

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

7 września 2020

Abstract

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Contents

| | | |
|----------|---|-----------|
| 1 | Introduction | 7 |
| 2 | Calculus | 9 |
| 3 | Type system | 11 |
| 3.1 | Subtyping | 13 |
| 3.2 | Parametricity | 13 |
| 3.3 | Principal type | 13 |
| 4 | Inference algorithm | 15 |
| 4.1 | Remarks on effect unification | 15 |
| 4.2 | Generating constraints | 16 |
| 4.3 | Solving constraints | 18 |
| 4.3.1 | Simple constraints | 18 |
| 4.3.2 | Interesting constraints | 19 |
| 4.4 | Illustrative example | 20 |
| 5 | Implementation | 21 |
| 5.1 | Representation | 21 |
| 5.2 | Project structure | 21 |
| 5.3 | Tutorial | 21 |
| 6 | Future work | 23 |
| | Bibliography | 25 |

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[5], Robin Milner, who presented the type inference algorithm W[4], 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.

The particular system presented by Biernacki et al[1] 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 enclosing expression by 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.

We transcribe the calculus given by Biernacki et al 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 is a subset of work by Biernacki et al[1]. It is an extension of standard *call-by-value lambda calculus with let* by effect handlers and operators. We adjusted it to match the style of *ML-the-calculus* and *ML-the-type-system*[2].

| | |
|--|------------------------|
| var $\ni x, \dots$ | (term variables) |
| qvar $\ni \alpha, \dots$ | (quantified variables) |
| tvar $\ni t, \dots$ | (type variables) |
| evar $\ni \epsilon, \dots$ | (effect variables) |
| ivar $\ni a, \dots$ | (instance variables) |
| type $\ni \tau ::= \alpha \mid t \mid \mathbf{Unit} \mid \mathbf{Int} \mid \tau \rightarrow_{\epsilon} \tau \mid \forall(a : \sigma). \tau$ | (types) |
| scheme $\ni \pi ::= \forall \alpha. \pi \mid \tau$ | (type schemes) |
| effect $\ni \epsilon ::= \alpha \mid \epsilon \mid a \mid \epsilon \circ \epsilon$ | (effects) |
| signature $\ni \sigma ::= \mathbf{Error} \mid \mathbf{State} \tau$ | (signatures) |
| expr $\ni e ::= x \mid () \mid n \mid \lambda x. e \mid \mathbf{fun} f x. e \mid e e \mid \mathbf{let} x = e \mathbf{in} e$ $\mid \lambda(a : \sigma). e \mid e a \mid op_a e \mid \mathbf{handle}_a e \{h; \mathbf{return} x.e\}$ | (terms) |
| operator $\ni op ::= \mathbf{Raise} \mid \mathbf{Get} \mid \mathbf{Put}$ | (operators) |
| handler $\ni h ::= [(\mathbf{Raise}, x, k.e)] \mid [(\mathbf{Get}, x, k.e); (\mathbf{Put}, x, k.e)]$ | (handlers) |

Terms (or expressions) are given by:

- variables, bound by abstractions, let-expressions, handlers, or environment Γ ,
- 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** _{a} $e \{h; \text{return } x.e_r\}$ 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 h , in which x gets bounds to value of e_x . Supplying the *continuation* k with some value v continues evaluation of e , with v substituted in place of operation call. 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. After e is evaluated, its value is bound to x in e_r expression, which is the final value returned by handler.

Instead of allowing effects of arbitrary signatures, we limited them to instances of either **Error** or **State** τ . Arbitrary signatures could be dealt with in similar fashion as *ADT's* (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 α variables). Accordingly we have dropped the type-lambdas (Λ) from calculus in favor of the polymorphic let-construct.

The semantics of calculus strictly follows the rules defined in Biernacki et al. For formal reduction rules and additional insight, see “Binders by Day, Labels by Night” [1].

Chapter 3

Type system

We have tweaked the type-and-effect system constructed by Biernacki et al[1] to match the restrictions required by classical methods of type inference. Thus, we introduce *unification variables* for both types and effects (written t and ϵ , respectively). We use these variables to denote “yet to be determined” types. 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 Γ (assigning types to variables), and Θ (assigning signatures to instances), term e *inhabits* syntactical type and effect τ/ϵ (which means that computing e would yield a value of type τ and possibly cause effect ϵ).

Typing terms:

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

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 τ_1 , computing the function’s body e produces some output of type τ_2 , possible causing effect ϵ

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

Accordingly, 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 $\vec{\alpha} = \text{freevars}(\tau) \setminus \text{freevars}(\Gamma)$. Here 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 different instances of same effect, i.e two cells of memory of **State** τ , there are lambda terms, which can be applied with instances bound by handlers or other instance lambdas. .

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

Then the operators of instance a and type $\tau_1 \rightarrow \tau_2$ if applied with some expression e of type τ_1 and effect ε are typed with $\tau_2 / a \circ \varepsilon$

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

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

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

Typing handlers:

$$\frac{\Gamma, (x : \mathbf{Unit}), (k : \tau' \rightarrow_\varepsilon \tau); \Theta \vdash e : \tau / \varepsilon}{\Gamma; \Theta \vdash [(\mathbf{Raise}, x, k.e)] : \mathbf{Error} \triangleright \tau / \varepsilon}$$

$$\frac{\Gamma, (x : \mathbf{Unit}), (k : \tau' \rightarrow_\varepsilon \tau); \Theta \vdash e_{\mathbf{Get}} : \tau / \varepsilon \quad \Gamma, (x : \tau'), (k : \mathbf{Unit} \rightarrow_\varepsilon \tau'); \Theta \vdash e_{\mathbf{Put}} : \tau / \varepsilon}{\Gamma; \Theta \vdash [(\mathbf{Get}, x, k.e_{\mathbf{Get}}); (\mathbf{Put}, x, k.e_{\mathbf{Put}})] : \mathbf{State} \tau' \triangleright \tau / \varepsilon}$$

There are two clauses for typing handlers, as we only have kinds of signatures. It could easily be extended for other signatures but that's not important to this work.

Typing operators:

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

3.1 Subtyping

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

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 relation on τ_1 and τ'_1 is reversed).

The subtyping plays a vital role in usability of the calculus. Consider a term f , a function that does some calculation, but allowed the function to fail, i.e.

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

but nothing stops us from applying some pure function in this place. On the other hand, it would be undesirable if we could supply a term expecting a pure function and with an effectful one. Clearly this property gives us flexibility, while keeping effects under control.

3.2 Parametricity

Our type system maintains predicative prenex polymorphism of ML, extended with universal quantification over effects because original paper maintains it.

3.3 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 π is called *principal* if any other type that could be given to e is an *instantiation* of it. For example, consider $e = \lambda x. x$. There's a few types that could be given to it: $\mathbf{Int} \rightarrow \mathbf{Int}$, but clearly any correct type we would think of would not be more general than $\forall \alpha. \alpha \rightarrow \alpha$.

Definition 3.3.1. Principal Type

In our type system, π is a principal type of e in environments Γ, Θ iff

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

For details about subtyping and principal type, see future work.

For now we will give an example of how to approximate principal type. Consider two type schemes that could be assigned to a term $compose = \lambda f. \lambda g. \lambda x. f(g\ x)$, representing function composition:

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

At first glance, it may look like π_1 is the “correct” type for $compose$, as it seems more natural, i.e. given functions f causing effect ε_f , and g causing ε_g , $compose\ f\ g$ is a function that would apply them both, so it clearly must be causing effect $\varepsilon_f \circ \varepsilon_g$.

While this reasoning is sound, we are interested in deriving the most concise type possible. Let’s see if there exists ε_h such that π_1 instantiated with arbitrary effects $\varepsilon_f, \varepsilon_g$ subtypes π_2 instantiated with ε_h :

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

Clearly, $\pi_1[\varepsilon_f, \varepsilon_g / \alpha, \beta] <: \pi_2[\varepsilon_h / \gamma]$ does not hold for ε_f and ε_g other than ι , thus π_1 cannot be a principal type of e . On the other hand, if we were to check if for arbitrary ε_h there exist ε_f and ε_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:

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

Clearly if we choose $\varepsilon_f = \varepsilon_g = \varepsilon_h$, it is satisfied, which means that π_2 is indeed more general π_1 . We designed our inference algorithm with this intuition in mind and π_2 is the desired result that our implementation actually infers.

Chapter 4

Inference algorithm

The algorithm we present loosely follows original algorithm \mathbb{W} , executing 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 .

$$\begin{aligned} \textbf{constraints} \ni C &::= \emptyset \mid \{\tau <: \tau\} \mid \{\varepsilon <: \varepsilon\} \mid C \cup C \\ \textbf{substitution} \ni S &::= id \mid [t \mapsto \tau] \mid [\varepsilon \mapsto \varepsilon] \mid S \circ S \end{aligned}$$

Constraints $\{x <: y\}$ 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 $a <: \epsilon_1 \circ \epsilon_2$, for some instance variable a and some effect unification variables ϵ_1, ϵ_2 . To resolve such constraint we have a few viable options:

1. Expand ϵ_1 with a .
2. Expand ϵ_2 with a .
3. Expand both ϵ_1 and ϵ_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:

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

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.

4.2 Generating constraints

As in 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 we shall interpret it.

We chose to present \vdash_{Gen} typing rules as close as possible to the practical inference algorithm. Hence, we heavily use *unification variables* (denoted t for types and ϵ for effects) and the rules are somewhat *algorithmic*, meaning the focus shifts from *checking* to *obtaining* a type.

Judgement \vdash_{Gen} is constructed in a very *algorithmic* way, meaning the focus shifts from *checking* typing derivation to *gathering* constraints and type.

Judgement \vdash_{Gen} should be treated as a set of rules for *generating* types and constraints (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 Γ (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 type generated by it (after we solve all the constraints and build the substitution) to the given term checks by the syntactic rules we defined in chapter 3. Proof of our algorithm is left for future work.

Inferring expressions:

$$\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} \\[10pt] \frac{(x : \pi) \in \Gamma \quad \text{inst}(\pi) = \tau}{\Gamma; \Theta \vdash_{\mathbf{Gen}} x : \tau/\iota \rightsquigarrow \emptyset; id} \\[10pt] \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} \\[10pt] \frac{\text{fresh}(t_1) \quad \text{fresh}(\epsilon) \quad \text{fresh}(t_2) \quad \Gamma, (f : t_1 \rightarrow_{\epsilon} t_2), (x : t_1); \Theta \vdash_{\mathbf{Gen}} e : \tau/\varepsilon \rightsquigarrow C; S}{\Gamma; \Theta \vdash_{\mathbf{Gen}} \mathbf{fun} \ f x. e : t_1 \rightarrow_{\epsilon} t_2/\iota \rightsquigarrow C \cup \{t_1 \rightarrow_{\varepsilon} \tau <: t_1 \rightarrow_{\epsilon} t_2\}; S} \\[10pt] \frac{\Gamma; \Theta \vdash_{\mathbf{Gen}} e_1 : \tau_1/\varepsilon_1 \rightsquigarrow C_1; S_1 \quad S_1\Gamma; S_1\Theta \vdash_{\mathbf{Gen}} e_2 : \tau_2/\varepsilon_2 \rightsquigarrow C_2; S_2 \quad \text{fresh}(t) \quad \text{fresh}(\epsilon) \quad C = C_1 \cup C_2 \cup \{\varepsilon_1 <: \epsilon, \varepsilon_2 <: \epsilon, \tau_1 <: \tau_2 \rightarrow_{\epsilon} t\}}{\Gamma; \Theta \vdash_{\mathbf{Gen}} e_1 e_2 : t/\epsilon \rightsquigarrow C; S_2 S_1} \\[10pt] \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}} \mathbf{let} \ x = e_1 \ \mathbf{in} \ e_2 : \tau_2/\varepsilon_2 \rightsquigarrow C \cup C_2; S_2 \circ S} \\[10pt] \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}} \mathbf{let} \ x = e_1 \ \mathbf{in} \ e_2 : \tau_2/\epsilon \rightsquigarrow C \cup C_2 \cup \{\varepsilon_1 <: \epsilon, \varepsilon_2 <: \epsilon\}; S_2 \circ S} \\[10pt] \frac{\Gamma; \Theta, (a : \sigma) \vdash_{\mathbf{Gen}} e : \tau/\varepsilon \rightsquigarrow C'; S' \quad \text{solve}'(\Gamma, C', S') = C; S \quad S\varepsilon = \iota \quad S\tau' = \tau}{\Gamma; \Theta \vdash_{\mathbf{Gen}} \lambda(a : \sigma). e : \forall(a : \sigma). \tau/\iota \rightsquigarrow C; S} \\[10pt] \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} \\[10pt] \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} \end{array}$$

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

And now typing handlers:

$$\frac{\Gamma, (x : \text{Unit}), (k : t \rightarrow_{\epsilon} \tau); \Theta \vdash_{\text{Gen}} e : \tau_{\text{Raise}}/\epsilon_{\text{Raise}} \rightsquigarrow C; S \quad \text{fresh}(t)}{\Gamma; \Theta \vdash_{\text{Gen}} [(\text{Raise}, x, k.e)] : \text{Error} \triangleright \tau/\epsilon \rightsquigarrow C \cup \{\tau_{\text{Raise}} <: \tau, \epsilon_{\text{Raise}} <: \epsilon\}; S}$$

$$\frac{\begin{array}{l} \Gamma, (x : \text{Unit}), (k : \tau' \rightarrow_{\epsilon} \tau); \Theta \vdash_{\text{Gen}} e_{\text{Get}} : \tau_{\text{Get}}/\epsilon_{\text{Get}} \rightsquigarrow C_{\text{Get}}; S_{\text{Get}} \\ S_{\text{Get}} \Gamma, (x : \tau'), (k : \text{Unit} \rightarrow_{\epsilon} \tau); S_{\text{Get}} \Theta \vdash_{\text{Gen}} e_{\text{Put}} : \tau_{\text{Put}}/\epsilon_{\text{Put}} \rightsquigarrow C_{\text{Put}}; S_{\text{Put}} \\ C = C_{\text{Get}} \cup \{\tau_{\text{Get}} <: \tau, \epsilon_{\text{Get}} <: \epsilon\} \cup C_{\text{Put}} \cup \{\tau_{\text{Put}} <: \tau, \epsilon_{\text{Put}} <: \epsilon\} \end{array}}{\Gamma; \Theta \vdash_{\text{Gen}} [(\text{Get}, x, k.e_{\text{Get}}); (\text{Put}, x, k.e_{\text{Put}})] : \text{State } \tau' \triangleright \tau/\epsilon \rightsquigarrow C; S_{\text{Put}} \circ S_{\text{Get}}}$$

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 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** $\Gamma \ (\tau, \epsilon) \ C \ S$
 deals with effect-constraints regarding effect unification variables (ϵ) occurring in Γ , τ and ϵ . 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 *irreducible* by **solve_simple** are of form $\epsilon_1 <: I * \epsilon_2$; other constraints are either solved or reduced to the *interesting* form.

```

let expand  $\epsilon$   $I$   $S$  =
  if  $I = \emptyset$  then  $S$ 
  else  $S[\epsilon \mapsto I * \epsilon']$  where 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  $C$   $S[t \mapsto \tau]$ 
    |  $\tau'_1 \rightarrow_{\epsilon_1} \tau''_1, \tau'_2 \rightarrow_{\epsilon_2} \tau''_2 \rightarrow$ 
      solve_simple  $\{\tau'_1 <: \tau'_2, \epsilon_1 <: \epsilon_2, \tau''_1 <: \tau''_2\} \cup C$   $S$ 
  |  $\{\epsilon_1 <: \epsilon_2\} \cup C \rightarrow$ 
    match  $S[\epsilon_1], S[\epsilon_2]$  with
    |  $\iota, \_ \rightarrow$  solve_simple  $C$   $S$ 
    |  $I_1, I_2$ 
    |  $I_1, I_2 * \epsilon \rightarrow$ 
      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_{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 $I_1 \subseteq I_2$, substitution $[\epsilon_1 \mapsto I_1]$ or $[\epsilon_1 \mapsto I_1 * \epsilon_2]$ obviously resolves 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` Γ τ/ε we keep information about *variance* of effect unification variables as a function V .

$$\text{variance} \ni v ::= \oplus \mid \odot \mid \ominus \mid \times$$

If variable ϵ appears in Γ , τ , and ε only in *covariant* (*positive*) positions, then $V(\epsilon) = \oplus$. If it appears only *contravariantly* (*negatively*) then $V(\epsilon) = \ominus$. If it appears *invariantly* (*both* positively and negatively), then $V(\epsilon) = \odot$. Finally, if ε 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 τ with some covariant effect ε_{\oplus} . Clearly, for any $\varepsilon'_{\oplus} <: \varepsilon_{\oplus}$ we have $\tau[\varepsilon'_{\oplus}/\varepsilon_{\oplus}] <: \tau$. Analogously, for any contravariant effect ε_{\ominus} and any effect ε'_{\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 $[\epsilon_1 \mapsto I_2 * \epsilon_2]$.

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

```

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

$$\text{solve}(\Gamma, \tau/\varepsilon, C, S) = \text{solve_constraints_within } \Gamma \tau/\varepsilon C' S'$$

$$\text{where } S', C' = \text{solve_simple } C S$$

4.4 Illustrative example

No what we have a complete picture of how the inference algorithm works, let's take a look how constraints for term would be generated and resolved.

Chapter 5

Implementation

Pure OCaml.

5.1 Representation

How calculus, type system, constraints and substitution are implemented.

5.2 Project structure

Which code does what

5.3 Tutorial

Some examples and how to run it.

Chapter 6

Future work

Let-polymorphism allows us to omit implicit type-lambdas (usually denoted by Λ) and type instantiation, making programmers' lives easier. Way of doing so for instance lambdas is yet to be found.

Ensuring that the inferred type is principal would most probably require a long and difficult proof. The way we resolve effect constraints is also a heuristic and not a formal method. Such constraints form a graph whose topology should be studied with according depth. There are works by Francois that do this in different subtyping relation.

Finally, resolving constraints *formally* would require us to study the topology of graphs constructed. Such constraints form a partially ordered set and finding *formally sound* substitution that satisfies it is beyond this work, so we present a *heuristic* approach, which produces desired results for all the examples we tried. However, it may be possible that there is type and constraint set that our method fails to generate the *most general* type, but given short time window for this work we were not able to find such example.

Bibliography

- [1] 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>
- [2] Pottier François and Didier Rémy. 2005. The Essence of ML Type Inference. MIT press, Chapter 10 of *Advanced topics in types and programming languages*.
- [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] Robin Milner, 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).
- [5] 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>