraints

The constraint solving algorithm we present is divided in two sub-procedures:

1. `solve_simple_constraints`  $C \ S$   
deals with both type- and effect-constraints behaving, much like the ordinary HM algorithm. As the rules for type subtyping are somewhat trivial, it can solve them efficiently and environment-agnostically. As we explain in the following subsection, there are some effect constraints that are *non-trivial* and cannot be resolved so easily, so they are dealt with by the second procedure:
2. `solve_constraints_within`  $\Gamma \ (\tau, \varepsilon) \ C \ S$   
deals with effect-constraints regarding effect unification variables ( $\epsilon$ ) occurring in  $\Gamma$ ,  $\tau$  and  $\varepsilon$ . The interesting constraints are of form  $\{\epsilon_1 <: I * \epsilon_2\}$ , as there are many substitutions that could satisfy one such constraint, but it is not obvious which one is *best* regarding bigger picture. In following subsections we argue how our approach approximates the most general type. For formal methods, see future work.

#### 4.3.1 Simple constraints

We call constraints solved by the first sub-procedure *simple* as the substitution they induce is *minimal*, meaning it is unambiguous that their premises *must* hold for the whole to be correct. Constraints that are *interesting* are *irreducible* by `solve_simple`. Other constraints are either solved or reduced to the *interesting* form. Notice that this procedure resolves all constraints that regard types (which are trivial) and we can use the substitution returned by it to find out *shapes* of types (and effects, to some extent) without *worrying* whether we would break some constraints by resolving them too quickly.

```

let expand  $\epsilon \ I \ S =$ 
  if  $I = \emptyset$  then  $S$ 
  else  $S[\epsilon \mapsto I * \epsilon']$  where  $\text{fresh}(\epsilon')$ 

```

```

let solve_simple C S =
  match C with
  |  $\emptyset$   $\rightarrow \emptyset; S$ 
  |  $\{\tau_1 <: \tau_2\} \cup C \rightarrow$ 
    match  $S[\tau_1], S[\tau_2]$  with
    |  $\tau'_1, \tau'_2$  when  $\tau'_1 = \tau'_2 \rightarrow$ 
      solve_simple C S
    |  $t, \tau$ 
    |  $\tau, t \rightarrow$ 
      solve_simple C  $S[t \mapsto \tau]$ 
    |  $\tau'_1 \rightarrow_{\varepsilon_1} \tau''_1, \tau'_2 \rightarrow_{\varepsilon_2} \tau''_2 \rightarrow$ 
      solve_simple  $\{\tau'_1 <: \tau'_2, \varepsilon_1 <: \varepsilon_2, \tau''_1 <: \tau''_2\} \cup C$  S
  |  $\{\varepsilon_1 <: \varepsilon_2\} \cup C \rightarrow$ 
    match  $S[\varepsilon_1], S[\varepsilon_2]$  with
    |  $I_1, I_2$  when  $I_1 \subseteq I_2 \rightarrow$  solve_simple C S
    |  $\iota, \_ \rightarrow$  solve_simple C S
    |  $I_1, I_2 * \epsilon \rightarrow$ 
      solve_simple C (expand  $(I_1 \setminus I_2) \epsilon$  S)
    |  $I_1 * \epsilon_1, I_2 * \epsilon_2 \rightarrow$ 
      if  $\epsilon_1 = \epsilon_2$  then solve_simple C (expand  $(I_1 \setminus I_2) \epsilon_2$  S)
      else let  $C'; S' =$  solve_simple C (expand  $(I_1 \setminus I_2) \epsilon_2$  S)
        in  $\{\epsilon_1 <: I_2 * \epsilon_2\} \cup C'; S'$ 
    |  $I_1 * \epsilon <: I_2$  when  $I_1 \subseteq I_2 \rightarrow$ 
      let  $C'; S' =$  solve_simple C S
      in  $\{\epsilon <: I_2\} \cup C'; S'$ 

```

We can now define function `solve'` from  $\vdash_{\text{Gen}}$  relation:

$$\text{solve}'(C, S) = \text{solve\_simple } C \ S$$

### 4.3.2 Interesting constraints

What makes constraints like  $(\epsilon_1 <: I_2 * \epsilon_2)$  interesting, is that there are many plausible substitutions that satisfy it. For every set of instances  $I_1 \subseteq I_2$ , substitutions  $[\epsilon_1 \mapsto I_1]$  and  $[\epsilon_1 \mapsto I_1 * \epsilon_2]$  obviously resolve the constraint, but clearly some substitutions are better than others.

As discussed in previous chapter,  $\rightarrow$  type constructor is *contravariant* to subtyping relation, which plays a great role in how we treat effect unification variables. During computation of `solve_constraints_within`  $\Gamma \ \tau/\varepsilon$  we keep information about *variance* of effect unification variables as a function  $V$ .

$$v ::= \oplus \mid \odot \mid \ominus \mid \times \quad (\text{variance})$$

If each occurrence of a variable  $\epsilon$  in  $\tau/\varepsilon$  is *covariant* (*positive*) and *contravariant* (*negative*) in  $\Gamma$ , then  $V(\epsilon) = \oplus$ . Similarly, if it appears only *contravariantly* in  $\tau/\varepsilon$  and only *covariantly* in  $\Gamma$ , then  $V(\epsilon) = \ominus$ . If it appears *invariantly* (*both* positively

and negatively), then  $V(\epsilon) = \odot$ . Finally, if  $\epsilon$  doesn't appear in type or environment at all, then  $V(\epsilon) = \times$ .

Because of the way we defined subtyping relation, we can *shrink* any covariant effect and *expand* any contravariant effect. Consider type  $\tau$  with some covariant effect  $\varepsilon_{\oplus}$ . Clearly, for any  $\varepsilon'_{\oplus} <: \varepsilon_{\oplus}$  we have  $\tau[\varepsilon'_{\oplus}/\varepsilon_{\oplus}] <: \tau$ . Analogously, for any contravariant effect  $\varepsilon_{\ominus}$  and any effect  $\varepsilon'_{\ominus}$  such that  $\varepsilon_{\ominus} <: \varepsilon'_{\ominus}$  we have  $\tau[\varepsilon'_{\ominus}/\varepsilon_{\ominus}] <: \tau$ .

With this intuition in mind, we can confidently resolve constraints  $\{\epsilon_1 <: I_2 * \epsilon_2\}$  if either  $V(\epsilon_1) = \oplus$  or  $V(\epsilon_2) = \ominus$  by assignment  $[\epsilon_1 \mapsto I_2 * \epsilon_2]$ , as it enables both covariant variables to be the smallest and contravariant to grow the largest possible.

It is important to include the non-appearing variables in our algorithm as well, as there are many unification variables that are generated along the way that do not appear in the examined type nor environment explicitly, but often do form *chains* of subtyping constraints like so (instances omitted for better visibility):

$$\epsilon_{\ominus} <: \epsilon_1 <: \dots <: \epsilon_n <: \epsilon_{\oplus} \text{ where } V(\epsilon_i) = \times$$

In such case, we want  $\epsilon_{\ominus}$  to be *the biggest* and  $\epsilon_{\oplus}$  to be *the smallest* possible, so we would like to deduce the substitution  $S$  such that  $S\epsilon_{\ominus} = S\epsilon_1 = \dots = S\epsilon_n = S\epsilon_{\oplus}$ , as it guarantees that it is in fact the case.

```

solve_within  $\Gamma \ \tau/\varepsilon \ C \ S =$ 
   $V := \text{gather\_free\_vars } S\Gamma \ S\tau/S\varepsilon$ 
  while  $\exists(I_1 * \epsilon_1 <: I_2 * \epsilon_2) \in S.C. \ \epsilon_1 \neq \epsilon_2 \wedge V(\epsilon_1), V(\epsilon_2) \text{ matches}$ 
    |  $\times, \oplus \mid \ominus, \times \mid \times, \odot \mid \odot, \times$ 
    |  $\oplus, \oplus \mid \ominus, \ominus \mid \odot, \odot \rightarrow$ 
       $C := C \setminus \{I_1 * \epsilon_1 <: I_2 * \epsilon_2\};$ 
       $S := (\text{expand } (I_1 \setminus I_2) \ \epsilon_2 \ S)[\epsilon_1 \mapsto I_2 * \epsilon_2];$ 
    |  $\ominus, \oplus \mid \ominus, \odot \mid \odot, \oplus \rightarrow$ 
       $C := C \setminus \{I_1 * \epsilon_1 <: I_2 * \epsilon_2\};$ 
       $S := (\text{expand } (I_1 \setminus I_2) \ \epsilon_2 \ S)[\epsilon_1 \mapsto I_2 * \epsilon_2];$ 
       $V := V[\epsilon_2 \mapsto \odot]$ 
  for  $(\epsilon, \oplus) \in V$ :
     $S := S[\epsilon \mapsto \iota]$ 
  return  $C, S$ 

```

We conclude the description of algorithm by defining function *solve* from  $\vdash_{\text{Gen}}$ :

$$\begin{aligned} \text{solve}(\Gamma, \tau/\varepsilon, C, S) &= \text{solve\_within } \Gamma \ \tau/\varepsilon \ C' \ S' \\ \text{where } S', C' &= \text{solve\_simple } C \ S \end{aligned}$$

## 4.4 Example

Now that we have a complete picture of how the inference algorithm works, let's take a look how constraints for *compose* would be generated and resolved.

$$\begin{array}{c}
\frac{(f : t_f) \in \Gamma_1 \quad \text{inst}(t_f) = t_f}{\Gamma_1 \vdash_{\mathbf{Gen}} f : t_f/\iota \rightsquigarrow \emptyset; id} \quad \frac{(f : t_g) \in \Gamma_1 \quad \text{inst}(t_g) = t_g}{\Gamma_1 \vdash_{\mathbf{Gen}} g : t_g/\iota \rightsquigarrow \emptyset; id} \quad \frac{(x : t_x) \in \Gamma_1 \quad \text{inst}(t_x) = t_x}{\Gamma_1 \vdash_{\mathbf{Gen}} x : t_x/\iota \rightsquigarrow \emptyset; id} \\
\hline
\Gamma_1 \vdash_{\mathbf{Gen}} (g \ x) : t_1/\epsilon_1 \rightsquigarrow C_1; id \\
\hline
\Gamma_1 \vdash_{\mathbf{Gen}} f (g \ x) : t_2/\epsilon_2 \rightsquigarrow C_1; id \\
\hline
(f : t_f), (g : t_g) \vdash_{\mathbf{Gen}} \lambda x. f (g \ x) : t_x \rightarrow_{\epsilon_2} t_2/\iota \rightsquigarrow C_1; id \\
\hline
(f : t_f) \vdash_{\mathbf{Gen}} \lambda g. \lambda x. f (g \ x) : t_g \rightarrow t_x \rightarrow_{\epsilon_2} t_2/\iota \rightsquigarrow C_1; id \\
\hline
\vdash_{\mathbf{Gen}} \lambda f. \lambda g. \lambda x. f (g \ x) : t_f \rightarrow t_g \rightarrow t_x \rightarrow_{\epsilon_2} t_2/\iota \rightsquigarrow C_1; id
\end{array}$$

$$C_1 = \{\iota <: \epsilon_1, \iota <: \epsilon_1, t_g <: t_x \rightarrow_{\epsilon_1} t_1\}$$

$$C_2 = \{\iota <: \epsilon_1, \iota <: \epsilon_1, t_g <: t_x \rightarrow_{\epsilon_1} t_1, \iota <: \epsilon_2, \epsilon_1 <: \epsilon_2, t_f <: t_1 \rightarrow_{\epsilon_2} t_2\}$$

$$\Gamma_1 = (f : t_f), (g : t_g), (x : t_x)$$

$$\text{solve\_simple } C_1 \text{ id} =$$

$$\text{solve\_simple } \{\epsilon_1 <: \epsilon_2, t_g <: t_x \rightarrow_{\epsilon_1} t_1, t_f <: t_1 \rightarrow_{\epsilon_2} t_2\} \text{ id} =$$

$$\{\epsilon_1 <: \epsilon_2\}; S_1$$

$$S_1 = [t_g \mapsto t_x \rightarrow_{\epsilon_1} t_1, t_f \mapsto t_1 \rightarrow_{\epsilon_2} t_2]$$

$$\text{solve\_within } \emptyset ((t_1 \rightarrow_{\epsilon_2} t_2) \rightarrow (t_x \rightarrow_{\epsilon_1} t_1) \rightarrow t_x \rightarrow_{\epsilon_2} t_2) \iota \{\epsilon_1 <: \epsilon_2\} S_1 =$$

$$V = [\epsilon_1 \mapsto \ominus, \epsilon_2 \mapsto \odot];$$

$$\text{return } \emptyset; S_2$$

$$S_2 = [\epsilon_1 \mapsto \epsilon_2, t_g \mapsto t_x \rightarrow_{\epsilon_1} t_1, t_f \mapsto t_1 \rightarrow_{\epsilon_2} t_2]$$

$$\text{solve}(\emptyset, t_f \rightarrow t_g \rightarrow t_x \rightarrow_{\epsilon_2} t_2/\iota, C_1, id) =$$

$$\emptyset; S_1[\epsilon_1 \mapsto \epsilon_2]$$

$$\tau_1 = (t_1 \rightarrow_{\epsilon_2} t_2) \rightarrow (t_x \rightarrow_{\epsilon_2} t_1) \rightarrow t_x \rightarrow_{\epsilon_2} t_2$$

$$\text{gen}(\emptyset, \tau_1) = \pi_1$$

$$\pi_1 = \forall \alpha, \beta, \gamma, \delta. (\alpha \rightarrow_{\delta} \beta) \rightarrow (\gamma \rightarrow_{\delta} \alpha) \rightarrow \gamma \rightarrow_{\delta} \beta$$

$$\Gamma_2 = (\text{compose} : \pi_1)$$

Notice that the inferred type is indeed the most general one, as discussed in Chapter 3.



## Chapter 5

# Implementation

The algorithm we described is accompanied by a software implementation, written in pure OCaml. Differences from  $\vdash_{\mathbf{Gen}}$  include that the built substitution is applied globally *on the fly*, rather than returned, as the unification variables are implemented using memory references that are updated throughout the process. By convention, all the occurrences *universal variables* (denoted  $'\tau\alpha$  or  $'\varepsilon\alpha$ ) are quantified *by default* (the  $\forall$ -quantifier of type schemes is implicit).

### 5.1 Project structure

Type system and expressions of the calculus are defined in `calculus.ml` file:

```
type identifier = string
type instance = identifier
type instance = identifier

type 'a univar = Free of identifier | Bound of 'a

type typ =
  | Unit
  | Int
  | Arrow    of typ * typ * effect                (*  $\tau_1 \rightarrow_\varepsilon \tau_2$  *)
  | TypVar   of typ univar ref                    (*  $t$  *)
  | GenTyp   of identifier                        (*  $'\tau\alpha$  *)
  | Forall   of instance * signature * typ        (*  $\forall(a : \sigma). \tau$  *)
  | IllTyped

and signature = Error | State of typ              (* Error | State  $\tau$  *)

and effect =
  | Fixed    of instance set                      (*  $I$  *)
  | Flexible of instance set * effect univar ref  (*  $I * \epsilon$  *)
  | GenEff   of instance set * identifier          (*  $I *' \varepsilon\alpha$  *)
```

```

type expr =
  | Nil                                     (* () *)
  | I           of int                     (* n *)
  | V           of var                     (* x *)
  | Lam         of var * expr              (*  $\lambda x. e$  *)
  | Fun         of var * var * expr        (* fun f, x. e *)
  | Let         of var * expr * expr       (* let x = e1 in e2 *)
  | App         of expr * expr             (* e1 e2 *)
  | Op         of instance * op * expr     (* opa e *)
  | Handle     of instance * expr * handler (* handlea e {h} *)
  | ILam       of instance * signature * expr (*  $\lambda(a:\sigma). e$  *)
  | IApp       of expr * instance          (* e a *)
  | UHandle    of expr * handler           (* handle e {h} *)
  | UOp        of op * expr               (* op e *)

and op = Raise | Get | Put

and handler = (op * var * var * expr) list * var * expr
              (* [(op, x, k, e)] ; return x. er *)

type op_type = typ * typ                  (*  $\tau_1 \Rightarrow \tau_2$  *)

type type_effect = typ * effect           (*  $\tau/\varepsilon$  *)

```

The UHandle and UOp stand for *unnamed* operators and handlers, used without declaring an instance. The way they're dealt with is by simple desugaring into named handlers with fresh instance.

Judgement  $\vdash_{\text{Gen}}$  is implemented in `inference.ml` by the `infer` sub-procedure of the `infer_type_with_env` function, which returns the type after resolving constraints, using procedures described in Chapter 4. Auxiliary procedures, such as the union-find operations and set operations are included in `utils.ml`.

```

type 'a constraints = ('a * 'a) list

type 'a environment = (identifier * 'a) list
type env = typ environment
type ienv = signature environment

val solve_simple : typ constraints -> effect constraints ->
  effect constraints

val solve_within : env -> type_effect -> effect constraints ->
  effect constraints

val infer_type_with_env : env -> ienv -> expr -> env *
  type_effect

```

The `env`, returned by `infer_type_with_env` contains types of variables bound by *let-constructs* in `expr`, included for meaningful examples printing.

## 5.2 Tutorial

File `examples.ml` contains a few meaningful examples of type inference. Simply running `ocamlbuild examples.byte` in project's root directory should build the executable and all its dependencies. By convention we print unresolved unification variables with `?` sign. Below, we have included output of selected examples.

```
$> ./examples.byte
⊢ λe:Error. λx. raise_e x : ∀e:Error. Unit -{e}-> ?τ1 / ι

⊢ handle
  put 21
  {put v k. k () | get () k. k 37 | return x. x} : Unit / ι

⊢ λx. λy. λz. (x z) (y z) : (?τ2 -{?ε2}-> ?τ4 -{?ε2}-> ?τ5) ->
  (?τ2 {?ε2}-> ?τ4) -> ?τ2 -{?ε2}-> ?τ5 / ι

⊢ let id = λx. x in id id : ?τ2 -> ?τ2 / ι
  (id : 'τa -> 'τa)

⊢ let apply =
  λf. λx. f x
  in apply (λx. x) : ?τ4 -> ?τ4 / ι
  (apply : ('τa -{'εb}-> 'τc) -> 'τa -{'εb}-> 'τc)

⊢ let fix =
  fun fix f. f (fix f)
  in fix (λx. λy. λz. 2) : ?τ6 -> ?τ7 -> Int / ι
  (fix : ('τa -{'εb}-> 'τa) -{'εb}-> 'τa)

⊢ let update =
  λs:State a. λf. put_s (f (get_s ()))
  in handle_b
    (λ(). get_b ()) ((update<b>) (λx. x))
  { get () k. λc. (k c) c
  | put v k. λc. (k ()) v
  | return x. λc. x} : ?τ15 -> ?τ15 / ι
  (update : ∀s:State a. (a -{s,'εa}-> a) -{s,'εa}-> Unit)

⊢ let move_map =
  λfrom:State a. λto:State b. λf. put_to (f (get_from ()))
  in 1 : Int / ι
  (move_map : ∀from:State a. ∀to:State b.
    (a -{from,'εa}-> b) -{from,to,'εa}-> Unit)
```



## Chapter 6

# Conclusions and future work

Based on the *original calculus*, we have proposed a subset calculus and type system that enabled us to implement ML-polymorphism for decidable type inference. The algorithm we have shown is an refinement of algorithm  $\mathbb{W}$  with simple subtyping relation, with extended calculus and type system for algebraic effects.

Let-polymorphism allows programmers to omit implicit *type-lambdas* (denoted by  $\Lambda$ ) and implicit type instantiating, making programmers' lives easier and the code more readable. Main focus of future work should be focused on how to achieve similar flexibility for instance lambdas. The problem is somewhat more difficult, as whether we use the *type-lambda* or not does not affect the semantics of expression, while the presence or absence of type lambdas certainly change it.

The way we resolve effect constraints is not a formally sound method. Presence of subtyping relation in type inference does not impose much added difficulty from the theoretical point of view, but it does make the practical constraint solving a much harder task. Because of this relation, constraints describing effects form a partially ordered set, and finding the most general substitution that satisfies it, is beyond this work. Proving this property is essential to arguing about the principal type, which is highly desirable for an inference algorithm, as it enables programmers to be sure that even if they skip all the type annotations, the inferred type will not be less general. This property is essential for full type inference.

A few works had been written the topic of type inference with subtyping, proving its difficulty, including the Ph.D. thesis of François Pottier[7]. However, our subtyping relation is much simpler than the one in his work, then maybe, with just the right amount of depth and a bit of ingenuity, studying the topological properties of subtyping relation, could lead to some simpler, yet *sound* solution.

Nonetheless, the *heuristic* approach we described produced desired results for all the examples we have tried. However, it may be possible that there is type and constraint set that our method fails to generate the *most general unifier*, but given short time window for this work we were not able to find such example.



# Bibliography

- [1] J. Roger Hindley. 1969. The Principal Type-Scheme of an Object in Combinatory Logic. *Transactions of the American Mathematical Society*, 146, 29-60. <https://doi.org/10.2307/1995158>
- [2] Robin Milner. 1978. A theory of type polymorphism in programming, *Journal of Computer and System Sciences*, Volume 17, Issue 3, 1978, Pages 348-375, ISSN 0022-0000, [https://doi.org/10.1016/0022-0000\(78\)90014-4](https://doi.org/10.1016/0022-0000(78)90014-4).
- [3] Luis Damas and Robin Milner. 1982. Principal type-schemes for functional programs. In *Proceedings of the 9th ACM SIGPLAN-SIGACT symposium on Principles of programming languages (POPL '82)*. Association for Computing Machinery, New York, NY, USA, 207–212. <https://doi.org/10.1145/582153.582176>
- [4] Jean Pierre Talpin, Pierre Jouvelot. 1992. The type and effect discipline, *Proceedings of the Seventh Annual IEEE Symposium on Logic in Computer Science*, Santa Cruz, CA, USA, 1992, pp. 162-173, <https://doi.org/10.1109/LICS.1992.185530>.
- [5] 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 (January 2020), 29 pages. <https://doi.org/10.1145/3371116>
- [6] Pottier François and Didier Rémy. 2005. The Essence of ML Type Inference. MIT press, Chapter 10 of *Advanced topics in types and programming languages*.
- [7] François Pottier. 1998. Type Inference in the Presence of Subtyping: from Theory to Practice. [Research Report] RR-3483, INRIA. [ffinria-00073205f](https://inria.fr/rr/rr-3483)