# **Capability-Flavoured Effects**

by

### Aaron Craig

A thesis
submitted to the Victoria University of Wellington
in fulfilment of the
requirements for the degree of
Bachelor of Science with Honours
in Computer Science.

Victoria University of Wellington 2017

#### **Abstract**

Privilege separation and least authority are principles for developing safe software, but existing languages offer insufficient techniques for allowing developers and architects to make informed design choices enforcing them. Languages adhering to the object-capability model impose constraints on the ways in which privileges are used and exchanged, giving rise to a form of lightweight effect-system. This effect-system allows architects and developers to make more informed choices about whether code from untrusted sources should be used. This paper develops an extension of the simply-typed lambda calculus to illustrate the ideas and proves it sound.

# Acknowledgments

# **Contents**

1	Intr	oductio	n		1				
2	Bacl	Background							
	2.1	Defini	ng a Formal Language		5				
	2.2	$\lambda^{\rightarrow}$ : S	imply-Typed $\lambda$ -Calculus		7				
	2.3								
	2.4	ETL: E	Effect-Typed Language		10				
	2.5	The C	apability Model		10				
	2.6	First-C	Class Modules		11				
3	Effe	Effect Calculi 13							
	3.1	$\lambda_{\pi}^{\rightarrow}$ : O	peration Calculus		13				
		3.1.1	Definition of $\lambda_{\pi}^{\rightarrow}$		13				
		3.1.2	Soundness of $\lambda_{\pi}^{\rightarrow}$		17				
	3.2	$\lambda_{\pi,\varepsilon}^{\rightarrow}$ : E	psilon Calculus		17				
		3.2.1	Definition		18				
		3.2.2	Soundness of $\lambda_{\pi,\varepsilon}^{\rightarrow}$		25				
4	App	Applications 29							
	4.1	Encod	ings		29				
		4.1.1	Unit		29				
		4.1.2	Let		29				
		4.1.3	Conditionals		30				
		4.1.4	Tuples		30				
	4.2	Examp	oles		30				
5	Eval	luation			33				
	5.1	Relate	d Work		33				
	5.2	Future	e Work		33				
	5.3	Concl	asion		34				
A	Proc	ofs			35				

vi *CONTENTS* 

# Chapter 1

### Introduction

Good software is distinguished from bad software by design qualities such as security, maintainability, and performance. One of these is the *principle of least authority*: that software components should only have access to the information and resources necessary for their purpose [14]. For example, a logger module, which need only append to a file, should not have arbitrary read-write access. Another is *privilege separation*, where the division of a program into components is informed by what resources are needed and how they are to be propagated [?].

Matters get complicated when a program is contains untrustworthy components. Such cases may arise in a development environment which adheres to *code ownership*, whereby groups of developers may function as particular experts over certain components [?]. When they interact with code sourced from outside their domain of expertise, they may make false assumptions or violate the internal constraints of other components. Applications may allow third-party plug-ins, in which case third-party code is sourced from an untrustworthy source. A web mash-up is a particular kind of software that brings together disparate applications into a central service, in which case the disparate applications may be untrustworthy.

When a codebase has untrustworthy code it may be impossible or infeasible for developers to verify that it adequately enforces separation of privileges and POLA. Often they may not have access to the original source code. This leaves developers to make a judgement call on whether this untrusted code should be used or executed based on the type interface and accompanying verbal contracts.

One approach to privilege separation is the *capability* model. A capability is an unforgeable token granting its bearer permission to perform some operation [4]. Resources in a program are only exercised though the capabilities granting them. Although the notion of a capability is an old one, there has been recent interest in the application of the idea to programming language. Miller has identified the ways in which capabilities should proliferate to encourage *robust composition* — a set of ideas summarised as "only connectivity begets connectivity" [10]. In his paradigm, the reference graph of a program

is the same as the access graph. This eliminates *ambient authority* — a pervasive enemy in determining by interface what privileges a component might exercise. Building on these ideas, Maffeis et. al. formalised *capability-safety* of a language, showing a subset of Caja (a JavaScript implementation) meets this notion [8].

While capabilities adequately separate privileges, understanding the way in which those privileges are exercised has received less attention. This area falls under the domain of effect systems, which extends type systems to include intensional information about the way in which a program executes [11]. For example, a logger's log method may have append as its effect, but a sloppy or malicious implementation may incur extra effects, such as write or close. This suggests the logger may be doing more than just logging, and knowing this guides the developer to a more informed decision about whether to use this particular implementation.

One obstacle to the adoption of effect systems is their verbosity. An effect system such as the Talpin-Jouvelot system requires the annotation of all values in a program [cite]. This requires the developer to be aware, at all points, of what effects are in scope. Minor alterations to the signatures and effects of a component require the labels on all interacting components to change in accordance. This overhead is something the developer must carry with them at all stages of programming, affecting the usability of effect systems. Successive works have focussed on reducing these issues through techniques such as effect-inference, but the benefit of capabilities for effect-based reasoning has received less attention.

This paper suggests that capability-safe languages might get a low-cost effect system with minimal overhead. To incur an effect requires one to possess a capability for the appropriate resource. By tracking the propagation of these capabilities, with help from the proliferation constraints imposed by capability-safety, we might better understand what are the possible effects of a piece of code, in a manner which facilitates modular development.

This paper's contribution is to develop an extension to the simply-typed lambda calculus  $\lambda^{\rightarrow}$  called the epsilon calculus  $\lambda^{\rightarrow}_{\pi,\varepsilon}$ .  $\lambda^{\rightarrow}_{\pi,\varepsilon}$  introduces a new import construct which selects those capabilities used in a piece of unlabelled code. A sound inference can be made about the unlabelled code by inspecting the type signatures of those functions in scope at the point of introduction. This is made possible by the restrictions imposed on the use and exchange of capabilities.

Chapter 2 covers some background information on capabilities and programming language semantics. It also establishes the various conventions used throughout. We identify some of the benefits obtained by capabilities and effects, and some of the drawbacks we set out to solve.

Chapter 3 introduces static and dynamic rules for  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ , developing and proving a formulation of soundness appropriate for the type-and-effect discipline.

Chapter 4 shows how  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  might solve these drawbacks, and try to convince the reader

that  $\lambda_{\pi,\varepsilon}^{\to}$  can be implemented in existing capability-safe languages in a routine manner.

# Chapter 2

# **Background**

In this section we cover some of the necessary concepts and existing work informing this paper. No prior knowledge is assumed.

#### 2.1 Defining a Formal Language

We will consider a programming language as three sets of rules: the grammar, the static rules, and the dynamic rules. To illustrate each we give rules for a toy language that evaluates basic arithmetic operations on  $\mathbb{N}$ .

The grammar specifies what strings are syntactically legal in the language. A grammar is specified by giving the different categories of terms, and specifying all the possible forms which instantiate that category. Metavariables range over the terms of the category for which they are named. The conventions for specifying a grammar are based on standard Backur-Naur form [1]. Figure 2.1. shows a simple grammar describing integer literals and arithmetic expressions on them.

Figure 2.1: Grammar for arithmetic expressions.

The static rules specify the type system and other constraints on terms with certain well-behavedness properties. In our case, we're interested in what makes a program well-typed, which is to say that execution of the program never gets stuck due to type-errors. For example, a well-typed program will never try to add two booleans, because addi-

tion only makes sense on numbers. A well-typed program will never try to evaluate an undefined variable, because variables must be defined.

Static rules are specified by *inference rules*. An inference rule is given as a set of premises above a dividing line which, if they hold, imply the result below the line. An application of an inference rule is called a *judgement*. Judgements take place in typing contexts, denoted by  $\Gamma$ , which map variables to types. A basic judgement like  $\Gamma \vdash e : \tau$  means that executing e will result in a term of type  $\tau$  (if it terminates). The contents of  $\Gamma$  are specified as a comma-separated sequence of variable-type pairs. The order is irrelevant. A consequence of this convention is that if  $\Gamma \vdash e : \tau$ , then  $\Gamma' \vdash e : \tau$ , where  $\Gamma$  is contained in  $\Gamma'$ . When a judgement can be proven from the empty context we leave the left of the turnstile blank, as in  $\vdash e : \tau$ .

Though our arithmetic language has no subtyping, most interesting languages do. This judgement is written  $\tau_1 <: \tau_2$  and it means that values of  $\tau_1$  may be provided anywhere instances of  $\tau_2$  are expected. Effect systems have another judgement, which ascribes a type and set of effects to an expression. We shall cover this in greater detail later.

$$\frac{\Gamma \vdash e : \tau}{\Gamma, x : \mathtt{Int} \vdash x : \mathtt{Int}} \ (\mathtt{T-VAR}) \ \frac{\Gamma \vdash e_1 : \mathtt{Int} \quad \Gamma \vdash e_2 : \mathtt{Int}}{\Gamma \vdash e_1 + e_2 : \mathtt{Int}} \ (\mathtt{T-ADD})$$

Figure 2.2: Inference rules for typing arithmetic expressions.

The dynamic semantics specifies what is the meaning of a legal term. There are different flavorus of dynamic semantics, but the one we use is called *small-step semantics*. This is a set of inference rules specifying how a program is executed. A single application of one of these rules is called a *reduction*.

$$e \longrightarrow e$$

$$\frac{e_1 \longrightarrow e_1'}{e_1 + e_2 \longrightarrow e_1' + e_2} \text{ (E-ADD1) } \frac{e_2 \longrightarrow e_2'}{l_1 + e_2 \longrightarrow l_1 + e_2'} \text{ (E-ADD2) } \frac{l_1 + l_2 = l_3}{l_1 + l_2 \longrightarrow l_3} \text{ (E-ADD3)}$$

Figure 2.3: Inference rules for reducing arithmetic expressions.

Almost all type systems in which we are interested are sound. Soundness is a property that holds between the static and dynamic rules of a language, which says that if a program e is considered well-typed by the static rules, then its reduction under the dynamic rules will never produce a runtime type-error. The exact definition of soundness depends on the semantics under consideration, but it is often split into two parts: progress and preservation. The progress theorem states every term, except those of a particular category called values, can always be reduced by applying some dynamic rule.

The preservation theorem is that programs remain well-typed under reduction. Adequate formulations of these two theorems for the language under consideration give us soundness.

In the language of arithmetic expressions, progress and preservation would be the following:

**Theorem 1** (Progress). *If*  $\Gamma \vdash e$ : Int and e is not an integer constant, then  $e \rightarrow e'$ .

**Theorem 2** (Preservation). *If*  $\Gamma \vdash e$ : Int and  $e \rightarrow e'$  then  $\Gamma \vdash e'$ : Int.

These theorems can be proven by structural induction on the typing judgement  $\Gamma \vdash e$ : Int. This is a common proof technique in formal semantics. A more thorough explanation of this sort of induction can be found in TAPL [13, p. 31].

Some languages extend the notion of a type system to a *type-and-effect system*. Effects describe intensional information about the way in which a program executes [11]. A judgement like  $\Gamma \vdash e : \tau$  with {File.write} can be interpreted as meaning that execution of e will result in a value of type  $\tau$  (if it halts), and during execution might perform the write effect on a File. This judgement is an upper-bound: it says that the only effects e might have when executed are in the set {File.write}. The upper-bound is not tight: if e executes, it may not incur every effect ascribed to it.

### 2.2 $\lambda^{\rightarrow}$ : Simply-Typed $\lambda$ -Calculus

The simply-typed  $\lambda$ -calculus  $\lambda^{\rightarrow}$  is a model of computation which incorporates a basic theory of types. It was first described by Alonzo Church [3]. In this section we will cover the rules of a version of  $\lambda^{\rightarrow}$  with subtyping, and summarise its basic properties.

e	::=		exprs:	au	-	::=		types:
		x	variable				B	$base\ type$
		e e	application				$\tau \to \tau$	$arrow\ type$
		v	value					
				I	7	::=		contexts:
v	::=		values:				Ø	$empty\ ctx.$
		$\lambda x : \tau.e$	abstraction				$\Gamma, x : \tau$	$var.\ binding$

Figure 2.4: Grammar for  $\lambda^{\rightarrow}$ .

Terms in  $\lambda^{\rightarrow}$  is based on function abstraction and function application. Types are either drawn from a set of base types B, or constructed using  $\rightarrow$  ("arrow"). Given types  $\tau_1$  and  $\tau_2$ ,  $\rightarrow$  can be used to compose a new type,  $\tau_1 \rightarrow \tau_2$ , which is the type of function taking  $\tau_1$ -typed terms as input to produce  $\tau_2$ -typed terms as output. For example, given

 $B = \{ \text{Bool}, \text{Int} \}$ , the following are examples of valid types: Bool, Int, Bool  $\rightarrow$  Bool, Bool  $\rightarrow$  Int, Bool  $\rightarrow$  (Bool  $\rightarrow$  Int).

Arrow is right-associative, so Bool  $\rightarrow$  Bool  $\rightarrow$  Int = Bool  $\rightarrow$  (Bool  $\rightarrow$  Int). The terms "arrow-type" and "function-type" are interchangeable.

$$\frac{\Gamma \vdash e : \tau}{\Gamma, x : \tau \vdash x : \tau} \text{ (T-VAR)} \quad \frac{\Gamma, x : \tau_1 \vdash e : \tau_2}{\Gamma \vdash \lambda x : \tau_1 . e : \tau_1 \to \tau_2} \text{ (T-Abs)}$$
 
$$\frac{\Gamma \vdash e_1 : \tau_2 \to \tau_3 \quad \Gamma \vdash e_2 : \tau_2}{\Gamma \vdash e_1 e_2 : \tau_3} \text{ (T-APP)} \quad \frac{\Gamma \vdash e : \tau_1 \quad \tau_1 <: \tau_2}{\Gamma \vdash e : \tau_2} \text{ T-Subsume}$$
 
$$\frac{\tau_1' <: \tau_1 \quad \tau_2 <: \tau_2'}{\tau_1 \to \tau_2 <: \tau_1' \to \tau_2'} \text{ (S-Arrow)}$$

Figure 2.5: Static rules for  $\lambda^{\rightarrow}$ .

Static rules for  $\lambda^{\rightarrow}$  are summarised in Figure 2.5. There are two T-VAR states that a variable bound in some context can be typed as its binding. T-ABS states that a function can be typed in  $\Gamma$  if  $\Gamma$  can type the body of the function when the function's argument has been bound. T-APP states that an application is well-typed if the left-hand expression is a function (has an arrow-type  $\tau_2 \rightarrow \tau_3$ ) and the right-hand expression has the same type as the function's input ( $\tau_2$ ).

T-SUBSUME is the rule which says you may a type a term more generally as any of its supertypes. For example, if we had base types Int and Real, and a rule specifying Int  $\langle : \text{Real}, \text{ a term of type Int can also be typed as Real}.$  This allows programs such as  $(\lambda x : \text{Real}.x)$  3 to type, because 3 can be widened to a Real by T-SUBSUMPTION.

The only subtyping rule we provide is S-ARROW, which describes when one function is a subtype of another. Note how the subtyping relation on the input types is reversed from the subtyping relation on the functions. This is called *contravariance*. Contrast this with the relation on the output type, which preserves the order. That is called *covariance*. Arrow-types are contravariant in their input and covariant in their output.

Because there are no axiomatic subtyping rules, there are no provable subtyping judgements in this  $\lambda$ -calculus. In practice we add extra subtyping judgements for the base-types we have chosen as primitive in our calculus. For example, if Int and Real were base-types, we might add Real <: Int as a rule. This is largely an implementation detail and particular to your chosen set of base-types, so we give no such rules here (but will when describing  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ ).

A useful principle in the design and understanding of subtyping rules is Liskov's

substitution principle, which states that if  $\tau_1 <: \tau_2$ , then instances of  $\tau_2$  can be replaced with instances of  $\tau_1$  without changing the semantics or correctness of the program [6]. Subtyping rules are not usually semantic-preserving, but we'll occasionally use this ideal to motivate certain rules.

substitution ::  $\hat{\mathbf{e}} \times \hat{\mathbf{e}} \times \hat{\mathbf{v}} \rightarrow \hat{\mathbf{e}}$ 

$$[\hat{v}/y]x = \hat{v}$$
, if  $x = y$   
 $[\hat{v}/y]x = x$ , if  $x \neq y$   
 $[\hat{v}/y](\lambda x : \hat{\tau}.\hat{e}) = \lambda x : \hat{\tau}.[\hat{v}/y]\hat{e}$ , if  $y \neq x$  and  $y$  does not occur free in  $\hat{e}$   
 $[\hat{v}/y](\hat{e}_1 \ \hat{e}_2) = ([\hat{v}/y]\hat{e}_1)([\hat{v}/y]\hat{e}_2)$ 

Figure 2.6: Substitution for  $\lambda^{\rightarrow}$ .

Before giving the dynamic rules we must define the substitution function.  $\operatorname{substitution}(e,e_1,x)$  replaces all free occurrences of x with  $e_1$  in e. The short-hand is  $[e_1/x]e1$ . When performing multiple substitutions we use the notation  $[e_1/x_1,e_2/x_2]e$  as shorthand for  $[e_2/x_2]([e_1/x_1]\hat{e})$ . Note how the order of the variables has been flipped; the substitutions occur as they are written, left-to-right.

To avoid accidental avriable capture we adopt the convention of  $\alpha$ -conversion, whereby we freely and implicitly interchange expressions which are equivalent up to the naming of bound variables [13, p. 71]. This elides some tedious bookkeeping in the definition of substitution. Consequently, we shall assume variables are (re-)named in this way to avoid accidental capture.

$$\frac{\hat{e}_{1} \longrightarrow \hat{e}'_{1} \mid \varepsilon}{\hat{e}_{1}\hat{e}_{2} \longrightarrow \hat{e}'_{1}\hat{e}_{2} \mid \varepsilon} \text{ (E-APP1)} \qquad \frac{\hat{e}_{2} \longrightarrow \hat{e}'_{2} \mid \varepsilon}{\hat{v}_{1}\hat{e}_{2} \longrightarrow \hat{v}_{1}\hat{e}'_{2} \mid \varepsilon} \text{ (E-APP2)}$$

$$\frac{(\lambda x : \hat{\tau}.\hat{e})\hat{v}_{2} \longrightarrow [\hat{v}_{2}/x]\hat{e} \mid \varnothing}{(\lambda x : \hat{\tau}.\hat{e})\hat{v}_{2} \longrightarrow [\hat{v}_{2}/x]\hat{e} \mid \varnothing} \text{ (E-APP3)}$$

$$\frac{\hat{e} \longrightarrow^* \hat{e} \mid \varnothing}{\hat{e} \longrightarrow^* \hat{e}' \mid \varepsilon_1} \text{ (E-MULTISTEP1)} \quad \frac{\hat{e} \longrightarrow \hat{e}' \mid \varepsilon}{\hat{e} \longrightarrow^* \hat{e}' \mid \varepsilon_1} \text{ (E-MULTISTEP2)}$$

$$\frac{\hat{e} \longrightarrow^* \hat{e}' \mid \varepsilon_1 \quad \hat{e}' \longrightarrow^* \hat{e}'' \mid \varepsilon_2}{\hat{e} \longrightarrow^* \hat{e}'' \mid \varepsilon_1 \cup \varepsilon_2} \text{ (E-MULTISTEP3)}$$

Figure 2.7: Dynamic rules for  $\lambda^{\rightarrow}$ .

There are two sorts of reduction rules: single-steps  $e \longrightarrow e$ , and multi-steps  $e \longrightarrow^*$ 

which consists of zero<sup>1</sup> or more single-steps. These semantics are call-by-value, because execution of a function (E-APP3) can only happen when the subexpressions have been reduced (by E-APP1 until the left-hand side is a value, and then E-APP2).

The soundness property for  $\lambda^{\rightarrow}$  is as follows.

**Theorem 3** ( $\lambda^{\rightarrow}$  Soundness). *If*  $\Gamma \vdash e_A : \tau_A$  *and*  $e_A \longrightarrow^* e_B$ , *then*  $\Gamma \vdash e_B : \tau_B$ , *where*  $\tau_B <: \tau_A$ .

 $\lambda^{\to}$  is also strongly-normalizing, meaning that well-typed terms always halt. A consequence is that  $\lambda^{\to}$  cannot express certain programs and is therefore Turing-incomplete. One routine way of making the language Turing-complete is to add a fix operator. Turing completenses is not overly relevant to the typing rules of  $\lambda^{\to}_{\pi,\varepsilon'}$ , so to simplify the presentation we shall leave our  $\lambda^{\to}$  a Turing-incomplete language.

Revisit this depending on how you encode types and stuff in  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ 

#### 2.3 Effect Systems

#### 2.4 ETL: Effect-Typed Language

#### 2.5 The Capability Model

A capability is a unique, unforgeable reference, giving its bearer permission to perform some operation [4]. A piece of code S has authority over a capability C if it can directly invoke the operations endowed by C; it has transitive authority if it can indirectly invoke the operations endowed by a capability C (for example, by deferring to another piece of code with authority over C).

In a capability model, authority can only proliferate in the following ways [10]:

- 1. By the initial set of capabilities passed into the program (initial conditions).
- 2. If a function or object is instantiated by its parent, the parent gains a capability for its child (parenthood).
- 3. If a function or object is instantiated by a parent, the parent may endow its child with any capabilities it possesses (endowment).
- 4. A capability may be transferred via method-calls or function applications (introduction).

The rules of authority proliferation are summarised as: "only connectivity begets connectivity".

<sup>&</sup>lt;sup>1</sup>We permit multi-step reductions of length zero to be consistent with Pierce, who defines multi-step reduction as a reflexive relation[13, p. 39].

Primitive capabilities are called *resources*. Resources model those initial capabilities passed into the runtime from the system environment. A capability is either a resource, or a function or object with (potentially transitive) authority over a capability. An example of a resource might be a particular file. A function which manipulates that file (for example, a logger) would also be a capability, but not a resource. Any piece of code which uses a capability, directly or indirectly, is called *impure*. For example,  $\lambda x : Int. x$  is pure, while  $\lambda f : File. f.log("error message")$  is impure.

A relevant concept in the design of capability-based programming languages is *ambient authority*. This is a kind of exercise of authority over a capability C which has not been explicitly [9]. Figure 2.4. gives an example in Java, where a malicious implementation of List.add attempts to overwrite the user's .bashrc file. MyList gains this capability by importing the java.io.File class, but its use of files is not immediate from the signature of its functions.

Ambient authority is a challenge to POLA because it makes it impossible to determine from a module's signature what authority is being exercised. From the perspective of Main, knowing that MyList.add has a capability for the user's .bashrc file requires one to inspect the source code of .bashrc; a necessity at odds with the circumstances which often surround untrusted code and code ownership.

A language is *capability-safe* if it satisfies this capability model and disallows ambient authority. Some examples include E, Js, and Wyvern. **Get citations.** 

#### 2.6 First-Class Modules

The exact way in which modules work is language-dependent, but we are particularly interested in languages with a first-class module systems. First-class modules are important in capability-safe languages because they mean capability-safe reasoning operates across module boundaries. Because modules are first-class, they must be instantiated like regular objects. They must therefore select their capabilities, and be supplied those capabilities by the proliferation rules of the capability model. In practice, first-class modules can be achieved by having module declarations desugar into an underlying lambda or object representation. This generally requires an "intermediate representation" of the language, which is simpler than the one in which programmers write.

Java is an example of a mainstream language whose modules are not first-class. Scala has first-class modules [12], but is not capability-safe. Smalltalk is a dynamically-typed capability-safe language with first-class modules [2]. Wyvern is a statically-typed capability-safe language with first-class modules [5].

```
import java.io.File;
2 (import java.io.IOException;
3 import java.util.ArrayList;
 class MyList<T> extends ArrayList<T> {
  O. @Override
 O. public boolean add(T elem) {
      try {
        File file = new File("$HOME/.bashrc");
        file.createNewFile();
  00.99 } catch (IOException e) {}
       return super.add(elem);
  0..9
 0}..99
14
import java.util.List;
3 class Main {
  D. public static void main(String[] args) {
 D.9 List<String> list = new MyList<String>();
 0.9 list.add(''doIt'');
7 (0...9)
8 (}...9)
```

Figure 2.8: Main exercises ambient authority over a File.

### Chapter 3

### **Effect Calculi**

### 3.1 $\lambda_{\pi}^{\rightarrow}$ : Operation Calculus

The operation calculus  $\lambda_{\pi}^{\rightarrow}$  is an extension of  $\lambda^{\rightarrow}$  with primitive capabilities (resources) and their operations. Effects are identified with operation-calls. A program's runtime effects are those operations which it calls during execution. The static rules approximate the runtime effects of an expression, forming a simple effect system. These rules are very simple, but formalising and studying them will introduce new notations and a new concept of effect-soundness on which we shall build the  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ .

#### **3.1.1 Definition of** $\lambda_{\pi}^{\rightarrow}$

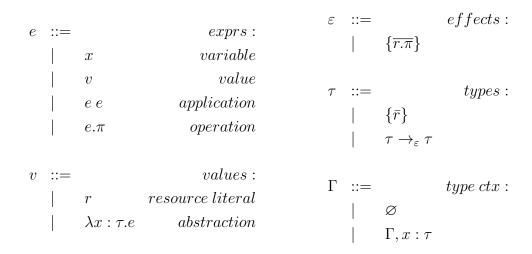


Figure 3.1: Grammar for  $\lambda_{\pi}^{\rightarrow}$ .

The base types of  $\lambda_{\pi}^{\rightarrow}$  are sets of resources, denoted by  $\{\bar{r}\}$ . Resources are drawn from a fixed set R of variables, and model those initial capabilities passed in from the system environment. They cannot be created at runtime. When a resource type is ascribed to a

program, as in the judgement  $\Gamma \vdash e : \{\bar{r}\}$ , it means that if e terminates it will result in a resource literal  $r \in \bar{r}$ . A resource literal r is a variable ranging over the elements of R.

An operation is a special action that can be invoked on a resource-typed expression. For example, we might invoke the open operation on a File resource. Operations are drawn from a fixed set  $\Pi$  of variables and cannot be created at runtime.

An effect is an operation performed on a resource. Formally, they are members of  $R \times \Pi$ , but for readability we write File.write over (File,write). A set of effects is denoted by  $\varepsilon$ . Effects and operations notationally look the same, but should be distinguished: an effect is some action upon a resource which may happen during runtime; an operation-call is an expression representing the actual invocation of an effect at runtime.

In practical applications, operations take arguments. For example, when writing to a file, we want to specify *what* is being written to the file, e.g. File.write("mymsg"). Because we are only concerned with the use and propagation of effects, and not the semantics of particular effects, we make the simplifying assumption that all operations are null-ary.

The only type constructor is  $\rightarrow_{\varepsilon}$ , where  $\varepsilon$  is a concrete set of effects.  $\tau_1 \rightarrow_{\varepsilon} \tau_2$  is the type of a function which takes inputs of type  $\tau_1$ , produces outputs of type  $\tau_2$ , and incurs no more effects than those contained in  $\varepsilon$ . For example, the type of a function which sends a message over a socket and returns a success flag could be  $\text{Str} \rightarrow_{\text{Socket.write}} \text{Bool.}$  From this signature we can tell this function will not open or close the socket, because the annotation on the arrow does not have those effects. A valid implementation of this function might not write to the Socket, because {Socket.write} is an upper-bound on the effects which can happen. Every function type in  $\lambda_{\pi}^{\rightarrow}$  must be annotated with an upper-bound in this way.

The static rules for  $\lambda_{\pi}^{\rightarrow}$  are summarised in Figure 3.2. The typing judgement  $\Gamma \vdash e$ :  $\tau$  with  $\varepsilon$  means that successive reductions on e will yield terms of type  $\tau$ , and collectively incur no more than those effects in  $\varepsilon$ . This  $\varepsilon$  is a conservative approximation to the runtime effects of executing e, so it may contain effects which don't happen at runtime.

The rules for variables and values are:  $\varepsilon$ -VAR,  $\varepsilon$ -RESOURCE, and  $\varepsilon$ -ABS. These are identical to the rules in  $\lambda^{\rightarrow}$ , except they approximate the set of effects as  $\varnothing$ . Although a fucntion and a resource literal both encapsulate capabilities, something must be done to them (apply the function, operate on the resource) to incur a runtime effect.

The effects of a lambda application are: the effects of evaluating its subexpressions, and the effects incurred by executing the body of the lambda to which the left-hand side evaluates. Those last effects are obtained from the label on the lambda's arrow-type in the first premise.

The effects of an operation call are: the effects of evaluating the subexpression, and the single effect incurred when the subexpression is reduced to a resource literal r, and operation  $\pi$  is invoked on it. It is not always possible to know statically which exact resource literal the subexpression reduces to (if it halts at all). For example, the program (if System.randomBool then File else Socket).close may either reduce to File.close

 $\Gamma \vdash e : \tau \; \mathtt{with} \; \varepsilon$ 

$$\overline{\Gamma, x : \tau \vdash x : \tau \text{ with } \varnothing} \ (\varepsilon\text{-VAR}) \ \overline{\Gamma, r : \{r\} \vdash r : \{r\} \text{ with } \varnothing} \ (\varepsilon\text{-Resource})$$

$$\frac{\Gamma, x : \tau_2 \vdash e : \tau_3 \text{ with } \varepsilon_3}{\Gamma \vdash \lambda x : \tau_2 . e : \tau_2 \to_{\varepsilon_3} \tau_3 \text{ with } \varnothing} \ (\varepsilon\text{-ABS}) \ \overline{\Gamma \vdash e_1 : \tau_2 \to_{\varepsilon} \tau_3 \text{ with } \varepsilon_1} \ \overline{\Gamma \vdash e_2 : \tau_2 \text{ with } \varepsilon_2} \ (\varepsilon\text{-APP})$$

$$\frac{\Gamma \vdash e : \{\bar{r}\} \quad \forall r \in \bar{r} \mid r : \{r\} \in \Gamma \quad \pi \in \Pi}{\Gamma \vdash e . \pi : \text{Unit with } \{\bar{r}.\pi\}} \ (\varepsilon\text{-OPERCALL})$$

$$\frac{\Gamma \vdash e : \tau \text{ with } \varepsilon \quad \tau <: \tau' \quad \varepsilon \subseteq \varepsilon'}{\Gamma \vdash e : \tau' \text{ with } \varepsilon'} \ (\varepsilon\text{-SUBSUME})$$

 $\Gamma \vdash e : au$  with arepsilon

$$\frac{\tau_1' <: \tau_1 \quad \tau_2 <: \tau_2' \quad \varepsilon \subseteq \varepsilon'}{\tau_1 \rightarrow_{\varepsilon} \tau_2 <: \tau_1' \rightarrow_{\varepsilon'} \tau_2'} \text{ (S-ARROW)} \qquad \frac{r \in r_1 \implies r \in r_2}{\{\bar{r}_1\} <: \{\bar{r}_2\}} \text{ (S-RESOURCE)}$$

Figure 3.2: Type-with-effect judgements in  $\lambda_{\pi}^{\rightarrow}$ .

or Socket.close. In such cases, the safe approximation is to type the conditional as {File, Socket}.  $\varepsilon$ -OPERCALL would then approximate the runtime effects of the operation call as {File.close, Socket.close}. Being able to type an expression as a (nonsingleton) set of resources requires an extra rule: S-RESOURCE. This says that a subset of resources is also a subtype.

Because we're not really interested in what exactly an operation call does, every operation is valid for every resource. This can give bizarre programs — Sensor.readTemp seems like a sensible operation call, but what about File.readTemp? Because we are not interested in the specific behaviours of operations, we permit any operation on any resource to simplify the static rules.

The other subtyping rule is S-ARROW, a modification of the rule from  $\lambda^{\rightarrow}$ . In addition to this rule being contravariant in the input and covariant in the output, it is also covariant in the effects. Justified in terms of the Liskov substitution principle, any possible effect which might be incurred by the subtype should be expected by the supertype, otherwise substitution of a supertype for a subtype would allow for the possibility of new effects not possible under the original.

Figure 3.2. shows the updated definition of substitution. In  $\lambda_{\pi}^{\rightarrow}$ , a variable may only be substituted for a value, making the function a partial one. This restriction is imposed because, if a variable can be replaced with an arbitrary expression, then we might also be introducing arbitrary effects — a situation which violates the preservation of effects under reduction.

 $\texttt{substitution} :: \texttt{e} \times \texttt{v} \times \texttt{v} \to \texttt{e}$ 

```
 [v/y]x = v, \text{ if } x = y \\ [v/y]x = x, \text{ if } x \neq y \\ [v/y](\lambda x : \tau.e) = \lambda x : \tau.[v/y]e, \text{ if } y \neq x \text{ and } y \text{ does not occur free in } e \\ [v/y](e_1 \ e_2) = ([v/y]e_1)([v/y]e_2) \\ [v/y](e_1.\pi) = ([v/y]e_1).\pi
```

Figure 3.3: Substitution function.

$$\frac{e_{1} \longrightarrow e'_{1} \mid \varepsilon}{e_{1}e_{2} \longrightarrow e'_{1} \mid e_{2} \mid \varepsilon} \text{ (E-APP1)} \quad \frac{e_{2} \longrightarrow e'_{2} \mid \varepsilon}{v_{1} \mid e_{2} \longrightarrow v_{1} \mid e'_{2} \mid \varepsilon} \text{ (E-APP2)} \quad \frac{(\lambda x : \hat{\tau}.\hat{e})\hat{v}_{2} \longrightarrow [\hat{v}_{2}/x]\hat{e} \mid \varnothing}{(\lambda x : \hat{\tau}.\hat{e})\hat{v}_{2} \longrightarrow [\hat{v}_{2}/x]\hat{e} \mid \varnothing} \text{ (E-APP3)}$$

$$\frac{\hat{e} \to \hat{e}' \mid \varepsilon}{\hat{e}.\pi \longrightarrow \hat{e}'.\pi \mid \varepsilon} \text{ (E-OPERCALL1)} \quad \frac{r \in R \quad \pi \in \Pi}{r.\pi \longrightarrow \text{unit} \mid \{r.\pi\}} \text{ (E-OPERCALL2)}$$

Figure 3.4: Single-step reductions.

Single-step reduction is now a relation on an expression e and a pair  $e \times \varepsilon$ , representing the reduced expression and all effects incurred during the single-step of computation. E-APP1 and E-APP2 incur whatever is the effect of reducing their subexpressions. E-APP3 incurs no effects.

The new single-step rules are E-OPERCALL1 and E-OPERCALL2. The former reduces the receiver of an operation-call, and the latter performs an operation on a resource literal. E-OPERCALL1 incurs whatever is the effect of reducing the subexpression. E-OPERCALL2, which reduces the operation-call  $r.\pi$ , incurs the effect  $r.\pi$ .

Operation calls reduce to unit (which is a derived form; see Encodings). unit is a value representing the absence of information (because it is the only value of its type). Because we are not interested in the semantics of effects, we choose unit as the result of reducing an operation-call.

A multi-step reduction consists of zero<sup>1</sup> or more single-step reductions. The resulting effect-set is the union of all the single-steps taken.

<sup>&</sup>lt;sup>1</sup>We permit multi-step reductions of length zero to be consistent with Pierce, who defines multi-step reduction as a reflexive relation[13, p. 39].

$$\begin{array}{c} \left[\hat{e}\longrightarrow^{*}\hat{e}\mid\varepsilon\right] \\ \\ \overline{\hat{e}\rightarrow^{*}\hat{e}\mid\varnothing} \text{ (E-MULTISTEP1)} & \frac{\hat{e}\rightarrow\hat{e}'\mid\varepsilon}{\hat{e}\rightarrow^{*}\hat{e}'\mid\varepsilon} \text{ (E-MULTISTEP2)} \\ \\ \frac{\hat{e}\rightarrow^{*}\hat{e}'\mid\varepsilon_{1}\quad\hat{e}'\rightarrow^{*}\hat{e}''\mid\varepsilon_{2}}{\hat{e}\rightarrow^{*}\hat{e}''\mid\varepsilon_{1}\cup\varepsilon_{2}} \text{ (E-MULTISTEP3)} \end{array}$$

Figure 3.5: Multi-step reductions in  $\lambda_{\pi}^{\rightarrow}$ .

#### **3.1.2** Soundness of $\lambda_{\pi}^{\rightarrow}$

Our goal is to show  $\lambda_{\pi}^{\rightarrow}$  is sound. This requires an appropriate notion of *effect-soudnness*, as we need both the type and approximated effects of an expression to be preserved under reduction. Intuitively, a reduction  $e_A \longrightarrow^* e_B \mid \varepsilon$  is sound if the type system's approximation of the effects of e contains the actual runtime effects  $\varepsilon$ , because then it has statically accounted for every possible runtime effect. Below is a definition, phrased in terms of single-step reduction.

**Theorem 4** (Soundness). *If*  $\hat{\Gamma} \vdash \hat{e}_A : \hat{\tau}_A$  with  $\varepsilon_A$  and  $\hat{e}_A$  is not a value, then  $e_A \longrightarrow e_B \mid \varepsilon$ , where  $\hat{\Gamma} \vdash e_B : \hat{\tau}_B$  with  $\varepsilon_B$  and  $\hat{\tau}_B <: \hat{\tau}_A$  and  $\varepsilon_B \cup \varepsilon \subseteq \varepsilon_A$ .

This definition is slightly stronger, because it relates the approximation after reduction to the approximation before reduction, by saying it can only get more precise. This is analogous to the type-safe component of the definition, which states that the types of successive terms under reduction can only get more precise.

Our road to proving soundness takes the standard approach of showing that progress and preservation hold of  $\lambda_{\pi}^{\rightarrow}$ . Progress follows immediately by observing some properties of the typing rules.

**Lemma 1** (Canonical Forms). *The following are true:* 

• If  $\hat{\Gamma} \vdash \hat{v} : \hat{\tau}$  with  $\varepsilon$  then  $\varepsilon = \emptyset$ . • If  $\hat{\Gamma} \vdash \hat{v} : \{\bar{r}\}$  then  $\hat{v} = r$  for some  $r \in R$  and  $\{\bar{r}\} = \{r\}$ .

**Theorem 5** (Progress). *If*  $\hat{\Gamma} \vdash \hat{e} : \hat{\tau}$  with  $\varepsilon$  and  $\hat{e}$  is not a value, then  $\hat{e} \longrightarrow \hat{e}' \mid \varepsilon$ .

*Proof.* By induction on  $\hat{\Gamma} \vdash \hat{e} : \hat{\tau}$  with  $\varepsilon$ , for  $\hat{e}$  not a value. If the rule is  $\varepsilon$ -SUBSUMPTION it follows by inductive hypothesis. If  $\hat{e}$  has a reducible subexpression then reduce it. Otherwise use one of  $\varepsilon$ -APP3 or  $\varepsilon$ -OPERCALL2.

### 3.2 $\lambda_{\pi,\varepsilon}^{\rightarrow}$ : Epsilon Calculus

The effect calculus is based on the simply-typed lambda calculus  $\lambda^{\rightarrow}$ . There is one type constructor,  $\rightarrow$ . The base types are sets of resources, denoted by  $\{\bar{r}\}$ . Although the calculus has no primitive notions of integers or booleans, we shall assume these may be

encoded as they are in the usual way (e.g. as Church numerals) and make free use of them in examples as though they were standard, for the sake of readability.

#### 3.2.1 Definition

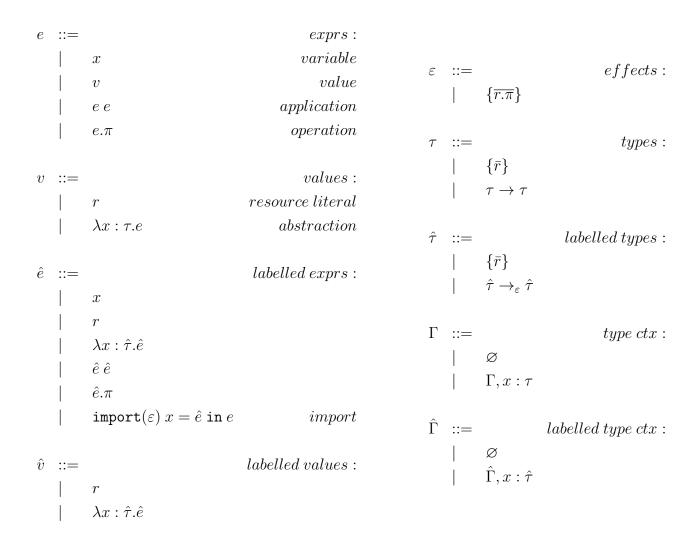


Figure 3.6: Effect calculus.

Resources are drawn from a fixed set R of variables, and model those initial capabilities passed in from the system environment. Resources cannot be created at runtime. When a resource type is ascribed to a program, as in the judgement  $\Gamma \vdash e : \{\bar{r}\}$ , it means that if e terminates it will result in a resource literal  $r \in \bar{r}$ .

A value v is either a resource literal r or a lambda abstraction  $\lambda x:\tau.e$ . The other forms of an expression are lambda application e e, variable x, and operation  $e.\pi$ . An operation is an action invoked on a resource. For example, we might invoke the open operation on a File resource. Operations are drawn from a fixed-set  $\Pi$  of variables. They cannot be created at runtime.

An effect is an operation performed on a resource. Formally, they are members of  $R \times \Pi$ , but for readability we write File.write over (File, write). A set of effects is denoted by  $\varepsilon$ . Effects and operations notationally look the same, but should be distinguished: an effect is some action upon a resource which may happen during runtime; an operation is the actual invocation of an effect at runtime.

In a practical language, operations should take arguments. For example, when writing to a file, we want to specify *what* is being written to the file, ala File.write("mymsg"). Because  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  is only concerned with the use and propagation of effects, and not the semantics of particular effects, we make the simplifying assumption that all operations are null-ary.

Expressions may be labelled with the set of effects they might incur during execution. This is achieved by annotating all arrow types inside the expression. If a metavariable represents a labelled expression, it will be written with a hat; if it represents an unlabelled expression, it will have no hat. Compare e and  $\hat{e}$ .

Labelling of an expression is *deep*. That is, every subterm of a labelled term is also labelled. Unlabelled terms are also deeply unlabelled. The only exception is the import expression, which is the only way to compose labelled and unlabelled code. import nests unlabelled code inside labelled code, and selects those capabilities  $\varepsilon$  over which the unlabelled code has authority. It is not possible to nest labelled code inside unlabelled code.

The distinction between labelled and unlabelled types and expressions requires us to have the notion of labelled and unlabelled contexts. Labelled contexts only bind variables to labelled types, whereas unlabelled contexts only bind variables to unlabelled types. There is no valid context which mixes labelled and unlabelled types.

A construct's labelled version is always denoted with a hat.

Given a piece of unlabelled code e and static effects  $\varepsilon$  we can produce a labelled piece of code  $\operatorname{annot}(e,\varepsilon)=\hat{e}$  by annotating every function with  $\varepsilon$ . In the reverse direction, given some labelled code  $\hat{e}$  we can produce an unlabelled piece of code  $\operatorname{erase}(\hat{e})=e$  by removing the labels on functions. Full definitions for these functions on expressions, types, and contexts are given in Figure 3.2. Note that  $\operatorname{erase}$  is undefined on  $\operatorname{import}$  expressions. We won't ever need to  $\operatorname{erase}$  import expressions, but it means the function is partial, so we need to be careful when we use it.

Annotation is not always safe. For instance,  $annot(\lambda l : Int \rightarrow_{File.read} Int. l 1, \varnothing)$  would overwrite the File.read effect permitted by l. annot is used in one place, in the dynamic rules, and for that limited use we will have to prove its safety.

We may wish to know what effects are encapsulated by a piece of labelled code. This is achieved by two functions,  $effects(\hat{e})$  and  $ho-effects(\hat{e})$ , which collectively compute the set of effects captured by  $\hat{e}$ . These are effects which may, directly or indirectly, be

```
annot :: e \times \varepsilon \rightarrow \hat{e}
           annot(r, \_) = r
           annot(\lambda x : \tau_1.e, \varepsilon) = \lambda x : annot(\tau_1, \varepsilon).annot(e, \varepsilon)
           annot(e_1 e_2, \varepsilon) = annot(e_1, \varepsilon) annot(e_2, \varepsilon)
           annot(e_1.\pi,\varepsilon) = annot(e_1,\varepsilon).\pi
annot :: 	au 	imes arepsilon 	o \hat{	au}
           \mathtt{annot}(\{\bar{r}\}, \_) = \{\bar{r}\}
           \operatorname{annot}(\tau \to \tau, \varepsilon) = \tau \to_{\varepsilon} \tau.
annot :: \Gamma \times \varepsilon \to \hat{\Gamma}
           annot(\emptyset, \_) = \emptyset
           annot(\Gamma, x : \tau, \varepsilon) = annot(\Gamma, \varepsilon), x : annot(\tau, \varepsilon)
erase :: \hat{	au} 
ightarrow 	au
           erase(\{\bar{r}\})
           erase(\hat{\tau}_1 \rightarrow_{\varepsilon} \hat{\tau}_2) = erase(\hat{\tau}_1) \rightarrow erase(\hat{\tau}_2)
erase :: \hat{e} \rightarrow e
           erase(r) = r
           erase(\lambda x : \hat{\tau}_1.\hat{e}) = \lambda x : erase(\hat{\tau}_1).erase(\hat{e})
           erase(e_1 e_2) = erase(e_1) erase(e_2)
           erase(e_1.\pi) = erase(e_1).\pi
```

Figure 3.7: Annotation functions.

invoked by  $\hat{e}$ . The difference between the two functions is in who supplies the effect.  $\mathsf{effect}(\hat{e})$  is the set of effects for which  $\hat{e}$  has direct authority, while ho-effects is the set of effects for which  $\hat{e}$  has (strictly) transitive authority. These higher-order effects are always supplied by some external environment.

For example, take the function which, given a file, reads and returns its contents (which are perhaps encoded as an integer). Its signature would be  $f: \{File\} \rightarrow_{File.read} Int$ . The effects(f) =  $\{File.read\} \cup effects(Int)$ , because any client using f will directly invoke the File.read operation and may use any resource encapsulated by the Int type. The ho-effects(f) =  $\{File.\pi \mid \pi \in \Pi\}$ , because to use f it must be supplied with a File literal from some outside source. Therefore, every possible effect on File is a higher-order effect.

The substitution function substitution( $\hat{e}, \hat{v}, x$ ) replaces all free occurrences of x with  $\hat{v}$  in  $\hat{e}$ . The short-hand is  $[\hat{v}/x]\hat{e}$ . When performing multiple substitutions we use the notation

21

$$\begin{split} \text{effects} &:: \hat{\pmb{\tau}} \to \pmb{\varepsilon} \\ &\quad \text{effects}(\{\bar{r}\}) = \{r.\pi \mid r \in \bar{r}, \pi \in \Pi\} \\ &\quad \text{effects}(\hat{\tau}_1 \to_{\varepsilon} \hat{\tau}_2) = \text{ho-effects}(\hat{\tau}_1) \cup \varepsilon \cup \text{effects}(\hat{\tau}_2) \end{split}$$
  $\text{ho-effects} &:: \hat{\pmb{\tau}} \to \pmb{\varepsilon} \\ &\quad \text{ho-effects}(\{\bar{r}\}) = \varnothing \\ &\quad \text{ho-effects}(\hat{\tau}_1 \to_{\varepsilon} \hat{\tau}_2) = \text{effects}(\hat{\tau}_1) \cup \text{ho-effects}(\hat{\tau}_2) \end{split}$ 

Figure 3.8: Effect functions.

substitution ::  $\hat{e} \times \hat{v} \times \hat{v} \rightarrow \hat{e}$ 

```
\begin{split} &[\hat{v}/y]x = \hat{v}, \text{ if } x = y \\ &[\hat{v}/y]x = x, \text{ if } x \neq y \\ &[\hat{v}/y](\lambda x : \hat{\tau}.\hat{e}) = \lambda x : \hat{\tau}.[\hat{v}/y]\hat{e}, \text{ if } y \neq x \text{ and } y \text{ does not occur free in } \hat{e} \\ &[\hat{v}/y](\hat{e}_1 \ \hat{e}_2) = ([\hat{v}/y]\hat{e}_1)([\hat{v}/y]\hat{e}_2) \\ &[\hat{v}/y](\hat{e}_1.\pi) = ([\hat{v}/y]e_1).\pi \\ &[\hat{v}/y](\text{import}(\varepsilon) \ x = \hat{e} \text{ in } e) = \text{import}(\varepsilon) \ x = [\hat{v}/y]\hat{e} \text{ in } e \end{split}
```

Figure 3.9: Substitution function.

 $[\hat{v}_1/x_1, \hat{v}_2/x_2]\hat{e}$  as shorthand for  $[\hat{v}_2/x_2]([\hat{v}_1/x_1]\hat{e})$ . Note how the order of the variables has been flipped; the substitutions occur as they are written, left-to-right.

Note that substitution is partial, because it is only defined when a free-variable is being replaced with a value. This is important for proving preservation, because if we replace variables with arbitrary expressions, then we might also be introducing arbitrary effects.

To avoid accidental variable capture we adopt the convention of  $\alpha$ -conversion, whereby we freely and implicitly interchange expressions which are equivalent up to the naming of bound variables [13, p. 71]. This elides some tedious bookkeeping. Consequently, we shall assume variables are (re-)named in this way to avoid accidental capture.

$$\Gamma \vdash e : \tau$$

$$\frac{\Gamma, x : \tau \vdash x : \tau}{\Gamma, x : \tau \vdash x : \tau} \text{ (T-VAR)} \qquad \frac{\Gamma, x : \tau_1 \vdash e : \tau_2}{\Gamma, x : \tau_1 \vdash e : \tau_1 \to \tau_2} \text{ (T-Abs)}$$

$$\frac{\Gamma \vdash e_1 : \tau_2 \to \tau_3 \quad \Gamma \vdash e_2 : \tau_2}{\Gamma \vdash e_1 e_2 : \tau_3} \text{ (T-APP)} \qquad \frac{\Gamma \vdash e : \{\bar{r}\} \quad \forall r \in \bar{r} \mid r \in R \quad \pi \in \Pi}{\Gamma \vdash e.\pi : \text{Unit}} \text{ (T-OPERCALL)}$$

Figure 3.10: Typing judgements in the epsilon calculus.

The first sort of static judgement ascribes a type to a piece of unlabelled code. T-VAR,

T-APP, and T-OPERCALL are the same as they are in  $\lambda^{\rightarrow}$ . T-RESOURCE is the same as T-VAR, but for variables representing primitive capabilities. T-OPERCALL is the rule for typing an operation call  $e_1.\pi$ . Such an expression is well-typed if  $e_1$  types to some valid resource, and  $\pi$  is a known operation.

$$\frac{\operatorname{safe}(\hat{\tau},\varepsilon)}{\operatorname{safe}(\{\bar{r}\},\varepsilon)} \text{ (SAFE-RESOURCE)} \quad \frac{\operatorname{safe}(\operatorname{Unit},\varepsilon)}{\operatorname{safe}(\operatorname{Unit},\varepsilon)} \text{ (SAFE-UNIT)}$$
 
$$\frac{\varepsilon \subseteq \varepsilon' \quad \operatorname{ho-safe}(\hat{\tau}_1,\varepsilon) \quad \operatorname{safe}(\hat{\tau}_2,\varepsilon)}{\operatorname{safe}(\hat{\tau}_1 \to_{\varepsilon'} \hat{\tau}_2,\varepsilon)} \text{ (SAFE-ARROW)}$$
 
$$\frac{\operatorname{ho-safe}(\hat{\tau},\varepsilon)}{\operatorname{ho-safe}(\{\bar{r}\},\varepsilon)} \text{ (HOSAFE-RESOURCE)} \quad \frac{\operatorname{ho-safe}(\operatorname{Unit},\varepsilon)}{\operatorname{ho-safe}(\operatorname{Unit},\varepsilon)} \text{ (HOSAFE-UNIT)}$$
 
$$\frac{\operatorname{safe}(\hat{\tau}_1,\varepsilon) \quad \operatorname{ho-safe}(\hat{\tau}_2,\varepsilon)}{\operatorname{ho-safe}(\hat{\tau}_1 \to_{\varepsilon'} \hat{\tau}_2,\varepsilon)} \text{ (HOSAFE-ARROW)}$$

Figure 3.11: Safety judgements in the epsilon calculus.

Before presenting the type-with-effect rules for labelled expressions, we first define a few safety predicates. Intuitively, the type  $\hat{\tau}$  is safe for  $\varepsilon$  if it has declared every (non higher-order) effect  $r.\pi \in \varepsilon$  in its signature.  $\hat{\tau}$  is ho-safe for  $\varepsilon$  if  $\hat{\tau}$  has declared every higher-order effect  $r.\pi \in \varepsilon$  in its signature. One way to think about these predicates is as a contract between caller and callee. If the caller supplies a set of capabilities  $\varepsilon$  to a piece of code typing to  $\hat{\tau}$ , it would violate the restriction on *ambient authority* if a capability was supplied that  $\hat{\tau}$  had not explicitly asked for. Therefore, safe $(\hat{\tau}, \varepsilon)$  holds when the (non higher-order) effects selected by  $\hat{\tau}$  include  $\varepsilon$ . ho-safe $(\hat{\tau}, \varepsilon)$  holds when the higher-order effects selected by  $\hat{\tau}$  include  $\varepsilon$ .

Because the implementation of  $\hat{\tau}$  might internally propagate capabilities, the definitions of safety and higher-order safety need to be transitive. Give an example of why this is so.

 $\lambda_{\pi,\varepsilon}^{\rightarrow}$  has a new kind of judgement.  $\Gamma \vdash \hat{e} : \hat{\tau}$  with  $\varepsilon$  can be read as saying that  $\hat{e}$ , if it halts, will produce a value of type  $\hat{\tau}$  and incur at most the set of effects  $\varepsilon$ . This judgement gives a conservative approximation as to what will happen; some of the effects in the typing judgement may not actually happen at runtime.

The simplest rules are those which type values as having no effect. Although a function and a resource literal can both capture capabilities, you must do something with them (apply the function, operate on the resource) to incur a runtime effect.

The effects of a lambda application are: the effects of evaluating its subexpressions,

 $\hat{\Gamma} \vdash \hat{e} : \hat{ au} \; \mathtt{with} \; arepsilon \, |$ 

$$\begin{split} \frac{\hat{\Gamma}, x : \tau \vdash x : \tau \text{ with } \varnothing}{\hat{\Gamma}, x : \tau \vdash x : \tau \text{ with } \varnothing} & (\varepsilon\text{-VAR}) \quad \frac{\hat{\Gamma}, r : \{r\} \vdash r : \{r\} \text{ with } \varnothing}{\hat{\Gamma}, r : \{r\} \vdash r : \{r\} \text{ with } \varnothing} & (\varepsilon\text{-RESOURCE}) \\ \frac{\hat{\Gamma}, x : \hat{\tau}_2 \vdash \hat{e} : \hat{\tau}_3 \text{ with } \varepsilon_3}{\hat{\Gamma} \vdash \lambda x : \tau_2. \hat{e} : \hat{\tau}_2 \to_{\varepsilon_3} \hat{\tau}_3 \text{ with } \varnothing} & (\varepsilon\text{-ABS}) \quad \frac{\hat{\Gamma} \vdash \hat{e}_1 : \hat{\tau}_2 \to_{\varepsilon} \hat{\tau}_3 \text{ with } \varepsilon_1 \quad \hat{\Gamma} \vdash \hat{e}_2 : \hat{\tau}_2 \text{ with } \varepsilon_2}{\hat{\Gamma} \vdash \hat{e} : \{\bar{r}\}} & \forall r \in \bar{r} \mid r : \{r\} \in \Gamma \quad \pi \in \Pi} \\ \frac{\hat{\Gamma} \vdash \hat{e} : \{\bar{r}\} \quad \forall r \in \bar{r} \mid r : \{r\} \in \Gamma \quad \pi \in \Pi}{\hat{\Gamma} \vdash \hat{e}.\pi : \text{Unit with } \{\bar{r}.\pi\}} & (\varepsilon\text{-OPERCALL}) \\ \\ \frac{\hat{\Gamma} \vdash e : \tau \text{ with } \varepsilon \quad \tau <: \tau' \quad \varepsilon \subseteq \varepsilon'}{\hat{\Gamma} \vdash e : \tau' \text{ with } \varepsilon'} & (\varepsilon\text{-SUBSUME}) \\ \\ \hat{\Gamma} \vdash \hat{e} : \hat{\tau} \text{ with } \varepsilon_1 \quad \varepsilon = \text{effects}(\hat{\tau}) \\ \\ \frac{\hat{\Gamma} \vdash \hat{e} : \hat{\tau} \text{ with } \varepsilon_1 \quad \varepsilon = \text{erase}(\hat{\tau}) \vdash e : \tau}{\hat{\Gamma} \vdash \text{import}(\varepsilon) \quad x = \hat{e} \text{ in } e : \text{annot}(\tau, \varepsilon) \text{ with } \varepsilon \cup \varepsilon_1} & (\varepsilon\text{-IMPORT}) \end{split}$$

Figure 3.12: Type-with-effect judgements.

and the effects incurred by executing the body of the lambda to which the left-hand side evaluates. Those last effects are pulled from the label on the lambda's arrow-type.

The effects of an operation call are: the effects of evaluating the subexpression, and the single effect incurred when the subexpression is reduced to a resource literal r, and operation  $\pi$  is invoked on it. It is not always possible to know statically which exact resource literal the subexpression reduces to (if it halts at all). Figure 3.8. shows such an example. The safe approximation is to say that the operation call  $\hat{e}.\pi$  incurs  $\pi$  on every possible resource to which  $\hat{e}$  might evaluate. In the case of Figure 3.8., this would be {File.write, Socket.write}.

Because we're not really interested in what exactly an operation call does, every operation is valid for every resource. This can give bizarre programs — Sensor.readTemp seems like a sensible operation call, but what about File.readTemp? — however, an adequate treatment is outside of the scope of  $\lambda_{\pi,\varepsilon}$ , so we gloss over such programs.

It actually might be possible to figure out the exact literal if the system's not Turing complete, since the simply-typed lambda calculus is strongly normalising (and this is basically that, with a few extras), so be careful about this claim

The most interesting rule is  $\varepsilon$ -Import. This rule is set up to ensure the interaction between labelled and unlabelled code is capability-safe. We type e with x: erase( $\hat{\tau}$ ). This elimi-

```
def getResource(b: Bool): { File, Socket } with Ø =

0.9
0.9if b then File else Socket

0.9
4 val boolVal: Bool = System.randomBool
5 getResource(boolVal).write
```

Figure 3.13: We cannot statically determine which branch will execute, so the safe approximation for getResource(boolVal).write is {File.write, Socket.write}.

nates ambient authority, because the only free variables in e will be those selected by the interface  $\hat{\tau}$ .

For our rule to be capability-safe, we need to ensure that any higher-order function in scope is expecting the set of capabilities in  $\hat{\tau}$ . If not, we could exercise ambient authority by passing that higher-order function a capability from  $\hat{\tau}$  which it hadn't selected. This is the purpose of ho-safe( $\hat{\tau}, \varepsilon$ ): all higher-order functions in scope need to be expecting any capability they might be passed.

In the conclusion of the rule we annotate the unlabelled code's effects as  $\mathsf{effects}(\hat{\tau})$ . Because this is the full set of capabilities over which e has access, and because this set is higher-order safe, we shall see this annotation is sound.

$$\frac{\left[\hat{\tau} <: \hat{\tau}\right]}{\varepsilon \subseteq \varepsilon' \quad \hat{\tau}_2 <: \hat{\tau}_2' \quad \hat{\tau}_1' <: \hat{\tau}_1 \\ \hat{\tau}_1 \to_{\varepsilon} \hat{\tau}_2 <: \hat{\tau}_1' \to_{\varepsilon'} \hat{\tau}_2' } \text{ (S-EFFECTS)} \quad \frac{r \in r_1 \implies r \in r_2}{\left\{\bar{r}_1\right\} <: \left\{\bar{r}_2\right\}} \text{ (S-RESOURCES)}$$

Figure 3.14: Subtyping judgements in the epsilon calculus.

In addition to the usual subtyping rules from  $\lambda^{\rightarrow}$  between  $\tau$  terms, we introduce two more for  $\hat{\tau}$  terms.

The rule for functions is contravariant in the input-type and covariant in the output-type (as in  $\lambda^{\rightarrow}$ ), and requires the effects of the super-type to be an upper-bound of the effects of the sub-type. We can think of this in terms of Liskov's substitution principle: if the subtype incurred an effect the supertype hadn't declared, it would violate the supertype's interface.

The rule for resources says that a superset of resources is a subtype.

A single-step reduction takes an expression to a pair consisting of an expression and a set of runtime effects. The rules E-APP1, E-APP2, E-OPERCALL1, E-IMPORT1 all reduce a single subexpression.

E-APP3 is the standard  $\lambda^{\rightarrow}$  rule for applying a value to a lambda, by performing substitution on the lambda body.

E-OPERCALL2 performs an operation on a resource literal. In this case it reduces to

$$\hat{e} \longrightarrow \hat{e} \mid \varepsilon$$

$$\begin{split} \frac{\hat{e}_{1} \longrightarrow \hat{e}'_{1} \mid \varepsilon}{\hat{e}_{1} \hat{e}_{2} \longrightarrow \hat{e}'_{1} \hat{e}_{2} \mid \varepsilon} \; & (\text{E-APP1}) \quad \frac{\hat{e}_{2} \longrightarrow \hat{e}'_{2} \mid \varepsilon}{\hat{v}_{1} \hat{e}_{2} \longrightarrow \hat{v}_{1} \hat{e}'_{2} \mid \varepsilon} \; (\text{E-APP2}) \quad \frac{(\lambda x : \hat{\tau}. \hat{e}) \hat{v}_{2} \longrightarrow [\hat{v}_{2}/x] \hat{e} \mid \varnothing}{(\lambda x : \hat{\tau}. \hat{e}) \hat{v}_{2} \longrightarrow [\hat{v}_{2}/x] \hat{e} \mid \varnothing} \; (\text{E-APP3}) \\ & \frac{\hat{e} \longrightarrow \hat{e}' \mid \varepsilon}{\hat{e}.\pi \longrightarrow \hat{e}'.\pi \mid \varepsilon} \; (\text{E-OPERCALL1}) \quad \frac{r \in R \quad \pi \in \Pi}{r.\pi \longrightarrow \text{unit} \mid \{r.\pi\}} \; (\text{E-OPERCALL2}) \\ & \frac{\hat{e} \longrightarrow \hat{e}' \mid \varepsilon'}{\text{import}(\varepsilon) \; x = \hat{e} \; \text{in} \; e \longrightarrow \text{import}(\varepsilon) \; x = \hat{e}' \; \text{in} \; e \mid \varepsilon'} \; (\text{E-IMPORT1}) \\ & \frac{\text{import}(\varepsilon) \; x = \hat{v} \; \text{in} \; e \longrightarrow [\hat{v}/x] \text{annot}(e, \varepsilon) \mid \varnothing}{\text{import}(\varepsilon) \; x = \hat{v} \; \text{in} \; e \longrightarrow [\hat{v}/x] \text{annot}(e, \varepsilon) \mid \varnothing} \; (\text{E-IMPORT2}) \end{split}$$

Figure 3.15: Single-step reductions.

unit (which is a derived form in our calculus; see 3.4. Encodings). This choice reflects the fact that  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  doesn't model the potentially varied return types of functions.

E-IMPORT2 performs module resolution. The (unlabelled) body of code is annotated with the set of effects captured by the interface, and then the value being imported is substituted into the body of code.

$$\begin{array}{c} \left| \hat{e} \longrightarrow^{*} \hat{e} \mid \varepsilon \right| \\ \\ \overline{\hat{e} \to^{*} \hat{e} \mid \varnothing} \text{ (E-MULTISTEP1)} & \frac{\hat{e} \to \hat{e}' \mid \varepsilon}{\hat{e} \to^{*} \hat{e}' \mid \varepsilon} \text{ (E-MULTISTEP2)} \\ \\ \frac{\hat{e} \to^{*} \hat{e}' \mid \varepsilon_{1} \quad \hat{e}' \to^{*} \hat{e}'' \mid \varepsilon_{2}}{\hat{e} \to^{*} \hat{e}'' \mid \varepsilon_{1} \cup \varepsilon_{2}} \text{ (E-MULTISTEP3)} \end{array}$$

Figure 3.16: Multi-step reductions.

A multi-step reduction consists of zero<sup>2</sup> or more single-step reductions. The resulting effect-set is the union of all the single-steps taken.

#### **3.2.2** Soundness of $\lambda_{\pi,\varepsilon}^{\rightarrow}$

Our goal is to show  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  is sound. This requires an appropriate notion of *effect-soundness*.

**Theorem 6** (Soundness). *If*  $\hat{\Gamma} \vdash \hat{e}_A : \hat{\tau}_A$  with  $\varepsilon_A$  and  $\hat{e}_A$  is not a value, then  $e_A \longrightarrow e_B \mid \varepsilon$ , where  $\hat{\Gamma} \vdash e_B : \hat{\tau}_B$  with  $\varepsilon_B$  and  $\hat{\tau}_B <: \hat{\tau}_A$  and  $\varepsilon_B \cup \varepsilon \subseteq \varepsilon_A$ .

<sup>&</sup>lt;sup>2</sup>We permit multi-step reductions of length zero to be consistent with Pierce, who defines multi-step reduction as a reflexive relation[13, p. 39].

This definition of soundness is the same as in  $\lambda^{\rightarrow}$  but for an extra conclusion:  $\varepsilon_B \cup \varepsilon \subseteq \varepsilon_A$ . Intuitively,  $\varepsilon_A$  is the approximation of what runtime effects the reduction of  $\hat{e}_A$  will incur,  $\varepsilon$  is the actual set of effects  $\hat{e}_A$  incurred (at most a singleton because we are working with single-step reduction), and  $\varepsilon_B$  is the approximation of what runtime effects the reduction of  $\hat{e}_B$  will incur. Evidently we want  $\varepsilon \subseteq \varepsilon_A$ ; an approximation which accounts for every runtime effect is a sound one. We also want  $\varepsilon_B \subseteq \varepsilon_A$ , so successive approximations only get better.

The soundness proof takes the standard approach of showing that progress and preservation hold of the calculus. This can be done immediately by observing some properties that follow immediately from the typing rules.

**Lemma 2** (Canonical Forms). *The following are true:* 

```
• If \hat{\Gamma} \vdash \hat{v} : \hat{\tau} with \varepsilon then \varepsilon = \emptyset.
• If \hat{\Gamma} \vdash \hat{v} : \{\bar{r}\} then \hat{v} = r for some r \in R and \{\bar{r}\} = \{r\}.
```

**Theorem 7** (Progress). *If*  $\hat{\Gamma} \vdash \hat{e} : \hat{\tau}$  with  $\varepsilon$  and  $\hat{e}$  is not a value, then  $\hat{e} \longrightarrow \hat{e}' \mid \varepsilon$ .

*Proof.* By induction on  $\hat{\Gamma} \vdash \hat{e} : \hat{\tau}$  with  $\varepsilon$ , for  $\hat{e}$  not a value. If the rule is  $\varepsilon$ -Subsumption it follows by inductive hypothesis. If  $\hat{e}$  has a reducible subexpression then reduce it. Otherwise use one of  $\varepsilon$ -APP3,  $\varepsilon$ -OPERCALL2, or  $\varepsilon$ -IMPORT2.

To prove preservation, we need to know types and effects are preserved under substitution. The substitution lemma gives us this result. It says that if x is bound to a type, and a value  $\hat{v}$  of that type is substituted into  $\hat{e}$ , then the type and effect of  $\hat{e}$  remain unchanged. Key to this property is that  $\hat{v}$  is a value, so by canonical forms it cannot introduce effects that weren't already in  $\hat{e}$ . Beyond this observation, the proof is routine.

**Lemma 3** (Substitution). If  $\hat{\Gamma}, x : \hat{\tau}' \vdash e : \hat{\tau}$  with  $\varepsilon$  and  $\hat{\Gamma} \vdash \hat{v} : \hat{\tau}'$  with  $\varnothing$  then  $\hat{\Gamma} \vdash [\hat{v}/x]e : \hat{\tau}$  with  $\varepsilon$ .

```
Proof. By induction on \hat{\Gamma}, x : \hat{\tau}' \vdash e : \hat{\tau} with \varepsilon.
```

The tricky case in preservation is when an import expression is resolved. To show the reduction  $\operatorname{import}(\varepsilon) \ x = \hat{v} \ \operatorname{in} \ e \longrightarrow [\hat{v}/x] \operatorname{annot}(e,\varepsilon) \mid \varnothing$  preserves soundness requires a few things. First, if  $\hat{\Gamma} \vdash \operatorname{import}(\varepsilon) \ x = \hat{v} \ \operatorname{in} \ e : \hat{\tau}_A \ \operatorname{with} \ \varepsilon_A$ , then we need to be able to type the reduced expression in the same context:  $\hat{\Gamma} \vdash [\hat{v}/x] \operatorname{annot}(e,\varepsilon) : \hat{\tau}_B \ \operatorname{with} \ \varepsilon_B$ . To be effect-sound, we need  $\varepsilon_B \subseteq \varepsilon_A$ . To be type-sound, we need  $\hat{\tau}_B <: \hat{\tau}_A$ . This motivates the next lemma, which relates a typing judgement of e to a typing judgement of e annot e.

**Lemma 4** (Annotation). *If the following are true:* 

```
• \hat{\Gamma} \vdash \hat{v} : \hat{\tau} \text{ with } \varnothing
• \Gamma, y : \operatorname{erase}(\hat{\tau}) \vdash e : \tau
• \varepsilon = \operatorname{effects}(\hat{\tau})
```

27

• ho-safe $(\hat{\tau}, \varepsilon)$ 

 $Then \ \hat{\Gamma}, \mathtt{annot}(\Gamma, \varepsilon), y : \hat{\tau} \vdash \mathtt{annot}(e, \varepsilon) : \mathtt{annot}(\tau, \varepsilon) \ \mathtt{with} \ \varepsilon \cup \mathtt{effects}(\mathtt{annot}(\Gamma, \varepsilon)).$ 

*Proof.* By induction on  $\Gamma, y : erase(\hat{\tau}) \vdash e : \tau$ .

The exact formulation of the Annotation lemma is very specific to the premises of  $\varepsilon$ -IMPORT2, but generalised slightly to accommodate a proof by induction. The generalisation is to allow e to be typed in any context  $\Gamma$  with a binding for y. We can think of  $\Gamma$  as encapsulating the ambient authority exercised by e. At the top-level of any program, we will always have  $\Gamma = \varnothing$ , because the typing judgement  $\varepsilon$ -IMPORT always types import expressions with just the authority being selected. However, inductively-speaking, there may be ambient capabilities. Consider  $(\lambda x : \{\text{File}\}. \text{ x.write})$  File. From the perspective of x.write, File is an ambient capability, and so if we were to inductively apply the Annotation lemma, at this point, File  $\in \Gamma$ . However, because the code encapsulating x.write selects File by binding it to x in the function declaration, this code is capability-safe.

Proving the annotation lemma requires an additional pair of lemmas, to relate  $\hat{\tau}$  and  $\mathtt{annot}(\mathtt{erase}(\hat{\tau}), \varepsilon)$ .

**Lemma 5.** If effects( $\hat{\tau}$ )  $\subseteq \varepsilon$  and ho-safe( $\hat{\tau}, \varepsilon$ ) then  $\hat{\tau} <: annot(erase(\hat{\tau}), \varepsilon)$ .

**Lemma 6.** If ho-effects( $\hat{\tau}$ )  $\subseteq \varepsilon$  and safe( $\hat{\tau}, \varepsilon$ ) then annot(erase( $\hat{\tau}$ ),  $\varepsilon$ ) <:  $\hat{\tau}$ .

*Proof.* By simultaneous induction on ho-safe and safe.

There is a close relation between these lemmas and the subtyping rule for functions. In a subtyping relation between functions, the input type is contravariant. Therefore, if  $\hat{\tau} = \hat{\tau}_1 \rightarrow_{\varepsilon'} \tau_2$  and we have  $\hat{\tau} <: \operatorname{annot}(\tau, \varepsilon)$ , then we need to know  $\operatorname{annot}(\tau_1) <: \hat{\tau}_1$ . This is why there are two lemmas, one for each direction.

Armed with the annotation lemma, we are now ready to prove the preservation theorem.

**Theorem 8** (Preservation). *If*  $\hat{\Gamma} \vdash \hat{e}_A : \hat{\tau}_A$  with  $\varepsilon_A$  and  $\hat{e}_A \longrightarrow \hat{e}_B \mid \varepsilon$ , then  $\hat{\tau}_B <: \hat{\tau}_A$  and  $\varepsilon_B \cup \varepsilon \subseteq \varepsilon_A$ .

*Proof.* By induction on  $\hat{\Gamma} \vdash \hat{e}_A : \hat{\tau}_A$  with  $\varepsilon_A$ , and then on  $\hat{e}_A \longrightarrow \hat{e}_B \mid \varepsilon$ .

Case:  $\varepsilon$ -APP Then  $e_A = \hat{e}_1 \hat{e}_2$  and  $\hat{e}_1 : \hat{\tau}_2 \to_{\varepsilon} \hat{\tau}_3$  with  $\varepsilon_1$  and  $\hat{\Gamma} \vdash \hat{e}_2 : \hat{\tau}_2$  with  $\varepsilon_2$ . If the reduction rule used was E-APP1 or E-APP2, then the result follows by applying the inductive hypothesis to  $\hat{e}_1$  and  $\hat{e}_2$  respectively.

Otherwise the rule used was E-APP3. Then  $(\lambda x: \hat{\tau}_2.\hat{e})\hat{v}_2 \longrightarrow [\hat{v}_2/x]\hat{e} \mid \varnothing$ . By inversion on the typing rule for  $\lambda x: \hat{\tau}_2.\hat{e}$  we know  $\Gamma, x: \hat{\tau}_2 \vdash \hat{e}: \hat{\tau}_3$  with  $\varepsilon_3$ . By canonical forms,  $\varepsilon_2 = \varnothing$  because  $\hat{e}_2 = \hat{v}_2$  is a value. Then by the substitution lemma,  $\hat{\Gamma} \vdash [\hat{v}_2/x]\hat{e}: \hat{\tau}_3$  with  $\varepsilon_3$ .

By canonical forms,  $\varepsilon_1 = \varepsilon_2 = \emptyset = \varepsilon_C$ . Therefore  $\varepsilon_A = \varepsilon_3 = \varepsilon_B \cup \varepsilon_C$ .

Case:  $\varepsilon$ -OPERCALL. Then  $e_A = e_1.\pi$  and  $\hat{\Gamma} \vdash e_1 : \{\bar{r}\}\$  with  $\varepsilon_1$ . If the reduction rule used was E-OPERCALL1 then the result follows by applying the inductive hypothesis to  $\hat{e}_1$ .

Otherwise the reduction rule used was E-OPERCALL2 and  $v_1.\pi \longrightarrow \text{unit} \mid \{r.\pi\}$ . By canonical forms,  $\hat{\Gamma} \vdash v_1 : \text{unit with } \{r.\pi\}$ . Also,  $\hat{\Gamma} \vdash \text{unit} : \text{Unit with } \emptyset$ . Then  $\tau_B = \tau_A$ . Also,  $\varepsilon_C \cup \varepsilon_B = \{r.\pi\} = \varepsilon_A$ .

Case:  $\varepsilon$ -IMPORT. Then  $e_A = \mathrm{import}(\varepsilon) \ x = \hat{e} \ \mathrm{in} \ e$ . If the reduction rule used was E-IMPORT1 then the result follows by applying the inductive hypothesis to  $\hat{e}$ .

Otherwise  $\hat{e}$  is a value and the reduction used was E-IMPORT2. The following are true:

```
1. e_A = \operatorname{import}(\varepsilon) \ x = \hat{v} \ \operatorname{in} \ e
2. \hat{\Gamma} \vdash e_A : \operatorname{annot}(\tau, \varepsilon) \ \operatorname{with} \ \varepsilon \cup \varepsilon_1
3. \operatorname{import}(\varepsilon) \ x = \hat{v} \ \operatorname{in} \ e \longrightarrow [\hat{v}/x] \operatorname{annot}(e, \varepsilon) \mid \varnothing
4. \hat{\Gamma} \vdash \hat{v} : \hat{\tau} \ \operatorname{with} \ \varnothing
5. \varepsilon = \operatorname{effects}(\hat{\tau})
6. \operatorname{ho-safe}(\hat{\tau}, \varepsilon)
7. x : \operatorname{erase}(\hat{\tau}) \vdash e : \tau
```

Apply the annotation lemma with  $\Gamma = \emptyset$  to get  $\hat{\Gamma}, x : \hat{\tau} \vdash \mathtt{annot}(e, \varepsilon) : \mathtt{annot}(\tau, \varepsilon)$  with  $\varepsilon$ . From assumption (4) we know  $\hat{\Gamma} \vdash \hat{v} : \hat{\tau}$  with  $\emptyset$ , and so the substitution lemma may be applied, giving  $\hat{\Gamma} \vdash [\hat{v}/x]\mathtt{annot}(e, \varepsilon) : \mathtt{annot}(\tau, \varepsilon)$  with  $\varepsilon$ . By canonical forms,  $\varepsilon_1 = \varepsilon_C = \emptyset$ . Then  $\varepsilon_B = \varepsilon = \varepsilon_A \cup \varepsilon_C$ . By examination,  $\tau_A = \tau_B = \mathtt{annot}(\tau, \varepsilon)$ .

Our statement of soundness combines the progress and preservation theorems into one.

**Theorem 9** (Soundness). *If*  $\hat{\Gamma} \vdash \hat{e}_A : \hat{\tau}_A$  with  $\varepsilon_A$  and  $\hat{e}_A$  is not a value, then  $e_A \longrightarrow e_B \mid \varepsilon$ , where  $\hat{\Gamma} \vdash e_B : \hat{\tau}_B$  with  $\varepsilon_B$  and  $\hat{\tau}_B <: \hat{\tau}_A$  and  $\varepsilon_B \cup \varepsilon \subseteq \varepsilon_A$ .

*Proof.* If  $\hat{e}_A$  is not a value then the reduction exists by the progress theorem. The rest follows by the preservation theorem.

Knowing that single-step reductions are sound, multi-step reductions can straightforwardly be be shown to also be sound. This is done by inductively applying single-step soundness to the length of the multi-step reduction.

**Theorem 10** (Multi-step Soundness). If  $\hat{\Gamma} \vdash \hat{e}_A : \hat{\tau}_A$  with  $\varepsilon_A$  and  $e_A \longrightarrow^* e_B \mid \varepsilon$ , where  $\hat{\Gamma} \vdash e_B : \hat{\tau}_B$  with  $\varepsilon_B$  and  $\hat{\tau}_B <: \hat{\tau}_A$  and  $\varepsilon_B \cup \varepsilon \subseteq \varepsilon_A$ .

*Proof.* By induction on the length of the multi-step reduction. If the length is 0 then  $e_A = e_B$  and the result holds vacuously. If the length is 1 the result holds by soundness of single-step reductions. if the length is n + 1, then the first n-step reduction is sound by inductive hypothesis and the last step is sound by single-step soundness, so the entire n + 1-step reduction is sound.

# Chapter 4

# **Applications**

### 4.1 Encodings

When writing practical examples it is useful to use higher-level constructs which have been derived from the base language. In this section we introduce some of the constructs that we use in examples. Because the core language is sound, any derived extension is also sound.

#### 4.1.1 Unit

Unit is a type inhabited by exactly one value. It conveys the absence of information. In our dynamic rules, unit is what an operation call on a resource literal is reduced to. We define unit  $\stackrel{\text{def}}{=} \lambda x : \varnothing.x$  and Unit  $\stackrel{\text{def}}{=} \varnothing \to_{\varnothing} \varnothing$ . Note that because there is no empty resource literal, unit cannot be applied to anything. Furthermore,  $\vdash$  unit : Unit with  $\varnothing$ , by  $\varepsilon$ -ABS, so any context can make this type judgement.

```
\frac{\Gamma \vdash e : \tau}{\hat{\Gamma} \vdash \hat{e} : \hat{\tau} \text{ with } \varepsilon}
\frac{\Gamma \vdash \text{unit} : \text{Unit}}{\Gamma \vdash \text{unit} : \text{Unit}} \text{ (T-UNIT)} \quad \frac{\hat{\Gamma} \vdash \text{unit} : \text{Unit with } \varnothing}{\hat{\Gamma} \vdash \text{unit} : \text{Unit with } \varnothing} \text{ ($\varepsilon$-UNIT)}
```

Figure 4.1: Derived Unit rules.

#### 4.1.2 Let

The expression let  $x = \hat{e}_1$  in  $\hat{e}_2$  first binds the value  $\hat{e}_1$  to the name x and then evaluates  $\hat{e}_2$ . We can generalise by allowing  $\hat{e}_1$  to be a non-value, in which case it must first be reduced to a value. If  $\Gamma \vdash \hat{e}_1 : \hat{\tau}_1$ , then let  $x = \hat{e}_1$  in  $\hat{e}_2 \stackrel{\text{def}}{=} (\lambda x : \hat{\tau}_1.\hat{e}_2)\hat{e}_1$ . Note that if  $\hat{e}_1$  is a non-value, we can reduce the let by E-APP2. If  $\hat{e}_1$  is a value, we may apply E-APP3,

which binds  $\hat{e}_1$  to x in  $\hat{e}_2$ . This is fundamentally a lambda application, so it can be typed using  $\varepsilon$ -APP (or T-APP, if the terms involved are unlabelled).

$$\begin{split} & \Gamma \vdash e : \tau \\ & \hat{\Gamma} \vdash \hat{e} : \hat{\tau} \text{ with } \varepsilon \\ & \hat{e} \to \hat{e} \mid \varepsilon \end{split}$$
 
$$& \frac{\Gamma \vdash e_1 : \tau_1 \quad \Gamma, x : \tau_1 \vdash e_2 : \tau_2}{\Gamma \vdash \text{let } x = e_1 \text{ in } e_2 : \tau_2} \text{ ($\varepsilon$-LET)} \\ & \frac{\hat{\Gamma} \vdash \hat{e}_1 : \hat{\tau}_1 \text{ with } \varepsilon_1 \quad \hat{\Gamma}, x : \hat{\tau}_1 \vdash \hat{e}_2 : \hat{\tau}_2 \text{ with } \varepsilon_2}{\hat{\Gamma} \vdash \text{let } x = \hat{e}_1 \text{ in } \hat{e}_2 : \hat{\tau}_2 \text{ with } \varepsilon_1 \cup \varepsilon_2} \text{ ($\varepsilon$-LET)} \\ & \frac{\hat{e}_1 \longrightarrow \hat{e}_1' \mid \varepsilon_1}{\text{let } x = \hat{e}_1 \text{ in } \hat{e}_2 \longrightarrow \text{let } x = \hat{e}_1' \text{ in } \hat{e}_2 \mid \varepsilon_1} \text{ ($\varepsilon$-LET1)} \\ & \frac{1}{\text{let } x = \hat{e}_1 \text{ in } \hat{e}_2 \longrightarrow \text{let } x = \hat{e}_1' \text{ in } \hat{e}_2 \mid \varepsilon_1} \text{ ($\varepsilon$-LET2)} \end{split}$$

Figure 4.2: Derived 1et rules.

#### 4.1.3 Conditionals

#### **4.1.4** Tuples

We need tuples to import multiple names.

### 4.2 Examples

EXAMPLE OF A RESOURCE LEAKING AND BREAKING CONFINEMENT EXAMPLE OF IMPORTING MULTIPLE CAPABILITIES, ONE GETS LEAKED AND PASSED SOMEWHERE IT HASN'T BEEN SELECTED

4.2. EXAMPLES 31

Figure 4.3: A logger client doesn't need to add effect labels. These can be inferred.

```
resource module Logger
require File
comparison of the state of th
```

Figure 4.4: This won't type because of a mismatch between the client's effects and the logger's effects.

# Chapter 5

### **Evaluation**

#### 5.1 Related Work

Fengyun Liu has approached the study of capability-based effect systems by developing a lambda calculus based around two type-constructors for building free and stoic functions [7]. Free functions may ambiently capture capabilities, but stoic functions may not; for a stoic function to have any effect, it must be explicitly given the capability for that effect. The resulting theory allows the type system to determine if a stoic function is pure or not by inspecting its parameters. If a function is known to be pure there are many optimisations that can be made (inlining, parallelisation). Liu's work is largely motivated by achieving such optimisations for Scala compilers.

By contrast, our work is motivated by the propagation and use of capabilities, and how language-design features might inform software design. Unlike Liu's System F-Impure,  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  has no effect-polymorphism. However, our work has more fine-grained detail about those effects incurred by a particular function — while System F-Impure can conclusively determine if a stoic function is pure, determining what particular effects an impure function has is outside of the scope of Liu's work.

### 5.2 Future Work

A major limitation to practical adoption of  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  is that it is not Turing complete — it has no general recursion, nor recursive types. Extending  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  to include these features would bring it up to par with real programming languages.

Miller's formulation of the capability-model is in terms of objects, and all of the capability-safe languages to which this paper has referred are object-oriented. It is worth investigating how the bridge between  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  and existing capability-safe languages might be bridged by investigating different object encodings, and determining which language extensions are needed to enable these. By extension, these languages have first-class modules, so a version of  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  which can reason about objects would immediately yield

module-level reasoning.

The biggest contribution that could be made to  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  would be to enrich it with a theory of polymorphic effects. As an example, consider  $\lambda x: \mathtt{Unit} \rightarrow_{\varepsilon} \mathtt{Unit}. x \mathtt{unit}$ , where  $\varepsilon$  is free. Invoking this particular function would incur every effect in  $\varepsilon$ , but allowing general. Currently  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  has no way to define such functions which are parametrised by effect-sets. Deveoping an extension which can handle polymorphic effects would be a valuable contribution, and improve the stock of  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  as a practical type-and-effect system.

### 5.3 Conclusion

 $\lambda_{\pi,\varepsilon}^{\rightarrow}$  is an extension to  $\lambda^{\rightarrow}$  which allows for the import of capabilities into unlabelled code. This importing is done in a capability-safe manner, which prohibits the exercise of ambient authority. As a result, we can safely bound the set of possible effects in the unlabelled code by inspecting those capabilities passed into it via the import expression.

Talk about examples given, mention any extensions needed to allow for things such as multiple imports.

There are some important limitations to  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ : it has no general recursion, and no recursive types; it is formulated in terms of the lambda calculus, whereas the capability model is stated in terms of objects; it has no way to express functions with polymorphic effects. These are all interesting avenues of future work that would enrich  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  and our collective understanding of the relation between effects and capabilities.

# Appendix A

### **Proofs**

**Lemma 7** (Canonical Forms). *The following are true:* 

```
• If \hat{\Gamma} \vdash \hat{v} : \hat{\tau} with \varepsilon then \varepsilon = \emptyset.
```

• If  $\hat{\Gamma} \vdash \hat{v} : \{\bar{r}\}$  then  $\hat{v} = r$  for some  $r \in R$  and  $\{\bar{r}\} = \{r\}$ .

**Theorem 11** (Progress). *If*  $\hat{\Gamma} \vdash \hat{e} : \hat{\tau}$  with  $\varepsilon$  and  $\hat{e}$  is not a value, then  $\hat{e} \longrightarrow \hat{e}' \mid \varepsilon$ .

*Proof.* By induction on  $\hat{\Gamma} \vdash \hat{e} : \hat{\tau}$  with  $\varepsilon$ , for  $\hat{e}$  not a value.

Case:  $\varepsilon$ -APP. Then  $\hat{e}=\hat{e}_1$   $\hat{e}_2$ . If  $\hat{e}_1$  is a non-value, then  $\hat{e}_1$   $\hat{e}_2$   $\longrightarrow$   $\hat{e}'_1$   $\hat{e}_2$  by E-APP1. If  $\hat{e}_1=\hat{v}_1$  is a value and  $\hat{e}_2$  is a non-value, then  $\hat{e}_1$   $\hat{e}_2$   $\longrightarrow$   $\hat{v}_1$   $\hat{e}'_2$  by E-APP2. Otherwise  $\hat{e}_1$  and  $\hat{e}_2$  are both values. By inversion,  $\hat{e}_1=\lambda x:\hat{\tau}.\hat{e}$ , so  $(\lambda x:\hat{\tau}.\hat{e})\hat{v}_2\longrightarrow [\hat{v}_2/x]\mid\varnothing$  by E-APP3.

Case:  $\varepsilon$ -OPER. Then  $\hat{e} = \hat{e}_1.\pi$ . If  $\hat{e}_1$  is a non-value, then  $\hat{e}_1.\pi \longrightarrow \hat{e}'_1.\pi \mid \varepsilon_1$  by E-OPERCALL1. Otherwise  $\hat{e}_1 = \hat{v}_1$  is a value. By canonical forms,  $\hat{v}_1 = r$  and  $\hat{\Gamma} \vdash v_1$ :  $\{r\}$  with  $\emptyset$ . Then  $r.\pi \longrightarrow \text{unit} \mid \{r.\pi\}$  by E-OPERCALL2.

Case:  $\varepsilon$ -Subsume. Then  $\hat{\Gamma} \vdash \hat{e} : \hat{\tau}'$  with  $\varepsilon'$ . By inversion,  $\hat{\Gamma} \vdash \hat{e} : \tau$  with  $\varepsilon$ , where  $\tau' <: \tau$  and  $\varepsilon' \subseteq \varepsilon$ . These are subderivations, so the result holds by inductive assumption.

Case:  $\varepsilon$ -MODULE. Then  $\hat{e} = \mathrm{import}(\varepsilon) \, x = \hat{e}' \, \mathrm{in} \, e$ . If  $\hat{e}'$  is a non-value then  $\mathrm{import}(\varepsilon) \, x = \hat{e}' \, \mathrm{in} \, e \longrightarrow \mathrm{import}(\varepsilon) \, x = \hat{e}'' \, \mathrm{in} \, e \mid \varepsilon' \, \mathrm{by} \, \mathrm{E}$ -MODULE1. Otherwise  $\hat{e}' = \hat{v}$  is a value. Then  $\mathrm{import}(\varepsilon) \, x = \hat{v} \, \mathrm{in} \, e \longrightarrow [\hat{v}/x] \mathrm{annot}(e,\varepsilon) \mid \varnothing \, \mathrm{by} \, \mathrm{E}$ -MODULE2.

**Lemma 8** (Substitution). If  $\hat{\Gamma}, x : \hat{\tau}' \vdash e : \hat{\tau}$  with  $\varepsilon$  and  $\hat{\Gamma} \vdash \hat{v} : \hat{\tau}'$  with  $\varnothing$  then  $\hat{\Gamma} \vdash [\hat{v}/x]e : \hat{\tau}$  with  $\varepsilon$ .

*Proof.* By induction on  $\hat{\Gamma}, x : \hat{\tau}' \vdash e : \hat{\tau}$  with  $\varepsilon$ .

Case:  $\varepsilon$ -VAR. Then  $\hat{e} = y$  and either y = x or  $y \neq x$ . If  $y \neq x$ . Then  $[\hat{v}/x]y = y$  and  $\hat{\Gamma} \vdash y : \hat{\tau}$  with  $\varnothing$ . Therefore  $\hat{\Gamma} \vdash [\hat{v}/x]y : \hat{\tau}$  with  $\varnothing$ . Otherwise y = x. By inversion on  $\varepsilon$ -VAR, the typing judgement from the theorem assumption is  $\hat{\Gamma}, x : \hat{\tau}' \vdash x : \hat{\tau}'$  with  $\varnothing$ . Since  $[\hat{v}/x]y = \hat{v}$ , and by assumption  $\hat{\Gamma} \vdash \hat{v} : \hat{\tau}'$  with  $\varnothing$ , then  $\hat{\Gamma} \vdash [\hat{v}/x]x : \hat{\tau}'$  with  $\varnothing$ .

Case:  $\varepsilon$ -RESOURCE. Because  $\hat{e}=r$  is a resource literal then  $\hat{\Gamma}\vdash r:\hat{\tau}$  with  $\varnothing$  by canonical forms. By definition  $[\hat{v}/x]r=r$ , so  $\hat{\Gamma}\vdash [\hat{v}/x]r:\hat{\tau}$  with  $\varnothing$ .

Case:  $\varepsilon$ -APP By inversion we know  $\hat{\Gamma}, x: \hat{\tau}' \vdash \hat{e}_1: \hat{\tau}_2 \to_{\varepsilon_3} \hat{\tau}_3$  with  $\varepsilon_A$  and  $\hat{\Gamma}, x: \hat{\tau}' \vdash \hat{e}_2: \hat{\tau}_2$  with  $\varepsilon_B$ , where  $\varepsilon = \varepsilon_A \cup \varepsilon_B \cup \varepsilon_3$  and  $\hat{\tau} = \hat{\tau}_3$ . By inductive assumption,  $\hat{\Gamma} \vdash [\hat{v}/x]\hat{e}_1: \hat{\tau}_2 \to_{\varepsilon_3} \hat{\tau}_3$  with  $\varepsilon_A$  and  $\hat{\Gamma} \vdash [\hat{v}/x]\hat{e}_2: \hat{\tau}_2$  with  $\varepsilon_B$ . By  $\varepsilon$ -APP we have  $\hat{\Gamma} \vdash ([\hat{v}/x]\hat{e}_1)([\hat{v}/x]\hat{e}_2): \hat{\tau}_3$  with  $\varepsilon_A \cup \varepsilon_B \cup \varepsilon_3$ . By simplifying and applying the definition of substitution, this is the same as  $\hat{\Gamma} \vdash [\hat{v}/x](\hat{e}_1\hat{e}_2): \hat{\tau}$  with  $\varepsilon$ .

Case:  $\varepsilon$ -OPERCALL By inversion we know  $\hat{\Gamma}, x: \hat{\tau}' \vdash \hat{e}_1: \{\bar{r}\}$  with  $\varepsilon_1$ , where  $\varepsilon = \varepsilon_1 \cup \{r.\pi \mid r.\pi \in \bar{r} \times \Pi\}$  and  $\hat{\tau} = \{\bar{r}\}$ . By applying the inductive assumption,  $\hat{\Gamma} \vdash [\hat{v}/x]\hat{e}_1: \{\bar{r}\}$  with  $\varepsilon_1$ . Then by  $\varepsilon$ -OPERCALL,  $\hat{\Gamma} \vdash ([\hat{v}/x]\hat{e}_1).\pi: \{\bar{r}\}$  with  $\varepsilon_1 \cup \{r.\pi \mid r.\pi \in \bar{r} \times \Pi\}$ . By simplifying and applying the definition of substitution, this is the same as  $\hat{\Gamma} \vdash [\hat{v}/x](\hat{e}_1.\pi): \hat{\tau}$  with  $\varepsilon$ .

Case:  $\varepsilon$ -Subsume By inversion we know  $\hat{\Gamma}, x : \hat{\tau}' \vdash \hat{e} : \hat{\tau}_2$  with  $\varepsilon_2$ , where  $\hat{\tau}_2 <: \hat{\tau}$  and  $\varepsilon_2 \subseteq \varepsilon$ . By inductive hypothesis,  $\hat{\Gamma} \vdash [\hat{v}/x]\hat{e} : \hat{\tau}_2$  with  $\varepsilon_2$ . Then by  $\varepsilon$ -Subsume we get  $\hat{\Gamma} \vdash [\hat{v}/x]\hat{e} : \hat{\tau}$  with  $\varepsilon$ .

Case:  $\varepsilon$ -MODULE Then  $\hat{\Gamma}, x : \hat{\tau}' \vdash \text{import}(:) = annot \text{ in } (\tau, \varepsilon) \text{ with } \varepsilon \cup \varepsilon_1$ . By inversion we know  $\hat{\Gamma}, x : \hat{\tau}' \vdash \hat{e} : \hat{\tau}_1 \text{ with } \varepsilon_1$ . By inductive assumption,  $\hat{\Gamma} \vdash [\hat{v}/x]\hat{e} : \hat{\tau}_1 \text{ with } \varepsilon_1$ . Then by  $\varepsilon$ -MODULE we have  $\hat{\Gamma} \vdash \text{import}(:) = annot \text{ in } (\tau, \varepsilon) \text{ with } \varepsilon \cup \varepsilon_1$ .

**Lemma 9.** If  $effects(\hat{\tau}) \subseteq \varepsilon$  and  $ho-safe(\hat{\tau}, \varepsilon)$  then  $\hat{\tau} <: annot(erase(\hat{\tau}), \varepsilon)$ .

**Lemma 10.** If ho-effects( $\hat{\tau}$ )  $\subseteq \varepsilon$  and safe( $\hat{\tau}, \varepsilon$ ) then annot(erase( $\hat{\tau}$ ),  $\varepsilon$ ) <:  $\hat{\tau}$ .

*Proof.* By simultaneous induction.

Case:  $\hat{\tau} = \{\bar{r}\}\ \text{Then } \hat{\tau} = \mathtt{annot}(\mathtt{erase}(\hat{\tau}), \varepsilon)$  and the results for both lemmas hold immediately.

Case:  $\hat{\tau} = \hat{\tau}_1 \rightarrow_{\varepsilon'} \hat{\tau}_2$ , effects $(\hat{\tau}) \subseteq \varepsilon$ , ho-safe $(\hat{\tau}, \varepsilon)$  It is sufficient to show  $\hat{\tau}_2 <:$  annot $(erase(\hat{\tau}_2), \varepsilon)$  and annot $(erase(\hat{\tau}_1), \varepsilon) <: \hat{\tau}_1$ , because the result will hold by S-EFFECTS. To achieve this we shall inductively apply lemma 2 to  $\hat{\tau}_2$  and lemma 3 to  $\hat{\tau}_1$ .

From  $\operatorname{effects}(\hat{\tau}) \subseteq \varepsilon$  we have  $\operatorname{ho-effects}(\hat{\tau}_1) \cup \varepsilon' \cup \operatorname{effects}(\hat{\tau}_2) \subseteq \varepsilon$  and therefore  $\operatorname{effects}(\hat{\tau}_2) \subseteq \varepsilon$ . From  $\operatorname{ho-safe}(\hat{\tau}, \varepsilon)$  we have  $\operatorname{ho-safe}(\hat{\tau}_2, \varepsilon)$ . Therefore we can apply lemma 2 to  $\hat{\tau}_2$ .

From  $\operatorname{effects}(\hat{\tau}) \subseteq \varepsilon$  we have  $\operatorname{ho-effects}(\hat{\tau}_1) \cup \varepsilon' \cup \operatorname{effects}(\hat{\tau}_2) \subseteq \varepsilon$  and therefore  $\operatorname{ho-effects}(\hat{\tau}_1) \subseteq \varepsilon$ . From  $\operatorname{ho-safe}(\hat{\tau}, \varepsilon)$  we have  $\operatorname{ho-safe}(\hat{\tau}_1, \varepsilon)$ . Therefore we can apply lemma 3 to  $\hat{\tau}_1$ .

Case:  $\hat{\tau} = \hat{\tau}_1 \rightarrow_{\varepsilon'} \hat{\tau}_2$ , ho-effects $(\hat{\tau}) \subseteq \varepsilon$ , safe $(\hat{\tau}, \varepsilon)$  It is sufficient to show annot(erase $(\hat{\tau}_2), \varepsilon) <$ :  $\hat{\tau}_2$  and  $\hat{\tau}_1 <$ : annot(erase $(\hat{\tau}_1), \varepsilon)$ , because the result will hold by S-EFFECTS. To achieve this we shall inductively apply lemma 3 to  $\hat{\tau}_2$  and lemma 2 to  $\hat{\tau}_1$ .

From ho-effects( $\hat{\tau}$ )  $\subseteq \varepsilon$  we have effects( $\hat{\tau}_1$ )  $\cup$  ho-effects( $\hat{\tau}_2$ )  $\subseteq \varepsilon$  and therefore ho-effects( $\hat{\tau}_2$ )  $\subseteq \varepsilon$ . From safe( $\hat{\tau}, \varepsilon$ ) we have safe( $\hat{\tau}_2, \varepsilon$ ). Therefore we can apply **lemma** 3 to  $\hat{\tau}_2$ .

From ho-effects( $\hat{\tau}$ )  $\subseteq \varepsilon$  we have effects( $\hat{\tau}_1$ )  $\cup$  ho-effects( $\hat{\tau}_2$ )  $\subseteq \varepsilon$  and therefore effects( $\hat{\tau}_1$ )  $\subseteq \varepsilon$ . From safe( $\hat{\tau}, \varepsilon$ ) we have ho-safe( $\hat{\tau}_1, \varepsilon$ ). Therefore we can apply **lemma** 2 to  $\hat{\tau}_1$ .

**Lemma 11** (Annotation). *If the following are true:* 

- $\bullet \ \hat{\Gamma} \vdash \hat{v} : \hat{\tau} \ \mathtt{with} \ \varnothing$
- $\bullet \ \Gamma, y : \mathtt{erase}(\hat{\tau}) \vdash e : \tau$
- $\varepsilon = \texttt{effects}(\hat{\tau})$
- ho-safe $(\hat{\tau}, \varepsilon)$

 $Then \ \hat{\Gamma}, \mathtt{annot}(\Gamma, \varepsilon), y : \hat{\tau} \vdash \mathtt{annot}(e, \varepsilon) : \mathtt{annot}(\tau, \varepsilon) \ \mathtt{with} \ \varepsilon \cup \mathtt{effects}(\mathtt{annot}(\Gamma, \varepsilon)).$ 

*Proof.* By induction on  $\Gamma, y : \text{erase}(\hat{\tau}) \vdash e : \tau$ .

*Case*: T-VAR Then e = x and  $\Gamma, y : erase(\hat{\tau}) \vdash x : \tau$ . Either x = y or  $x \neq y$ .

**Subcase 1:** x=y. Then by  $\varepsilon\textsc{-VAR}$  we get  $\hat{\Gamma}, \mathtt{annot}(\Gamma, \varepsilon), y: \hat{\tau} \vdash x: \hat{\tau} \ \mathtt{with} \ \varnothing$ . First note that  $\mathtt{annot}(x,\varepsilon) = x$  in this case. Therefore  $\Gamma, y: \mathtt{erase}(\hat{\tau}) \vdash \mathtt{annot}(\mathtt{erase}(x),\varepsilon): \hat{\tau} \ \mathtt{with} \ \varnothing$ . We know by assumption that  $\mathtt{effects}(\hat{\tau}) = \varepsilon$  and  $\mathtt{ho\textsc{-safe}}(\hat{\tau},\varepsilon)$ . Applying **Lemma 2** we know  $\hat{\tau} <: \mathtt{annot}(\mathtt{erase}(\hat{\tau}),\varepsilon)$ . Lastly, by  $\varepsilon\textsc{-Subsume}$  we have  $\Gamma, y: \mathtt{erase}(\hat{\tau}) \vdash \mathtt{annot}(\mathtt{erase}(x),\varepsilon): \mathtt{annot}(\mathtt{erase}(x),\varepsilon)$  with  $\varepsilon \cup \mathtt{effects}(\mathtt{annot}(\Gamma,\varepsilon))$ .

**Subcase 2:**  $x \neq y$ . Then  $x : \tau \in \Gamma$ . Together with the definition  $\mathrm{annot}(x,\varepsilon) = x$ , we know  $x : \mathrm{annot}(\tau,\varepsilon) \in \mathrm{annot}(\Gamma,\varepsilon)$ . By  $\varepsilon$ -VAR we have  $\hat{\Gamma}$ ,  $\mathrm{annot}(\Gamma,\varepsilon)$ ,  $y : \hat{\tau} \vdash \mathrm{annot}(x,\varepsilon) : \mathrm{annot}(\tau,\varepsilon)$  with  $\varnothing$ . Lastly, by  $\varepsilon$ -Subsume we have  $\Gamma,y:\mathrm{erase}(\hat{\tau}) \vdash \mathrm{annot}(\mathrm{erase}(x),\varepsilon) : \mathrm{annot}(\mathrm{erase}(x),\varepsilon)$  with  $\varepsilon \cup \mathrm{effects}(\mathrm{annot}(\Gamma,\varepsilon))$ .

Case: T-RESOURCE Then  $\Gamma, y : \operatorname{erase}(\hat{\tau}) \vdash r : \{r\}$ . By definition,  $\operatorname{annot}(r, \varepsilon) = r$  and  $\operatorname{annot}(\{r\}, \varepsilon)$ . By  $\varepsilon$ -RESOURCE  $\hat{\Gamma}$ ,  $\operatorname{annot}(\Gamma, \varepsilon), y : \hat{\tau} \vdash r : \{r\}$  with  $\varnothing$ . By  $\varepsilon$ -SUBSUME,  $\hat{\Gamma}$ ,  $\operatorname{annot}(\Gamma, \varepsilon), y : \hat{\tau} \vdash r : \{r\}$  with  $\varepsilon \cup \operatorname{effects}(\operatorname{annot}(\Gamma, \varepsilon))$ .

Case: T-ABS Then  $\Gamma, y: \operatorname{erase}(\hat{\tau}) \vdash \lambda x: \tau_1.e_{body}: \tau_1 \to \tau_2$ . By inversion, we get the sub-derivation  $\Gamma, y: \operatorname{erase}(\hat{\tau}), x: \tau_1 \vdash e_2: \tau_2$ . By definition,  $\operatorname{annot}(e, \varepsilon) = \operatorname{annot}(\lambda x: \tau_1.e_2, \varepsilon) = \lambda x: \operatorname{annot}(\tau_1, \varepsilon).\operatorname{annot}(e_2, \varepsilon)$  and  $\operatorname{annot}(\tau_1, \varepsilon) = \operatorname{annot}(\tau_1, \varepsilon) \to_{\varepsilon} \operatorname{annot}(\tau_2, \varepsilon)$ .

To apply the inductive assumption to  $e_2$  we use the unlabelled context  $\Gamma, x:\tau_1$ . The inductive assumption tells us  $\hat{\Gamma}$ , annot $(\Gamma, \varepsilon), y:\hat{\tau}, x:$  annot $(\tau_1, \varepsilon) \vdash$  annot $(e_2, \varepsilon):$  annot $(\tau_2, \varepsilon)$  with  $\varepsilon \cup$  effects(annot $(\Gamma, \varepsilon)$ )  $\cup$  effects(annot $(\tau_1, \varepsilon)$ ). Call this last effect-set  $\varepsilon'$ . By  $\varepsilon$ -ABS, we get  $\hat{\Gamma}$ , annot $(\Gamma, \varepsilon), y:\hat{\tau} \vdash \lambda x:$  annot $(\tau_1, \varepsilon)$ .annot $(e_2, \varepsilon):$  annot $(\hat{\tau}_1) \to_{\varepsilon'}$  annot $(\hat{\tau}_2)$  with  $\varnothing$ . Then by  $\varepsilon$ -Subsume, we get  $\hat{\Gamma}$ , annot $(\Gamma, \varepsilon), y:\hat{\tau} \vdash$  annot $(e, \varepsilon):$  annot $(\hat{\tau}_1) \to_{\varepsilon}$  annot $(\hat{\tau}_2)$  with  $\varepsilon \cup$  effects(annot $(\Gamma), \varepsilon$ ).

Case: T-APP Then  $\Gamma, y: \operatorname{erase}(\hat{\tau}) \vdash e_1 \ e_2 : \tau_3$ , where  $\Gamma, y: \operatorname{erase}(\hat{\tau}) \vdash e_1 : \tau_2 \to \tau_3$  and  $\Gamma, y: \operatorname{erase}(\hat{\tau}) \vdash e_2 : \tau_2$ . By applying the inductive assumption to  $e_1$  and  $e_2$ , we get  $\hat{\Gamma}$ , annot $(\Gamma, \varepsilon), y: \hat{\tau} \vdash \operatorname{annot}(e_1, \varepsilon): \operatorname{annot}(\tau_1, \varepsilon)$  with  $\varepsilon$  and  $\hat{\Gamma}$ , annot $(\Gamma, \varepsilon), y: \hat{\tau} \vdash \operatorname{annot}(e_2, \varepsilon): \operatorname{annot}(\tau_2, \varepsilon)$  with  $\varepsilon$ . Simplifying,  $\hat{\Gamma}$ , annot $(\Gamma, \varepsilon), y: \hat{\tau} \vdash \operatorname{annot}(e_1, \varepsilon): \operatorname{annot}(\tau_2, \varepsilon) \to_{\varepsilon} \operatorname{annot}(\tau_3, \varepsilon)$  with  $\varepsilon$ . Then by  $\varepsilon$ -APP, we get  $\hat{\Gamma}$ , annot $(\Gamma, \varepsilon), y: \hat{\tau} \vdash \operatorname{annot}(e_1 e_2, \varepsilon): \operatorname{annot}(\tau_3, \varepsilon)$  with  $\varepsilon$ .

Case: T-OPERCALL Then  $\Gamma, y: \mathtt{erase}(\hat{\tau}) \vdash e_1.\pi: \mathtt{Unit}$ . By inversion we get the subderivation  $\Gamma, y: \mathtt{erase}(\hat{\tau}) \vdash e_1: \{\bar{r}\}$ . By definition,  $\mathtt{annot}(\{\bar{r}\}, \varepsilon) = \{\bar{r}\}$ . By inductive assumption,  $\hat{\Gamma}, \mathtt{annot}(\Gamma, \varepsilon), y: \hat{\tau} \vdash e_1: \{\bar{r}\} \text{ with } \varepsilon \cup \mathtt{effects}(\mathtt{annot}(\Gamma, \varepsilon))$ . By  $\varepsilon$ -OPERCALL,  $\hat{\Gamma}, \mathtt{annot}(\Gamma, \varepsilon), y: \hat{\tau} \vdash e_1.\pi: \{\bar{r}\} \mathtt{with } \varepsilon \cup \{\bar{r}.\pi\}$ .

It remains to show  $\{\bar{r}.\pi\}\subseteq \varepsilon$ . We shall do this by considering where r must have come from (which subcontext left of the turnstile).

**Subcase 1.**  $r = \hat{\tau}$ . As  $\varepsilon = \text{effects}(\hat{\tau})$ , then  $r.\pi \in \text{effects}(\hat{\tau})$ .

**Subcase 2.**  $r:\{r\}\in\Gamma.$  As annot $(r,\varepsilon)=r$ , then  $r.\pi\in\mathrm{annot}(\Gamma,\varepsilon).$ 

**Subcase 3.**  $r:\{r\}\in \hat{\Gamma}$ . Then because  $\Gamma,y:\operatorname{erase}(\hat{\tau})\vdash e_1:\{\bar{r}\}$ , then  $r\in \Gamma$  or

 $r = \mathtt{erase}(\hat{\tau}) = \hat{\tau}$  and one of the above subcases must also hold.

**Theorem 12** (Preservation). *If*  $\hat{\Gamma} \vdash \hat{e}_A : \hat{\tau}_A$  with  $\varepsilon_A$  and  $e_A \longrightarrow e_B \mid \varepsilon_C$ , then  $\hat{\Gamma} \vdash e_B : \tau_B$  with  $\varepsilon_B$ , where  $e_B <: e_A$  and  $\varepsilon \cup \varepsilon_B \subseteq \varepsilon_A$ .

*Proof.* By induction on  $\hat{\Gamma} \vdash \hat{e}_A : \hat{\tau}_A$  with  $\varepsilon_A$ , and then on  $e_A \longrightarrow e_B \mid \varepsilon$ .

Case:  $\varepsilon$ -VAR,  $\varepsilon$ -RESOURCE,  $\varepsilon$ -UNIT,  $\varepsilon$ -ABS. Then  $e_A$  is a value and cannot be reduced, so the theorem holds vacuously.

Case:  $\varepsilon$ -APP. Then  $e_A = \hat{e}_1 \ \hat{e}_2$  and  $\hat{e}_1 : \hat{\tau}_2 \to_{\varepsilon} \hat{\tau}_3$  with  $\varepsilon_1$  and  $\hat{\Gamma} \vdash \hat{e}_2 : \hat{\tau}_2$  with  $\varepsilon_2$ .

**Subcase:** E-APP1. Todo. **Subcase:** E-APP2. Todo.

**Subcase:** E-APP3. Then  $(\lambda x:\hat{\tau}_2.\hat{e})\hat{v}_2 \longrightarrow [\hat{v}_2/x]\hat{e} \mid \varnothing$ . By inversion on the typing rule for  $\lambda x:\hat{\tau}_2.\hat{e}$  we know  $\Gamma, x:\hat{\tau}_2 \vdash \hat{e}:\hat{\tau}_3$  with  $\varepsilon_3$ . By canonical forms,  $\varepsilon_2 = \varnothing$  because  $\hat{e}_2 = \hat{v}_2$  is a value. Then by the substitution lemma,  $\hat{\Gamma} \vdash [\hat{v}_2/x]\hat{e}:\hat{\tau}_3$  with  $\varepsilon_3$ . By canonical forms,  $\varepsilon_1 = \varepsilon_2 = \varnothing = \varepsilon_C$ . Therefore  $\varepsilon_A = \varepsilon_3 = \varepsilon_B \cup \varepsilon_C$ .

*Case:*  $\varepsilon$ -OPERCALL.

**Subcase:** E-OPERCALL1.

**Subcase:** Otherwise the reduction rule used was E-OPERCALL2 and  $v_1.\pi \longrightarrow \text{unit} \mid \{r.\pi\}$ . By canonical forms,  $\hat{\Gamma} \vdash v_1$ : unit with  $\{r.\pi\}$ . Also,  $\hat{\Gamma} \vdash \text{unit}$ : Unit with  $\varnothing$ . Then  $\tau_B = \tau_A$ . Also,  $\varepsilon_C \cup \varepsilon_B = \{r.\pi\} = \varepsilon_A$ .

Case:  $\varepsilon$ -MODULE Then  $e_A = \text{import}(\varepsilon) \ x = \hat{e} \ \text{in } e$ .

**Subcase:** E-MODULE1 If the reduction rule used was E-MODULECALL1 then the result follows by applying the inductive hypothesis to  $\hat{e}$ .

**Subcase:** E-MODULE2 Otherwise  $\hat{e}$  is a value and the reduction used was E-MODULECALL2. The following are true:

- 1.  $e_A = import(\varepsilon) x = \hat{v} in e$
- 2.  $\hat{\Gamma} \vdash e_A : \mathtt{annot}(\tau, \varepsilon) \ \mathtt{with} \ \varepsilon \cup \varepsilon_1$
- 3.  $import(\varepsilon) x = \hat{v} in e \longrightarrow [\hat{v}/x] annot(e, \varepsilon) \mid \varnothing$
- 4.  $\Gamma \vdash \hat{v} : \hat{\tau} \text{ with } \varnothing$
- 5.  $\varepsilon = \text{effects}(\hat{\tau})$
- 6. ho-safe( $\hat{\tau}, \varepsilon$ )
- 7.  $x : erase(\hat{\tau}) \vdash e : \tau$

Apply the annotation lemma with  $\Gamma = \emptyset$  to get  $\hat{\Gamma}, x : \hat{\tau} \vdash \mathtt{annot}(e, \varepsilon) : \mathtt{annot}(\tau, \varepsilon)$  with  $\varepsilon$ . From **4.** we have  $\hat{\Gamma} \vdash \hat{v} : \hat{\tau}$  with  $\emptyset$ , so we can apply the substitution lemma, giving  $\hat{\Gamma} \vdash [\hat{v}/x]\mathtt{annot}(e, \varepsilon) : \mathtt{annot}(\tau, \varepsilon)$  with  $\varepsilon$ . By canonical forms,  $\varepsilon_1 = \varepsilon_C = \emptyset$ . Then  $\varepsilon_B = \varepsilon = \varepsilon_A \cup \varepsilon_C$ . By examination,  $\tau_A = \tau_B = \mathtt{annot}(\tau, \varepsilon)$ .

**Theorem 13** (Soundness). If  $\hat{\Gamma} \vdash \hat{e}_A : \hat{\tau}_A$  with  $\varepsilon_A$  and  $\hat{e}_A$  is not a value, then  $e_A \longrightarrow e_B \mid \varepsilon$ , where  $\hat{\Gamma} \vdash e_B : \hat{\tau}_B$  with  $\varepsilon_B$  and  $\hat{\tau}_B <: \hat{\tau}_A$  and  $\varepsilon_B \cup \varepsilon \subseteq \varepsilon_A$ .

*Proof.* If  $\hat{e}_A$  is not a value then the reduction exists by the progress theorem. The rest follows by the preservation theorem.

**Theorem 14** (Multi-step Soundness). If  $\hat{\Gamma} \vdash \hat{e}_A : \hat{\tau}_A$  with  $\varepsilon_A$  and  $e_A \longrightarrow^* e_B \mid \varepsilon$ , where  $\hat{\Gamma} \vdash e_B : \hat{\tau}_B$  with  $\varepsilon_B$  and  $\hat{\tau}_B <: \hat{\tau}_A$  and  $\varepsilon_B \cup \varepsilon \subseteq \varepsilon_A$ .

*Proof.* By induction on the length of the multi-step reduction.

Case: Length 0. Then  $e_A = e_B$ , and therefore  $\tau_A = \tau_B$  and  $\varepsilon = \emptyset$  and  $\varepsilon_A = \varepsilon_B$ .

Case: Length 1. Then the result follows by single-step soundness.

Case: Length n+1. Then by inversion the multi-step can be split into a multi-step of length n, which is  $\hat{e}_A \longrightarrow^* \hat{e}_C \mid \varepsilon'$  and a single-step of length 1, which is  $e_C \longrightarrow e_B \mid \varepsilon''$ , where  $\varepsilon = \varepsilon' \cup \varepsilon''$ . By inductive assumption and preservation theorem,  $\hat{\Gamma} \vdash \hat{e}_C : \hat{\tau}_C$  with  $\varepsilon_C$  and  $\hat{\Gamma} \vdash \hat{e}_B : \hat{\tau}_B$  with  $\varepsilon_B$ . By inductive assumption,  $\hat{\tau}_C <: \hat{\tau}_A$  and  $\hat{\varepsilon}_C \cup \varepsilon' \subseteq \varepsilon_A$ . By single-step soundness,  $\hat{\tau}_B <: \hat{\tau}_C$  and  $\hat{\varepsilon}_B \cup \varepsilon'' \subseteq \varepsilon_C$ . Then by transitivity,  $\hat{\tau}_B <: \hat{\tau}$  and  $\hat{\varepsilon}_B \cup \varepsilon' \cup \varepsilon'' = \varepsilon_B \cup \varepsilon \subseteq \varepsilon_A$ .

# **Bibliography**

- [1] AHO, A. V., SETHI, R., AND ULLMAN, J. D. Compilers: Principles, Techniques, and Tools. Addison-Wesley, Reading, MA, USA, 1986.
- [2] Bracha, G., von der Ahé, P., Bykov, V., Kashai, Y., Maddox, W., and Miranda, E. Modules as Objects in Newspeak. In *European Conference on Object-Oriented Programming* (2010).
- [3] CHURCH, A. A formulation of the simple theory of types. *American Journal of Mathematics* 5 (1940), 56–68.
- [4] DENNIS, J. B., AND VAN HORN, E. C. Programming Semantics for Multiprogrammed Computations. *Communications of the ACM 9*, 3 (1966), 143–155.
- [5] KURILOVA, D., POTANIN, A., AND ALDRICH, J. Modules in wyvern: Advanced control over security and privacy. In *Symposium and Bootcamp on the Science of Security* (2016). Poster.
- [6] LISKOV, B. Keynote address data abstraction and hierarchy. In *Addendum to the Proceedings on Object-oriented Programming Systems, Languages and Applications (Addendum)* (New York, NY, USA, 1987), OOPSLA '87, ACM, pp. 17–34.
- [7] LIU, F. A study of capability-based effect systems. Master's thesis, École Polytechnique Fédérale de Lausanne, 2016.
- [8] MAFFEIS, S., MITCHELL, J. C., AND TALY, A. Object Capabilities and Isolation of Untrusted Web Applications. In *IEEE Symposium on Security and Privacy* (2010).
- [9] MILLER, M., YEE, K.-P., AND SHAPIRO, J. Capability myths demolished. Tech. rep., 2003.
- [10] MILLER, M. S. Robust Composition: Towards a Unified Approach to Access Control and Concurrency Control. PhD thesis, Johns Hopkins University, 2006.
- [11] NIELSON, F., AND NELSON, H. R. Type and Effect Systems. pp. 114–136.
- [12] ODERSKY, M., ALTHERR, P., CREMET, V., DUBOCHET, G., EMIR, B., HALLER, P., MICHELOUD, S., MIHAYLOV, N., MOORS, A., RYTZ, L., SCHINZ, M., STENMAN,

42 BIBLIOGRAPHY

- E., AND ZENGER, M. Scala Language Specification. http://scala-lang.org/files/archive/spec/2.11/. Last accessed: Nov 2016.
- [13] PIERCE, B. C. *Types and Programming Languages*. The MIT Press, Cambridge, MA, USA, 2002.
- [14] SALTZER, J. H. Protection and the Control of Information Sharing in Multics. *Communications of the ACM 17*, 7 (1974), 388–402.