## **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	on and the second of the secon	1						
2	Bacl	ackground							
	2.1	-							
		2.1.1	Grammar	6					
		2.1.2	Dynamic Rules	6					
		2.1.3	Static Rules	9					
		2.1.4	Soundness	11					
	2.2	$\lambda^{\rightarrow}$ : S	imply-Typed $\lambda$ -Calculus	12					
	2.3	Effect Systems							
		2.3.1	ETL: Effect-Typed Language	16					
	2.4	The C	apability Model	16					
	2.5	First-C	Class Modules	17					
3	Effe	ct Calc	uli	19					
	3.1	Exam	ples	19					
	3.2	$\lambda_{\pi}^{\rightarrow}$ : O	peration Calculus	19					
		3.2.1	$\lambda_{\pi}^{\rightarrow}$ Grammar	20					
		3.2.2	$\lambda_\pi^{\rightarrow}$ Dynamic Rules	21					
		3.2.3	$\lambda_{\pi}^{\rightarrow}$ Static Rules	22					
		3.2.4	Soundness of $\lambda_{\pi}^{\rightarrow}$	23					
	3.3	$\lambda_{\pi,\varepsilon}^{\rightarrow}$ : E	psilon Calculus	27					
		3.3.1	$\lambda_{\pi,\varepsilon}^{ o}$ Grammar	27					
		3.3.2	$\lambda_{\pi,\varepsilon}^{\rightarrow}$ Dynamic Rules	29					
		3.3.3	$\lambda_{\pi,\varepsilon}^{\rightarrow}$ Static Rules	30					
		3.3.4	Soundness of $\lambda_{\pi,\varepsilon}^{\rightarrow}$	34					
4	App	applications 3							
	4.1	Transl	ations and Encodings	39					
		4.1.1	Unit	40					
		4.1.2	Let	40					
		4.1.3	Modules and Objects	41					

vi *CONTENTS* 

В	$\lambda_{\pi,arepsilon}^{ o}$ :	Proofs		61				
A	$\lambda_\pi^{ o}$ I	Proofs		57				
	5.3	Conclu	asion	56				
	5.2	Future	Work	55				
	5.1	Relate	d Work	55				
5	Eval	uation		55				
		4.2.8	Resource Leak	47				
		4.2.7	Hidden Authority 2	47				
		4.2.6	Hidden Authority 2	47				
		4.2.5	Hidden Authority	47				
		4.2.4	API Violation	47				
		4.2.3	API Violation	47				
		4.2.2	API Violation	46				
		4.2.1	Unannotated Client	44				
	4.2	Examp	Examples					

## Chapter 1

### Introduction

Good software is distinguished from bad software by design qualities such as security, maintainability, and performance. We are interested in how the design of a programming language and its type system can help achieve these qualities.

It is difficult to determine if a piece of code is trustworthy. Several scenarios can give rise to the issue of trust, such as in a development environment adhering to *code own-ership*. In this setting, groups of developers may function as experts over certain components. When their components must interact with code sourced from outside their domain of expertise, they can make false assumptions or violate the internal constraints of other components by using them incorrectly. Another setting involves applications which 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.

A range of methods might be employed to help us determine whether we can trust a piece of code. Sandboxing is one, where the suspect code is executed in a safe, separate environment from the trusted code, but this approach has many shortcomings [?]. Verification techniques allow for robust analyses on code, but are heavyweight and require the developers to have a deep understanding of the techniques being employed [?]. Lightweight analyses, such as those present in type system, enjoy the benefit of being easy for the developer to use, but existing languages do not provide adequate controls for detecting and isolating untrustworthy components [?].

A qualitative approach to trustworthiness is to develop according to particular guidelines considered to be in good practice. One such guideline is the *principle of least authority*: that software components should only have access to the information and resources necessary for their purpose [15]. 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 [?]. This report is interested in how type systems might enforce more static constraints, putting developers in a more-informed position to make qualitative judgements about the trustworthiness of code.

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 design. Miller has identified the ways in which capabilities should proliferate to encourage *robust composition* — a set of ideas summarised as "only connectivity begets connectivity" [11]. In his paradigm, the reference graph of a program is the same as the access graph. This eliminates *ambient authority*, whereby a privilege is exercised without being explicitly declared. This enables one to reason about what privileges a component might exercise by examining its interface. Building on these ideas, Maffeis et. al. formalised *capability-safety* of a language, showing a subset of Caja (a JavaScript implementation) meets this notion [9].

In addition to realising privilege separation, capabilities also encapsulate the source of *effects*. An *effect* describes some intensional information about the way in which a program executes [12]. For example, a logger's method might append to a file, and so executing this method would incur {File.append}. To be able to do this in a capability-safe language, the logger must have a capability for the File. Therefore, the constraints imposed by how capabilities proliferate, also impose constraints on what effects the components of a program may incur.

Capabilities can offer fine-grained control over the way in which privileges are exercised via *confinement*. For example, while we expect the logger's log method to have the append effect, a sloppy or malicious implementation may incur extra effects on the File, such as write or close. Knowing the logger might incur these extra effects may inform the decision a developer makes about whether or not to trust this particular component.

An effect-system is an extension to a type-system, which track what effects a program might incur, and where. They have been used to do this, this, and that. Some have criticised their verbosity lightweight polymorphic effects paper as a point against their practical adoption. An effect system such as the Talpin-Jouvelot ETL system ATAPL? requires the annotation of all values in a program. This requires the developer to be aware, at all points, of what effects are in scope. Is this true of all, or even most, effect systems? Minor alterations to the signatures and effects of one component might 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 development, affecting their usability in large 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.

Because capabilities encapsulate effectful privileges, and because capability-safe languages impose constraints on how these privileges can spread throughout the system, this considerably simplifies effect-based reasoning. To incur an effect requires one to possess a capability for the appropriate resource, and whether this resource is captured by a component can be determined by inspecting the type-signatures of the component. The developer need not look at the source code. This is the key contribution of this report: that capability-safety facilitates a low-cost effect-based reasoning with minimal user overhead.

We begin this paper by discussing preliminary concepts involving the formal definition of programming languages (2.1.), effect systems (2.2.), and the capability model (2.3.). Along the way we summarise some existing languages to illustrate these points.

Chapter 3 introduces a pair of languages called  $\lambda_{\pi}^{\rightarrow}$  and  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ .  $\lambda_{\pi}^{\rightarrow}$  is a typed lambda calculus with a simple notion of capabilities and runtime effects. Every function in  $\lambda_{\pi}^{\rightarrow}$  is annotated, which gives a simple, sound system for determining what effects a piece of code might incur.  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  allows for unannotated code by introducing an import construct. At the point of interaction between labelled and unlabelled code, a capability-based reasoning enables us to make a safe inference about what effects the unlabelled code might incur.

Chapter 4 shows how  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  might be used practically, and we try to convince the reader that  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  can be implemented in existing capability-safe languages in a routine manner. We finish with a literature comparison.

## Chapter 2

## **Background**

In this section we cover some of the necessary concepts and existing work informing this report. First we cover the process of formally defining a programming language and proving some of its properties. For this purpose, we present EBL. We then summarise the rules and properties of a variation of the simply-typed lambda calculus  $\lambda^{\rightarrow}$ .  $\lambda^{\rightarrow}$  is an historically important model of computation and serves as a basis for many functional programming language. It is also the basis of  $\lambda^{\rightarrow}_{\pi,\varepsilon}$ , so a preliminary understanding of  $\lambda^{\rightarrow}$  will help us understand  $\lambda^{\rightarrow}_{\pi,\varepsilon}$ .

 $\lambda_{\pi,\varepsilon}^{\rightarrow}$  is a capability-based language with an effect system. To understand what that means we cover some existing work on effect systems and discuss Miller's capability model.

### 2.1 Formally Defining a Programming Language

A programming language can be defined formally by supplying three sets of rules: a grammar, which defines syntactically legal terms; static rules, which determine whether programs meet certain well-formedness properties; and dynamic rules, which express the meaning of a program by defining how they are executed. When a language has been defined we want to know its rules are mathematically correct.

We illustrate these concepts by presenting EBL, which is a simple, type-safe language based around boolean and arithmetic expressions. Like every language in this report, it is expression-based, meaning that valid programs can always be evaluated to yield a value. Although EBL is not very interesting, the process of defining it and proving its rules correct illustrate the general approach this report will take towards formally specifying programming languages.

#### 2.1.1 Grammar

The grammar of a language specifies what strings are syntactically legal. It 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. A syntactically valid string is called a term.

A string like 3 + (x + 2) should be seen as a short-hand for the corresponding abstract syntax tree (AST), whose structure is given by the rules of the grammar. A diagram might be nice here. Sometimes the AST is ambiguous, as in 3 + x + 2 which might be parsed as 3 + (x + 2) or as (3 + x) + 2. How we parse and disambiguate is an implementation detail, so throughout this report we only consider strings which unambiguously correspond to a valid AST.

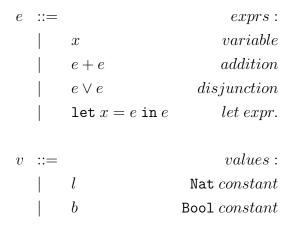


Figure 2.1: Grammar for EBL expressions.

#### 2.1.2 Dynamic Rules

The dynamic rules of a language specify the meaning of syntactically-valid terms. There are different approaches, but the one we use is called *small-step semantics*, where the meaning of a program is given by how it is executed. This is specified as a set of *in-ference 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. If an inference rule has no premises it is called an *axiom*. An instantiation of a particular inference rule is called a judgement.

$$e \longrightarrow e$$

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

$$\frac{e_1 \longrightarrow e_1'}{e_1 \vee e_2 \longrightarrow e_1' \vee e_2} \text{ (E-OR1) } \frac{}{\text{true} \vee e_2 \longrightarrow \text{true}} \text{ (E-OR2) } \frac{}{\text{false} \vee e_2 \longrightarrow e_2} \text{ (E-OR3)}$$

$$\frac{e_1 \longrightarrow e_1'}{\operatorname{let} x = e_1 \text{ in } e_2 \longrightarrow \operatorname{let} x = e_1' \text{ in } e_2} \text{ (E-LET1)} \quad \frac{1}{\operatorname{let} x = v \text{ in } e_2 \longrightarrow [v/x]e_2} \text{ (E-LET2)}$$

Figure 2.2: Inference rules for single-step reductions.

Figure 2.4. gives the dynamic rules for EBL. Their conclusions specify members of a binary relation  $\longrightarrow$ , representing a single computational step. When the relation holds of a particular pair, we say the judgement  $e \longrightarrow e'$  holds, and that e reduces to e'.

If a non-variable expression is irreducible under the dynamic rules, it is called a value. The grammar of EBL also specifies a category of terms called "value". As we shall see, these two definitions correspond. 3, true, and false are examples of irreducible expressions.

A disjunction is reduced by first reducing the left-hand side to a value (E-OR1). If the left-hand side is the boolean literal true, then we can reduce the expression to true (because true  $\lor Q = \text{true}$ ). Otherwise if the left-hand side is the boolean literal false, we can reduce the expression to the right-hand side  $e_2$  (because false  $\lor Q = Q$ ). This particular formulation of the rules encodes short-circuiting behaviour into  $\lor$ , meaning that if the left-hand side is true, the expression evaluates to true without checking the right-hand side.

An addition expression is reduced by first reducing the left-hand side to a value (E-ADD1) and then the right-hand side (E-ADD2) to a value. When both sides are integer literals, the expression reduces to whatever is the sum of those literals.

A let expression is reduced by first reducing the subexpression being bound (E-LET1). When that is a value, we substitute the variable x for the value  $v_1$  in the body  $e_2$  of the let expression. The notation for this is  $[v_1/x]e_2$ . For example, let x = 1 in x + 1 reduces to 1 + 1 by an application of E-LET2.

Consider let x = 1 + 1 in x + 1. According to the rules, 1 + 1 would first be reduced to 2 before the substitution is made into x + 1. This strategy of reducing expressions before they are bound to variable names is *call-by-value*.

Formally, substitution is a function operating on expressions. A definition is given in

Figure 2.5. The notation  $[e_1/x]e$  is short-hand for substitution  $(e, e_1, x)$ . When performing multiple substitutions we use the notation  $[e_1/x_1, e_2/x_2]e$  as shorthand for  $[e_2/x_2]([e_1/x_1]e)$ . Note how the order of the variables has been flipped; the substitutions occur as they are written, left-to-right.

```
substitution :: e × e × v \rightarrow e  [e'/y]l = l   [e'/y]b = b   [e'/y]x = v \text{, if } x = y   [e'/y]x = x \text{, if } x \neq y   [e'/y](e_1 + e_2) = [e'/y]e_1 + [e'/y]e_2   [e'/y](e_1 \vee e_2) = [e'/y]e_1 \vee [e'/y]e_2   [e'/y](\text{let } x = e_1 \text{ in } e_2) = \text{let } x = [e'/y]e_1 \text{ in } [e'/y]e_2 \text{, if } y \neq x \text{ and } y \text{ does not occur }  free in e_1 or e_2
```

Figure 2.3: Substitution for EBL.

A robust definition of the substitution function is surprisingly tricky. Consider the program let x = 1 in (let x = 2 in x + z). It contains two different variables with the same name x, with the inner one "shadowing" the outer one. Neither variable occurs "free", because both have been introduced in the body of the program (one for each let). Such variables are called bound variables. By contrast, z is a free variable because it has no definition in the program. A robust substitution should not accidentally conflate two different variables with identical names, and it should not do anything to bound variables.

To illustrate the solution, consider let x=1 in (let x=2 in x+z). In some sense, this is an equivalent program to let x=1 in (let y=2 in y+z). Because the names of variables are arbitrary, changing them will not change the semantics of the program. Therefore, we freely and implicitly interchange expressions which are equivalent up to the naming of bound variables. This process is called  $\alpha$ -conversion [14, p. 71]. Consequently, we assume variables are (re-)named in this way to avoid these problems and to play nicely with the definition of substitution.

Given a single-step reduction relation, we may define a multi-step reduction relation as a sequence of zero<sup>1</sup> or more single-steps. This is written  $e \longrightarrow^* e'$ . For example, if  $e_1 \longrightarrow e_2$  and  $e_2 \longrightarrow e_3$ , then  $e_1 \longrightarrow^* e_3$ . Figure 2.4. shows how multi-step reduction can be defined with a set of inference rules.

<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 [14, p. 39].

$$\frac{e \longrightarrow^* e}{e \longrightarrow^* e} \text{ (E-MultiStep1)} \quad \frac{e \longrightarrow e'}{e \longrightarrow^* e'} \text{ (E-MultiStep2)}$$

$$\frac{e \longrightarrow^* e' \quad e' \longrightarrow^* e''}{e \longrightarrow^* e''} \text{ (E-MultiStep3)}$$

Figure 2.4: Dynamic rules.

#### 2.1.3 Static Rules

If you try to reduce some terms you either end up with nonsense or get stuck in a situation where no rule applies. For example,  $(1+1) + false \longrightarrow 2 + false$  by E-ADD1, but then you are stuck. false  $\vee$  3  $\longrightarrow$  3 by E-OR3, which is strange.

When designing a language we often want to consider those syntactically legal terms satisfying certain *well-behavedness* properties. One such property is that of being *well-typed*: if a program is well-typed then during execution it will never get *stuck* due to type-errors. Another useful well-formedness property says that every variable used in a program must be declared beforehand. The examples above are *not* well-typed, because they are applying operators to arguments of the wrong type. We want rules to help us determine if this is the case without having to execute the program.

The static rules of EBL describe a basic type system which let us determine, without executing a program, whether it contains type errors. The relevant constructs for EBL are given as a grammar in Figure 2.5. There are two types: Nat and Bool. Furthermore, there is the notion of a *typing context*, which maps variables to their types. This is needed in the case of a program like let x = 1 in x + 1. In trying to determine whether x + 1 is well-typed, we need to know what is the type of x. To do this, when we see the binding x = 1 we extend the context to say that x has type Nat.

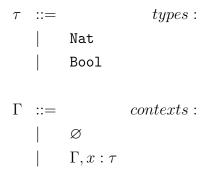


Figure 2.5: Grammar for arithmetic expressions.

Figure 2.6. summarises the static rules of EBL. Note that every judgement holds in a particular typing context. For example, the judgement  $x : Int \vdash x + 1 : Int$  is a claim about a particular property (x + 1) has the type Int in a particular context (the

context where x is an Int). If a judgement can be derived from the empty context, we conventionally write it as  $\vdash e : \tau$  instead of  $\varnothing \vdash e : \tau$ .

T-BOOL and T-NAT are rules which say that constants always type to Bool or Nat. T-VAR says that a variable types to whatever the context binds it to. T-OR types a disjunction if the arguments are both Bool. T-ADD types a sum if the arguments are both Nat. The most interesting rule is T-LET, where the context gains a binding for x when type-checking the body of the let expression. The type of a let expression is the type of its body.

$$\begin{array}{c} \hline \Gamma \vdash e : \tau \\ \hline \\ \hline \hline \Gamma, x : \mathtt{Int} \vdash x : \mathtt{Int} \end{array} (\mathtt{T-VAR}) & \overline{\vdash b : \mathtt{Bool}} \end{array} (\mathtt{T-BOOL}) & \overline{\vdash l : \mathtt{Nat}} \end{array} (\mathtt{T-NAT}) \\ \hline \\ \frac{\Gamma \vdash e_1 : \mathtt{Bool}}{\Gamma \vdash e_1 : \mathtt{Bool}} & \Gamma \vdash e_2 : \mathtt{Bool} \end{array} (\mathtt{T-OR}) & \frac{\Gamma \vdash e_1 : \mathtt{Nat}}{\Gamma \vdash e_1 : e_2 : \mathtt{Nat}} \end{array} (\mathtt{T-ADD}) \\ \hline \\ \frac{\Gamma \vdash e_1 : \tau_1}{\Gamma \vdash \mathtt{let}} & \frac{\Gamma, x : \tau_1 \vdash e_2 : \tau_2}{\Gamma \vdash \mathtt{let}} & (\mathtt{T-LET}) \end{array}$$

Figure 2.6: Inference rules for typing arithmetic expressions.

There are some pesky technicalities about typing contexts which need to be addressed. Although we have defined  $\Gamma$  as a *sequence* of variable-type mappings, the order shouldn't really be significant: x: Int, y: Int is really the same thing as y: Int, x: Int. We essentially want to treat it as a set.

Formally, we can specify their equivalence by giving structural rules which say that a judgement holds in  $\Gamma$  if it holds in any permutation of  $\Gamma$ . Another convention is that any judgement which holds in a context  $\Gamma$  should hold in any bigger context  $\Gamma'$ , where  $\Gamma \subseteq \Gamma'$ . For example,  $x: \operatorname{Int} \vdash x: \operatorname{Int}$ , but it is also true that  $x: \operatorname{Int}, y: \operatorname{Int} \vdash x: \operatorname{Int}$ . In practice, the notation for contexts and the rules for how to manipulate them are so conventional that, beyond the quick summary in Figure 2.3., we will not bother to mention them again.

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

$$\frac{\Gamma \vdash e : \tau \quad \Gamma' \ is \ a \ permutation \ of \ \Gamma}{\Gamma' \vdash e : \tau} \ (\Gamma \text{-Permute}) \ \frac{\Gamma \vdash e : \tau \quad x \notin \Gamma}{\Gamma, x : \tau' \vdash e : \tau} \ (\Gamma \text{-WIDEN})$$

Figure 2.7: Structural rules for typing contexts.

Though EBL has no subtyping, most interesting languages do. This judgement is written form  $\tau_1 <: \tau_2$  and it means that expressions of  $\tau_1$  may be provided anywhere expressions of  $\tau_2$  are expected, and the program will still be well-typed. 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 program's semantic properties [7]. Subtyping rules are not usually totally semantic-preserving, but we'll occasionally use this idea to justify why certain subtyping rules are sensible.

#### 2.1.4 Soundness

Having defined a type system, we want to know it is *sound*: that if the type system says a program is well-typed, the program will not run into type errors during execution. Soundness is a guarantee that our typing judgements, as we intuitively understand them, are mathematically correct. The exact definition of soundness depends on the language under consideration, but is often split into two parts called progress and preservation.

**Theorem 1** (Progress). *If*  $\vdash e : \tau$  *and* e *is not a value, then*  $e \longrightarrow e'$ .

Progress states that any well-typed, non-value term can be reduced i.e. it will not get stuck due to type errors. It also says that the definition of a value, as a non-variable irreducible expression, is coincident with the definition of a value as a particular category of terms in the grammar.

**Theorem 2** (Preservation). *If*  $\vdash e : \tau$  *and*  $e \longrightarrow e'$  *then*  $\vdash e' : \tau$ .

Preservation states that a well-typed term is still well-typed after it has been reduced. This means a sequence of reductions will produce intermediate terms that are also well-typed and do not get stuck. Note that in this particular formulation of preservation for EBL, the type of the term after reduction is the same as the type of the term before reduction.

By combining progress and preservation, we know that a runtime type-error can never occur as the result of a single-step reduction. This is *small-step soundness*. Once this has been established, we may extend this to multi-step reductions by inducting on the length of the multi-step and appealing to the soundness of single-step reductions. This yields the following result.

**Theorem 3** (Soundness). *If*  $\vdash e : \tau$  *and*  $e \longrightarrow e'$  *then*  $\vdash e' : \tau$ .

These theorems are proven by structural induction on the typing rule used  $\Gamma \vdash e : \tau$  or on the reduction rule used  $e \rightarrow e'$ .

In order to prove certain cases of progress and preservation there are two common lemmas needed. The first is canonical forms, which outlines a set of observations that follow immediately by observing the typing rules. The second is the substitution lemma, which says if a term is well-typed in a context  $\Gamma, x: \tau' \vdash e: \tau$ , and you replace variable x with an expression e' of type  $\tau$ , then  $\Gamma \vdash [e'/x]e: \tau$ . In EBL, this lemma is needed to show that the reduction step in E-LET2 preserves soundness. The other languages considered in this report will have similar reduction steps.

A precise formulation of these two lemmas for EBL is given below.

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

```
If Γ ⊢ v : Int, then v = l is a Nat constant.
If Γ ⊢ v : Bool, then b = l is a Bool constant.
```

```
Lemma 2 (Substitution). If \Gamma, x : \tau' \vdash e : \tau and \Gamma \vdash e' : \tau' then \Gamma \vdash [e'/x]e : \tau.
```

Proofs for these lemmas and theorems can be found in Appendix A.

To summarise, soundness is a property which says, generally, that if a type system says a program is well-typed, it will not encounter a runtime type-error. The corollary of this is also interesting to consider: if a program has no runtime type-error, will the type system accept it? This property is called *completeness*, and almost no (interesting) type-systems are complete. This means a type system may reject type-safe programs. However, soundness guarantees that a type system will *always* reject programs which are *not* type-safe. Consider Figure 2.7., which demonstrates a type-safe Java program rejected by Java's type-system: the body of double is type-safe, because the conditional will always execute returnx + x. However, Java will reject this program.

```
public int double(int x) {
   if (true) return x + x;
   else return true;
}
```

Figure 2.8: A type-safe Java method which does not typecheck.

Throughout this report we will only be concerned with sound type systems, but it is important to recognise that these type systems are all *conservative* because they may reject type-safe programs. One view of type-systems is that they "calculate a kind of static approximation to the run-time behaviours of the terms in a program" [14, p. 2]. In order to approximate, simplifying assumptions must be made, and these simplifying assumptions are what make the type-system sound; but assumptions which are too generalising can make the system too conservative and of less practical use. This is an important trade-off we discuss in motivating  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ .

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

The simply-typed  $\lambda$ -calculus  $\lambda^{\to}$  is a model of computation, first described by Alonzo Church [3], based on the definition and application of functions. In this section we present a variation of  $\lambda^{\to}$  with subtyping and summarise its basic properties. Various  $\lambda$ -calculi serve as the basis for numerous functional programming languages, including  $\lambda^{\to}_{\pi,\varepsilon}$ . This section gives us an opportunity to familiarise ourself with  $\lambda^{\to}$  to help introduce  $\lambda^{\to}_{\pi,\varepsilon}$ .

Figure 2.9: Grammar for  $\lambda^{\rightarrow}$ .

Types in  $\lambda^{\rightarrow}$  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). "Arrow-type" and "function-type" will be used interchangeably.

In addition to variables, there are function definitions ("abstraction") and the application of a function to an expression ("application"). For example,  $\lambda x: \mathrm{Int.} x$  is the identity function on integers.  $(\lambda x: \mathrm{Int.} x)$ 3 is the application of the identity function to the integer literal 3.  $(\lambda x: \mathrm{Int.} x)$ true is the application of the identity function to a boolean literal, which is syntactically valid, but as we'll see is not well-typed. A more drastic example is true 3, which is trying to apply true to 3. Again, this is a syntactically valid term, but not well-typed because true is not a function.

$$\begin{array}{c} \Gamma \vdash e : \tau \\ \hline \\ \overline{\Gamma, x : \tau \vdash x : \tau} \end{array} \text{ (T-VAR)} \quad \frac{\Gamma, x : \tau_1 \vdash e : \tau_2}{\Gamma \vdash \lambda x : \tau_1.e : \tau_1 \to \tau_2} \text{ (T-Abs)} \\ \hline \\ \frac{\Gamma \vdash e_1 : \tau_2 \to \tau_3 \quad \Gamma \vdash e_2 : \tau_2}{\Gamma \vdash e_1 \; e_2 : \tau_3} \text{ (T-App)} \quad \frac{\Gamma \vdash e : \tau_1 \quad \tau_1 <: \tau_2}{\Gamma \vdash e : \tau_2} \text{ T-Subsume} \\ \hline \\ \frac{\tau_1' <: \tau_1 \quad \tau_2 <: \tau_2'}{\tau_1 \to \tau_2 <: \tau_1' \to \tau_2'} \text{ (S-Arrow)} \\ \hline \end{array}$$

Static rules for  $\lambda^{\rightarrow}$  are summarised in Figure 2.8. 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-

Figure 2.10: Static rules for  $\lambda^{\rightarrow}$ .

APP states that an application is well-typed if the left-hand expression is a function (has an arrow-type  $\tau_2 \to \tau_3$ ) 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 <: Real, a term of type Int can also be typed as Real. This allows programs such as  $(\lambda x : \text{Real}.x)$  3 to type, as shown in Figure 2.9.

```
\frac{\frac{}{x: \mathtt{Real} \vdash x: \mathtt{Real}} \ (\mathtt{T-VAR})}{\vdash \lambda x: \mathtt{Real} \cdot x: \mathtt{Real} \rightarrow \mathtt{Real}} \ (\mathtt{T-ABS}) \qquad \frac{\vdash 3: \mathtt{Int} \quad \mathtt{Int} <: \mathtt{Real}}{\vdash 3: \mathtt{Real}} \ (\mathtt{T-SUBSUME})}{\vdash (\lambda x: \mathtt{Real} \cdot x) \ 3: \mathtt{Real}} \ (\mathtt{T-APP})
```

Figure 2.11: Derivation tree showing how T-SUBSUME can be used.

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.

This presentation has no subtyping rules without premises (axioms), which means there is no way to actually prove a particular subtyping judgement. In practice, we add subtyping axioms for the base-types we have chosen as primitive in our calculus. For example, given base types Int and Real, we might add Real <: Int as a rule. This is largely an implementation detail particular to your chosen set of base-types, so we give no subtyping axioms here (but will later when describing  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ ).

```
substitution :: e \times e \times v \rightarrow e
```

```
[v/y]x = v, if x = y

[v/y]x = x, if x \neq y

[v/y](\lambda x : \tau \cdot e) = \lambda x : \tau \cdot [v/y]e, if y \neq x and y does not occur free in e

[v/y](e_1 e_2) = ([v/y]e_1)([v/y]e_2)
```

Figure 2.12: Substitution for  $\lambda^{\rightarrow}$ .

Substitution in  $\lambda^{\rightarrow}$  follows the same conventions as it does in EBL. Substitution on an application is the same as substitution on its sub-expressions. Substitution on a function involves substitution on the function body.

Applications are the only reducible expressions in  $\lambda^{\rightarrow}$ . Such an expression is reduced by first reducing the left subexpression (E-APP1). For a well-typed expression, this will always be a function. Once that is a value, the right subexpression is reduced (E-APP2). When both subexpressions are values, the right subexpression replaces the formal argument of the function via substitution. The multi-step rules for  $\lambda^{\rightarrow}$  are identical to those in EBL.

15

$$e \longrightarrow e$$

$$\frac{e_{1} \longrightarrow e'_{1} \mid \varepsilon}{e_{1}e_{2} \longrightarrow e'_{1}e_{2} \mid \varepsilon} \text{ (E-APP1)} \quad \frac{e_{2} \longrightarrow e'_{2} \mid \varepsilon}{v_{1}e_{2} \longrightarrow v_{1}e'_{2} \mid \varepsilon} \text{ (E-APP2)}$$

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

Figure 2.13: Dynamic rules for  $\lambda^{\rightarrow}$ .

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

**Theorem 4** (
$$\lambda^{\rightarrow}$$
 Soundness). *If*  $\Gamma \vdash e_A : \tau_A$  *and*  $e_A \longrightarrow^* e_B$ , *then*  $\Gamma \vdash e_B : \tau_B$ , *where*  $\tau_B <: \tau_A$ .

Note how with the inclusion of subtyping rules, the type after reduction can get more specific than the type before reduction, but never less specific. This is in contrast to EBL, where the type remains the same.

 $\lambda^{\rightarrow}$  is also strongly-normalizing, meaning that well-typed terms always halt i.e. they eventually yield a value. As a consequence it is *not* Turing complete, meaning there are certain computer programs which cannot be written in  $\lambda^{\rightarrow}$ . By comparison, the *untyped*  $\lambda$ -calculus is known to be Turing complete [5]. The essential ingredient missing from  $\lambda^{\rightarrow}$  is a means of general recursion. In mainstream languages such as Java, this is realised by constructs like the while loop; in the untyped  $\lambda$ -calculus by the Y-combinator.  $\lambda^{\rightarrow}$  can be made Turing-complete by adding a fix operator which mimics the Y-combinator.

Turing-completeness is an essential property for practical, general-purpose programming languages. However, the key contribution of this report is in the static rules of  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ , and not the expressive power of its dynamic semantics. Therefore we acknowledge this practical short-coming, but leave the basis of  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  as a Turing-incomplete language to reduce the number of rules and simplify its presentation.

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

#### 2.3 Effect Systems

In the previous section we looked at how the static rules of a language might make a judgement such as  $\Gamma \vdash e : \tau$ , which ascribes the type  $\tau$  to program e. This expresses a certain about what the runtime behaviour of the program is: namely that successive reductions of e will produce terms of type  $\tau$ , and that this sequence of reductions will never get stuck due to a runtime type-error.

One extension to classical type systems is to incorporate a theory of *effects*. A *type-and-effect* system can ascribe a type and an effect to a piece of code, the effect component of which specifies intensional information about what will happen during the execution of the program [12]. For example, a judgement like  $\Gamma \vdash e : \tau$  with {File.write} means

that successive reductions of e will result in terms of type  $\tau$ , and during execution might write to a file. This tells us extra information about what can happen during runtime.

Talk about use of effect systems. Mention the basic effect system we will introduce.

#### 2.3.1 ETL: Effect-Typed Language

#### 2.4 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 [11]:

- 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".

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 [10]. 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.

```
import java.io.File;
  import java.io.IOException;
  import java.util.ArrayList;
  class MyList<T> extends ArrayList<T> {
5
     @Override
    public boolean add(T elem) {
       try {
         File file = new File("$HOME/.bashrc");
         file.createNewFile();
10
       } catch (IOException e) {}
11
       return super.add(elem);
12
13
14
  import java.util.List;
  class Main {
3
    public static void main(String[] args) {
4
       List<String> list = new MyList<String>();
       list.add(''doIt'');
     }
7
```

Figure 2.14: Main exercises ambient authority over a File capability.

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

#### 2.5 First-Class Modules

8 }

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 [13], 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 [6].

## Chapter 3

### **Effect Calculi**

In this section we give examples to motivate the practical benefits of an effect system. We then describe a pair of languages:  $\lambda_{\pi}^{\rightarrow}$  and  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ . In  $\lambda_{\pi}^{\rightarrow}$ , every function's input type is labelled by the effects that can be incurred by values of that type. This enables reasoning about the effects that might be incurred by a piece of code.

We then explore what happens when we drop the requirement that every part of a program be labelled with its effects. This leads to the description of  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ , where labelled and unlabelled code can interact via an import construct. import enables a programmer to nest unlabelled code inside labelled code. The primary result of this chapter is that capability-safety enables a simple, effect-safe inference about the unlabelled code in an import expression.

#### 3.1 Examples

Be motivational here, explaining the need for effect-based reasoning using examples from Chapter 4

### 3.2 $\lambda_{\pi}^{\rightarrow}$ : Operation Calculus

The operation calculus  $\lambda_{\pi}^{\rightarrow}$  is an extension of  $\lambda^{\rightarrow}$  with primitive capabilities (resources), on which operations may be invoked. An effect is an operation invoked on a resource. Every function-type is annotated by what effects may be incurred during execution of the function body. The static rules of  $\lambda_{\pi}^{\rightarrow}$  can inspect this information and ascribe a set of effects to a piece of code, which conservatively approximates what will happen at runtime.

The results of this chapter are straightforward and unsurprising, but  $\lambda_{\pi}^{\rightarrow}$  contains new notations and a new concept of effect-soundness, and forms the basis of the more interesting  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ , which we discuss in the next chapter.

#### 3.2.1 $\lambda_{\pi}^{\rightarrow}$ Grammar

The grammar for  $\lambda_{\pi}^{\rightarrow}$  and its meta-theory are summarised in Figure 3.1. Expressions are the same as they are in  $\lambda^{\rightarrow}$ , except for two new forms: resource literals and operation-calls.

A resource literal r is a variable drawn from a fixed set R. They cannot be created or destroyed at runtime. The resources in R model those initial capabilities passed into the program, perhaps from the system environment. For example, a File or a Socket would be an example of a resource literal.

An operation is a special action that be 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; like resources, they cannot be created or destroyed 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.open instead of (File, open). A set of effects is denoted by  $\varepsilon$ . Effects and operations look notationally similar, but should be distinguished: an effect is some description of runtime behaviour in the meta-theory of  $\lambda_{\pi}^{\rightarrow}$ ; an operation-call is an expression inside an  $\lambda_{\pi}^{\rightarrow}$  which actually invokes that runtime behaviour.

Realistically, operations should 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"). However, we shall see the rules of  $\lambda_{\pi}^{\rightarrow}$  are about tracking potential resource-use in a system, and so the exaxct behaviour of a particular operation call is of little interset. For this reason we make the simplifying assumption that all operations are null-ary: File.write instead of File.write("mymsg").

The base types of  $\lambda_{\pi}^{\rightarrow}$  are sets of resources, denoted by  $\{\bar{r}\}$ . If an expression e is given type  $\{\bar{r}\}$ , then evaluating e will reduce to one of the resource literals in  $\bar{r}$  (if e terminates).

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. Because functions can only be typed with this annotated arrow-type, and because the only way to incur an effect at the top-level is to be supplied a capability for it, we say every function in  $\lambda_{\pi}^{\rightarrow}$  is "annotated" by what effects they can incur.

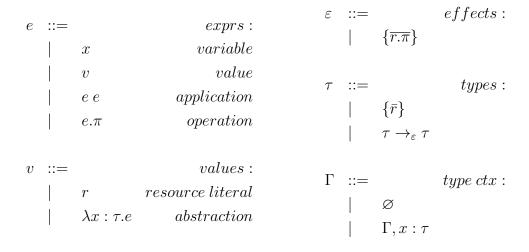


Figure 3.1: Grammar for  $\lambda_{\pi}^{\rightarrow}$ .

#### 3.2.2 $\lambda_\pi^ o$ Dynamic Rules

Before giving dynamic rules, Figure 3.2. first shows the updated definition of substitution. It is straight-forward, but in  $\lambda_{\pi}^{\rightarrow}$  we make the restriction that a variable may only be substituted for a value. 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. Because our dynamic rules will employ a call-by-value this tightening of substitution is no problem.

```
substitution :: e × v × v \rightarrow 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 3.2: Substitution function in  $\lambda_{\pi}^{\rightarrow}$ .

Rules for single-step reductions are given in Figure 3.3. Single-step reduction now takes the form  $e \longrightarrow e \mid \varepsilon$ , with the resulting pair being the expression after reduction, and the set of effects incurred during the single-step of computation (which in the case of single-step reduction is at most a singleton).

E-APP1 and E-APP2 incur whatever is the effect of reducing their subexpressions. E-APP3 incurs no effects when it performs substitution (and proving the safety of this reduction depends on our narrowed definition of substitution).

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). An important property is that unit is the only value of its type (which is called Unit). Because of this, it is used to signify the absence of information. As we have chosen not to model the semantics of operation-calls, we choose unit as a sensible result of reducing an operation-call.

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

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

Figure 3.3: Single-step reductions in  $\lambda_{\pi}^{\rightarrow}$ .

A multi-step reduction consists of zero or more single-step reductions. The resulting effect-set is the union of the effect-sets produced by all the intermediate single-steps. Rules are given in Figure 3.4.

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

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

Figure 3.4: Multi-step reductions in  $\lambda_{\pi}^{\rightarrow}$ .

#### 3.2.3 $\lambda_{\pi}^{\rightarrow}$ Static Rules

The static rules for  $\lambda_{\pi}^{\rightarrow}$  are summarised in Figure 3.5. There is the standard subtyping judgement form  $\tau <: \tau$ , and a new form judgement,  $\Gamma \vdash e : \tau$  with  $\varepsilon$ . The new form ascribes a type-and-effect to a piece of code e, meaning successive reductions of e will yield terms of type  $\tau$ , and collectively incur no more than those effects in  $\varepsilon$ . These rules give a conservative approximation to the runtime effects of executing e, so the static  $\varepsilon$  ascribed to e may include effects which don't actually 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 runtime-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: their 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 or Socket.close. In such cases, the safe approximation is to type the conditional as  $\{\text{File}, \text{Socket}\}$ .  $\varepsilon\text{-OPERCALL}$  would then approximate the runtime effects of the operation call as  $\{\text{File}.close}, \text{Socket.close}\}$ .

The rules of  $\lambda_{\pi}^{\rightarrow}$  permit any operation to be performed on any resource. This can give bizarre programs — Sensor.readTemp seems like a sensible operation call, but what about File.readTemp? We acknowledge that this allows for strange programs, but because  $\lambda_{\pi}^{\rightarrow}$  does not model the semantics of particular operatino-calls, we ignore it.

Being able to type an expression as a (non-singleton) set of resources requires the subtyping rule S-RESOURCE. This says that a subset of resources is also a subtype. To justify this rule, consider  $\{\bar{r}\} <: \{\bar{r}_2\}$ . Any value with type  $\{\bar{r}_1\}$  can reduce to any resource literal in  $\bar{r}_1$ , so to be complatible with type  $\{\bar{r}_2\}$ , the resource literals in  $\bar{r}_1$  must also be in  $\bar{r}_2$ , hence the definition.

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. This is because any possible effect which might be incurred by the subtype should be expected by the supertype, otherwise substitution of a supertyped value for a subtyped value would allow the introduction of new effects not possible under the original.

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

The goal of this section is to show  $\lambda_{\pi}^{\rightarrow}$  is sound, but this requires an appropriate notion of *effect soundness*. Intuitively, if a static judgement like  $\Gamma \vdash e : \tau$  with  $\varepsilon$  were correct, it could be read as saying that successive reductions on e will never produce effects not in the approximation  $\varepsilon$ . By adding this to our notion of soundness, we get the following first definition:

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

 $\Gamma \vdash e : \tau \text{ with } arepsilon$ 

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

Figure 3.5: Static rules of  $\lambda_{\pi}^{\rightarrow}$ .

**Theorem 5** (Soundness 1). *If*  $\Gamma \vdash e_A : \tau_A$  with  $\varepsilon_A$  and  $e_A$  is not a value, then  $e_A \longrightarrow e_B \mid \varepsilon$ , where  $\Gamma \vdash e_B : \tau_B$  with  $\varepsilon_B$  and  $\tau_B <: \tau_A$  and  $\varepsilon \subseteq \varepsilon_A$ .

In this formulation,  $\varepsilon_A$  approximates the effects of the term  $e_A$  in the context  $\Gamma$ .  $e_A$  can be reduced to  $e_B$ , incurring the runtime effects in  $\varepsilon$ . The same context can also approximate the runtime effects of  $e_B$  as  $\varepsilon_B$ , meaning the term after reduction can be typed, but no additional information about  $\varepsilon_B$  is stipulated.

Our approach to proving that multi-step reduction is sound will be to inductively appeal to the soundness of single-step reductions. This is tricky under the given definition of Soundness because it only relates the runtime effects to the approximation of the runtime effects *before* reduction. There is constraint on the runtime effects *after* reduction. To accommodate a proof of multi-step soundness, we need a stronger version of soundness which relates the approximated effects before reduction ( $\varepsilon_A$ ) to the approximated effects after reduction ( $\varepsilon_B$ ).

First consider the analogous relation for the types of temrs before and after reduction. In  $\lambda$ -calculi, the type after reduction can be the same or more specific (i.e.  $\tau_B <: \tau_A$ ) than the type before reduction. But it can never be less specific. Similarly, we shall require the approximated effects of a type can get more specific after reduction, but never less-specific.

To illustrate why the approximated effects might get more specific, consider the function get  $= \lambda x$ : {File, Socket}.x and the program (f File).write. In the context  $\Gamma$  = File: {File}, the rule  $\varepsilon$ -APP can be used to approximate the effects of (f File).write as {File.write, Socket.write}. By E-APP3 we have the reduction (get File).write  $\longrightarrow$ 

File.write  $| \varnothing$ . The same context can use  $\varepsilon$ -OPERCALL to approximate the reduced expression File.write as {File.write}; note how the approximation of effects is more precise after reduction. This example shows why the approximation after reduction ( $\varepsilon_B$ ) should be a subset of the approximation before reduction ( $\varepsilon_A$ ).

We have our final definition of soundness:

**Theorem 6** (Soundness). *If*  $\Gamma \vdash e_A : \tau_A$  with  $\varepsilon_A$  and  $e_A$  is not a value, then  $e_A \longrightarrow e_B \mid \varepsilon$ , where  $\Gamma \vdash e_B : \tau_B$  with  $\varepsilon_B$  and  $\tau_B <: \tau_A$  and  $\varepsilon_B \cup \varepsilon \subseteq \varepsilon_A$ .

We take the standard road to proving soundness by showing that progress and preservation hold of  $\lambda_{\pi}^{\rightarrow}$ , which in turn rely on modified versions of Canonical Forms and the Substitution Lemma.

Canonical Forms for  $\lambda_{\pi}^{\rightarrow}$  states that resource-typed values are resource literals, and any typing judgement of a value will approximate the runtime effects as  $\varnothing$ . This result is not true if the rule used is  $\varepsilon$ -SUBSUME, so the lemma statement excludes judgements which use that rule. Progress follows from Canonical Forms.

**Lemma 3** (Canonical Forms). *Unless the rule used is*  $\varepsilon$ -SUBSUME, the following are true:

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

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

*Proof.* By induction on  $\Gamma \vdash e : \tau$  with  $\varepsilon$ , for e not a value. If the rule is  $\varepsilon$ -Subsumption it follows by inductive hypothesis. If e has a reducible subexpression then reduce it. Otherwise use one of  $\varepsilon$ -App3 or  $\varepsilon$ -OperCall2.

To show preservation holds we need to know that type-and-effect safety, as it has been formulated in the definition of soundness, is preserved by the substitution in E-APP3. As noted earlier, variables can only be substituted for values in  $\lambda_{\pi}^{\rightarrow}$ . Canonical Forms tells us that any value will have its effects approximated as  $\varnothing$  (except when  $\varepsilon$ -SUBSUME is used). This leads to  $\lambda_{\pi}^{\rightarrow}$  having a slightly stronger formulation than the  $\lambda^{\rightarrow}$  equivalent. The proof is routine.

**Lemma 4** (Substitution). *If*  $\Gamma, x : \tau' \vdash e : \tau$  with  $\varepsilon$  and  $\Gamma \vdash v : \tau'$  with  $\varnothing$  then  $\Gamma \vdash [v/x]e : \tau$  with  $\varepsilon$ .

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

With this lemma, we are ready to prove the preservation theorem.

**Theorem 8** (Preservation). *If*  $\Gamma \vdash e_A : \tau_A \text{ with } \varepsilon_A \text{ and } e_A \longrightarrow e_B \mid \varepsilon, \text{ then } \tau_B <: \tau_A \text{ and } \varepsilon_B \cup \varepsilon \subseteq \varepsilon_A.$ 

*Proof.* By induction on  $\Gamma \vdash e_A : \tau_A$  with  $\varepsilon_A$ , and then on  $e_A \longrightarrow e_B \mid \varepsilon$ . Since  $e_A$  can be reduced, we need only consider those rules which apply to non-values and non-variables.

Case:  $\varepsilon$ -APP Then  $e_A = e_1 \ e_2$  and  $e_1 : \tau_2 \to_{\varepsilon} \tau_3$  with  $\varepsilon_1$  and  $\Gamma \vdash e_2 : \tau_2$  with  $\varepsilon_2$ . If the reduction rule used was E-APP1 or E-APP2, then the result follows by applying the inductive hypothesis to  $e_1$  and  $e_2$  respectively.

Otherwise the rule used was E-APP3. Then  $(\lambda x: \tau_2.e)v_2 \longrightarrow [v_2/x]e \mid \varnothing$ . By inversion on the typing rule for  $\lambda x: \tau_2.e$  we know  $\Gamma, x: \tau_2 \vdash e: \tau_3$  with  $\varepsilon_3$ . By canonical forms,  $\varepsilon_2 = \varnothing$  because  $e_2 = v_2$  is a value. Then by the substitution lemma,  $\Gamma \vdash [v_2/x]e: \tau_3$  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. Then  $e_A = e_1.\pi$  and  $\Gamma \vdash e_1 : \{\bar{r}\}$  with  $\varepsilon_1$ . If the reduction rule used was E-OPERCALL1 then the result follows by applying the inductive hypothesis to  $e_1$ .

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

Our single-step soundness theorem now holds immediately by joining the progress and preservation theorems into one.

**Theorem 9** (Soundness). *If*  $\Gamma \vdash e_A : \tau_A$  with  $\varepsilon_A$  and  $e_A$  is not a value, then  $e_A \longrightarrow e_B \mid \varepsilon$ , where  $\Gamma \vdash e_B : \tau_B$  with  $\varepsilon_B$  and  $\tau_B <: \tau_A$  and  $\varepsilon_B \cup \varepsilon \subseteq \varepsilon_A$ .

*Proof.* If  $e_A$  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, the soundness of multi-step reductions can be shown by inductively applying single-step soundness on the length of a multi-step reduction.

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

*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.

This concludes the soundness proof for  $\lambda_{\pi}^{\rightarrow}$ . As we have seen,  $\lambda_{\pi}^{\rightarrow}$  builds upon  $\lambda^{\rightarrow}$  with resources and operations. Every function must have its input type labelled with the

effects any values of that type might incur. This allows us to easily effect-check a piece of code to ascertain what runtime effects it might have when executed. And we have just proven it sound.

### 3.3 $\lambda_{\pi,\varepsilon}^{\rightarrow}$ : Epsilon Calculus

 $\lambda_{\pi}^{\rightarrow}$  requires every function to have its input type annotated — if we relax this requirement, can our type system say anything useful about pieces of unannotated code? In this section we introduce  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ , which leverages capability-safe design to enable our type-system to approximate what effects might be incurred by unannotated code.

There are practical reasons to permit unannotated code in an effect-conscious language. Previous effect systems have been criticised for their verbosity, **citation needed** which might disincline a programmer from bothering to use them. Allowing the structured mixture of annotated and unannotated also permits one to rapidly prototype software in the unannotated sublanguage and then incrementally add effect annotations as they are needed, giving a trade-off between convenience and safety.

In general, reasoning about unannotated code is difficult because there are no constraints on what effects they might incur. Figure 3.6 demonstrates the issue with a Wyvern-like snippet of code. someMethod takes a function as input and executes it; but the effects of f depend on what particular implementation is passed to someMethod. Without more information, whether by means of a more complex type-system or more static annotations, there is no general way to know what effects might be incurred by someMethod.

```
1 def someMethod(f: Unit → Unit):
2 f()
```

Figure 3.6: What effects might be incurred by someMethod?

Therefore,  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  permits the nesting of unannotated code when it is nested inside annotated code with the import expression. The typing rule for import expression,  $\varepsilon$ -IMPORT, imposes a set of restrictions on what authority can be exercised by the unannotated code. This enables the type system to safely approximate the effects of unannotated code as those effects captured by the capabilities selected by enclosing import expressions. This is the key result of  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ , which we formalise and prove in this section.

#### 3.3.1 $\lambda_{\pi,\varepsilon}^{\rightarrow}$ Grammar

The grammar of  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  is essentially split into rules for annotated code and analogous rules for unannotated code. To distinguish the two, annotated types, expressions, and contexts always have a hat above them:  $\hat{e}$ ,  $\hat{\tau}$ , and  $\hat{\Gamma}$  instead of e,  $\tau$ , and  $\Gamma$ .

The unannotated portion consists of programs made from the same building blocks as  $\lambda_{\pi}^{\rightarrow}$ , but with the regular type-constructor  $\rightarrow$  from  $\lambda^{\rightarrow}$ . Unannotated programs have no labels on their functions at all. They are also *deeply* unannotated: if a term is unannotated at the top-level, then every sub-term is also unannotated, so you cannot nest annotated code inside unannotated code. The corresponding grammar of types contains unannotated types  $\tau$ , which are built from resource-stes and  $\rightarrow$ , and unannotated contexts  $\Gamma$ . An unannotated context can only map variables to unannotated types.

Except for the new import expression, the analogous rules for annotated programs and their surrounding meta-theory is the same as  $\lambda_{\pi}^{\rightarrow}$ . Annotated types  $\hat{\tau}$  are those built using the type constructors  $\rightarrow_{\varepsilon}$ , where  $\varepsilon$  is a concrete set of effects. The category of annotated contexts is  $\hat{\Gamma}$ , which binds variables to annotated types.

The interesting new form is import, which belongs to the annotated sublanguage. import introduces a name x with annotated definition  $\hat{e}$  into a body of unannotated code e.  $\varepsilon$  is the set of effects selected by the unannotated code, so any resources and operation calls used in e must be declared in  $\varepsilon$ . import is the only means of nesting unannotated subterms inside annotated terms.

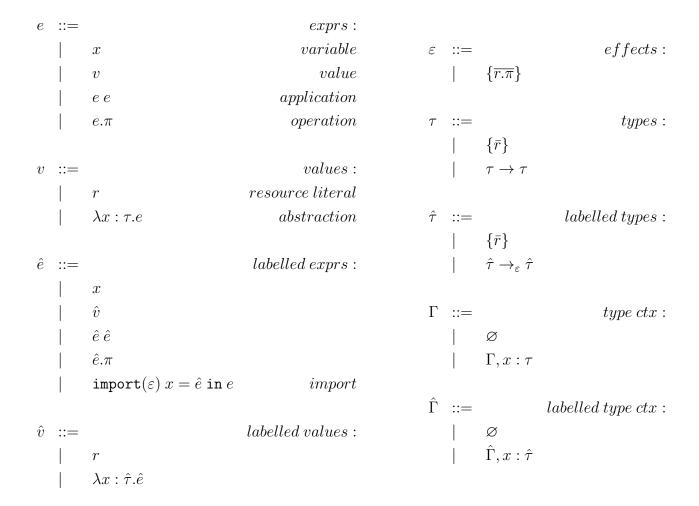


Figure 3.7: Effect calculus.

### **3.3.2** $\lambda_{\pi,\varepsilon}^{\rightarrow}$ Dynamic Rules

Different approaches can be taken to defining the execution of an  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ . One way is to define reductions for both annotated and unannotated programs, but this results in a lot of uninteresting rules which clutter the formalism. Another way is to define reductions for either annotated or unannotated programs, and translate programs into the appropriate form before executing them. We choose this approach, but the transformation happens during execution of the program, rather than before the program is executed. The idea is that whenever a piece of unlabelled code is encountered, the import expression surrounding it will have been evaluated to the point where we know what effects  $\varepsilon$  are being selected by the unannotated body e. At this point, we can annotate e with  $\varepsilon$  and continue executing this result.

To this end we define annot in Figure 3.8. This function takes a piece of unlabelled code e and a set of effects  $\varepsilon$  and produces  $\hat{e}$ , obtained by labelling every arrow-type with  $\varepsilon$ . A version of this function is given for expressions, types, and contexts. Its dual is erase, which deletes all annotations from code. We will need erase in the definition of  $\varepsilon$ -IMPORT, so we give its definition here. Like annot, there are versions which delete the annotations from expressions, contexts, and types. However, erase is partial because it is not defined on import expressions.

It is worth mentioning that annot and erase operate on a purely syntactic level. Their definitions state no meaningful correspondence between the types and effects of a program before and after it has been annotated or erased. It remains for us to prove there is a meaningful correspondence in the particular way these functions are used in the rules for import expressions.

Finally, before giving the dynamic rules we must define a version of substitution for  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ . As our dynamic rules are defined on annotated expressions, so too will substitution be defined. The definition, given in 3.9, is otherwise straight-forward. It has been updated for the new annotated notation and the new import expression. It retains the same restriction from  $\lambda_{\pi}^{\rightarrow}$  where substitution is only well-defined when a variable is replaced with a value.

The multi-step rules of  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  are identical to  $\lambda_{\pi}^{\rightarrow}$ . Every single-step rule in  $\lambda_{\pi}^{\rightarrow}$  is also a single-step rule in  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ . The only difference is that the new annotated notation in  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  would have a hat above every expression and type. For brevity, we omit those rules which are (basically) identical to those of  $\lambda_{\pi}^{\rightarrow}$ .

The two new single-step reductions are given in 3.10. Both reduce import expressions. E-IMPORT1 reduces the definition of the capability being imported into the unannotated code. The more interesting rule is E-IMPORT2, which applies when the capability being impoted is a value. The unlabelled body e is annotated with the authority  $\varepsilon$  it selects; this is  $\operatorname{annot}(e,\varepsilon)$ . The name x of the capability is then replaced with its actual definition  $\hat{v}$ ; this is  $[\hat{v}/x]\operatorname{annot}(e,\varepsilon)$ . This reduction incurs no effects.

```
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}\}
           \mathtt{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}\})
           \mathtt{erase}(\hat{\tau}_1 \to_\varepsilon \hat{\tau}_2) = \mathtt{erase}(\hat{\tau}_1) \to \mathtt{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 3.8: Definitions of annot and erase.

### 3.3.3 $\lambda_{\pi,\varepsilon}^{\rightarrow}$ Static Rules

Our goal in this section is to introduce  $\varepsilon$ -IMPORT, which approximates the effects of unannotated code by inspecting its selected authority. The rule is complicated, so we build up to it.

First, since programs in  $\lambda_{\pi,\varepsilon}^{\to}$  can be annotated or unannotated, we need to be able to recognise when either kind is well-typed. Since the annotated subset of  $\lambda_{\pi,\varepsilon}^{\to}$  contains  $\lambda_{\pi}^{\to}$ , all the  $\lambda_{\pi}^{\to}$  still apply, but the notation is different: we put hats on everything to signify that a typing judgement is being made about annotated code inside an annotated context. This looks like  $\hat{\Gamma} \vdash \hat{e} : \hat{\tau}$  with  $\varepsilon$ . Except for this change in notation the rules are the same, so we shall not repeat them.

The rules for typing unannotated pieces of code take the form  $\Gamma \vdash e : \tau$ . The subtyping judgement for unannotated code takes the form  $\tau <: \tau$ . A summary of these typing and subtyping rules is given in 3.11; each is analogous to some rule in  $\lambda_{\pi}^{\rightarrow}$ , but the parts relating to effects have been removed.

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

```
\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: Definition of substitution.

$$\begin{split} \frac{\hat{e} \longrightarrow \hat{e} \mid \varepsilon}{\text{import}(\varepsilon) \; x = \hat{e} \; \text{in} \; e \longrightarrow \text{import}(\varepsilon) \; x = \hat{e}' \; \text{in} \; e \mid \varepsilon'} \; \text{(E-IMPORT1)} \\ \\ \frac{\text{import}(\varepsilon) \; x = \hat{e} \; \text{in} \; e \longrightarrow \text{import}(\varepsilon) \; x = \hat{e}' \; \text{in} \; e \mid \varepsilon'}{\text{import}(\varepsilon) \; x = \hat{v} \; \text{in} \; e \longrightarrow [\hat{v}/x] \\ \text{annot}(e, \varepsilon) \mid \varnothing} \; \text{(E-IMPORT2)} \end{split}$$

Figure 3.10: New single-step reductions in  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ .

Again, there are no rules which directly approximate the effects of unannotated code. The only mechanism for doing this is to encapsulate that code in an import expression, which restricts what authority it can exercise. This means the rules can only tell us something interesting about annotated code and unannotated code which is nested inside an import expression. For the rest of this section, we are going to build up to the definition of  $\varepsilon$ -IMPORT.

To begin, typing  $\operatorname{import}(\varepsilon)$   $x=\hat{e}$  in e requires us to know that the capability  $\hat{e}$  being imported is well-typed, so we add the premise  $\hat{\Gamma}\vdash\hat{e}:\hat{\tau}$  with  $\varepsilon_1$ . Whatever authority is exercised by e must have been selected by the import expression, so you should be able to type e with a binding for  $x:\hat{\tau}$ . However,  $\hat{\tau}$  is annotated and e is unannotated, so we erase the labels on it. This gives our second premise,  $x:\operatorname{erase}(\hat{\tau})\vdash e:\tau$ . These observations give our first definition of  $\varepsilon$ -IMPORT in Figure 3.12. Since E-IMPORT2 labels the unannotated code code with the selected authority  $\varepsilon$ , an import expression should type to  $\operatorname{annot}(\tau,\varepsilon)$ . The safe approximation of the unannotated code's effects is  $\varepsilon_1 \cup \varepsilon$ ; the former comes from reducing the imported capability (which happens prior to the execution of the unannotated code) and the latter contains all the authority which the unannotated code may use.

On the surface, this rule may seem overly restrictive because import only allows one to import a single capability. What about unannotated code which uses multiple capabilities? One solution is to reformulate the rules (and grammar) to permit any number of imports. Another is to import multiple capabilities by importing a tuple of capabilities; for example  $import(\varepsilon)$  x = (File, Socket) in e. We have not presented tuples as a part of

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

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

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

 $\tau <: \tau$ 

$$\frac{\tau_1' <: \tau_1 \quad \tau_2 <: \tau_2'}{\tau_1 \rightarrow \tau_2 <: \tau_1' \rightarrow \tau_2'} \text{ (S-ARROW)} \quad \frac{\{\bar{r}_1\} \subseteq \{\bar{r}_2\}}{\{\bar{r}_1\} <: \{\bar{r}_2\}} \text{ (S-RESOURCES)}$$

Figure 3.11: (Sub)typing judgements for the unannotated sublanguage of  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ 

 $\tau <: \tau$ 

$$\frac{\hat{\Gamma} \vdash \hat{e} : \hat{\tau} \; \mathtt{with} \; \varepsilon_1 \quad x : \mathtt{erase}(\hat{\tau}) \vdash e : \tau}{\hat{\Gamma} \vdash \mathtt{import}(\varepsilon) \; x = \hat{e} \; \mathtt{in} \; e : \mathtt{annot}(\tau, \varepsilon) \; \mathtt{with} \; \varepsilon \cup \varepsilon_1} \; (\varepsilon\text{-IMPORT1})$$

Figure 3.12: A first rule for type-and-effect checking import expressions.

the base language, but they could be encoded as a derived form, or added as a language extension. Both approaches are straightforward and extending  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  in these ways would be routine — but only allowing the import of a single capability reduces clutter in the rules and simplifies the presentation, so we stick to this convention.

The first version of the rule has some issues. First of all, there is no formal relation between  $\varepsilon$  and whatever effects are captured by  $\hat{\tau}$ . Consider  $\hat{e} = \mathrm{import}(\varnothing) \ x = \mathrm{File} \ \mathrm{in} \ \mathrm{x.write}$ , which imports a File and writes to it, but declares its authority as  $\varnothing$ . According to  $\varepsilon$ -IMPORT1,  $\vdash \hat{e} : \mathrm{Unit} \ \mathrm{with} \ \varnothing$ , but this is clearly wrong. We need to add a constraint to say that whatever effects are captured by  $\hat{\tau}$  should be selected in  $\varepsilon$ . To this end we define a function, effects, which collects the set of effects that an annotated type captures. A first definition is given in Figure 3.13. With it, we can add an extra premise effects( $\hat{\tau}$ )  $\subseteq \varepsilon$  which formalises the idea that any capability exercised by e must be selected in  $\varepsilon$ . The updated rule is given in Figure 3.14.

effects  $:: \hat{ au} 
ightarrow arepsilon$ 

$$\begin{split} & \texttt{effects}(\{\bar{r}\}) = \{r.\pi \mid r \in \bar{r}, \pi \in \Pi\} \\ & \texttt{effects}(\hat{\tau}_1 \to_{\varepsilon} \hat{\tau}_2) = \texttt{effects}(\hat{\tau}_1) \cup \varepsilon \cup \texttt{effects}(\hat{\tau}_2) \end{split}$$

Figure 3.13: A first definition of effects.

The counterexample which defeated  $\varepsilon$ -IMPORT1 is now rejected by  $\varepsilon$ -IMPORT2, but

```
\frac{\hat{\Gamma} \vdash \hat{e} : \hat{\tau} \; \mathtt{with} \; \varepsilon_1 \quad x : \mathtt{erase}(\hat{\tau}) \vdash e : \tau \quad \mathtt{effects}(\hat{\tau}) \subseteq \varepsilon}{\hat{\Gamma} \vdash \mathtt{import}(\varepsilon) \; x = \hat{e} \; \mathtt{in} \; e : \mathtt{annot}(\tau, \varepsilon) \; \mathtt{with} \; \varepsilon \cup \varepsilon_1} \; (\varepsilon\text{-IMPORT2})
```

Figure 3.14: A second rule for type-and-effect checking import expressions.

we lose some of the precision gained by effect annotations. To illustrate, consider  $\hat{e} = \mathtt{import}(\mathtt{File.write}) \ x = ((\lambda \mathtt{y} : \{\mathtt{File}\}.\lambda\mathtt{z} : \mathtt{Unit.y.write}) \ \mathtt{File}) \ \mathtt{in} \ \mathtt{x} \ \mathtt{unit}.$  The capability being imported is a function which, when given unit, performs a write operation. The stated

A precise definition of effects is given in Figure 3.14. along with an associated function, ho-effects. The difference between the two is the difference between direct and transitive authority. If  $r.\pi \in \texttt{effects}(\hat{\tau})$ , then values of  $\hat{\tau}$  have the authority to directly incur  $r.\pi$ . If  $r.\pi \in \texttt{ho-effects}(\hat{\tau})$ , then values of  $\hat{\tau}$  can incur  $r.\pi$  by deferring to another function. The two are mutually recursive, with resource-types as a base-case. A resource captures every operation on itself, because resource have the authority to call any operation on themselves. A resource captures no higher-order effects.

```
\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 3.15: Effect functions.

arepsilon-IMPORT2 is better, but still puts no constraints on what higher-order effects might be captured by  $\hat{\tau}$ . At the moment, so long as  $\hat{e}$  can only directly invoke those effects selected by  $\varepsilon$ , it is allowed to nest any other type with arbitrary authority. On the other hand, it is safe to pass in something of type  $\hat{\tau}$  if every callable function inside e expects the captured higher-order effects. Intuitively, the functions of e need to have selected ho-effects( $\hat{\tau}$ ), otherwise we may be exceeding their authority in giving them access to  $\hat{\tau}$ . This motivates the predicates  $\mathrm{safe}(\hat{\tau},\varepsilon)$  and ho- $\mathrm{safe}(\hat{\tau},\varepsilon)$ . A type is safe for  $\varepsilon$  if every directly invocable function selects the authority in  $\varepsilon$ . A type is higher-order safe for  $\varepsilon$  if every indirectly function selects the authority in  $\varepsilon$ . 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,  $\mathrm{safe}(\hat{\tau},\varepsilon)$  holds when the (non higher-order) effects selected by  $\hat{\tau}$  include  $\varepsilon$ . ho- $\mathrm{safe}(\hat{\tau},\varepsilon)$  holds when the higher-order effects selected by  $\hat{\tau}$  include  $\varepsilon$ . 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. Formal definitions are given in Figure 3.15.

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

Figure 3.16: Safety judgements in the epsilon calculus.

The final version of  $\varepsilon$ -IMPORT is given in Figure 3.16. It requires the import  $\hat{e}$  to be well-typed. The effects captured by the import  $\operatorname{effects}(\hat{\tau})$  must be the same as those effects selected by the body of the import. The import must be higher-order safe; that is, every possible function that could be invoked by the import must be expecting the effects declared in  $\varepsilon$ . Lastly, the body of the import must be well-formed with only a binding for the imported name.

```
\begin{split} \hat{\Gamma} \vdash \hat{e} : \hat{\tau} \text{ with } \varepsilon \\ & \qquad \qquad \hat{\Gamma} \vdash \hat{e} : \hat{\tau} \text{ with } \varepsilon_1 \quad \varepsilon = \text{effects}(\hat{\tau}) \\ & \qquad \qquad \qquad \frac{\text{ho-safe}(\hat{\tau}, \varepsilon) \quad x : \text{erase}(\hat{\tau}) \vdash e : \tau}{\hat{\Gamma} \vdash \text{import}(\varepsilon) \ x = \hat{e} \text{ in } e : \text{annot}(\tau, \varepsilon) \text{ with } \varepsilon \cup \varepsilon_1} \end{split}  Figure 3.17: Type-with-effect judgements.
```

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

Soundness in  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  is much the same as it is in  $\lambda_{\pi}^{\rightarrow}$ , but only for annotated programs. A definition is given below.

**Theorem 11** (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  $\varepsilon_B \cup \varepsilon \subseteq \varepsilon_A$ .

Because the rules of  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  are the same as  $\lambda_{\pi}^{\rightarrow}$ , we simply extend the existing proofs to cover the case where the typing rule used is  $\varepsilon$ -IMPORT. Canonical Forms remains unchanged. The Substitution Lemma gains an extra case, but the proof is routine.

**Lemma 5** (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\}.
```

**Lemma 6** (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  $\varepsilon$ .

The Progress Theorem now has an extra case: when the typing rule used is  $\varepsilon$ -IMPORT. The result follows by considering whether the imported  $\hat{e}$  in import( $\varepsilon$ )  $x=\hat{e}$  in e is an expression or not.

**Theorem 12** (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.* If the rule is  $\varepsilon$ -IMPORT then  $e = \mathsf{import}(\varepsilon) \ x = \hat{e} \ \mathsf{in} \ e$ . If  $\hat{e}$  is a non-value then it reduces by inductive assumption and the import reduces via  $\varepsilon$ -IMPORT1. Otherwise  $\hat{e}$  is a value and the import reduces via  $\varepsilon$ -IMPORT2.

Likewise, the preservation theorem gains an extra case when  $\varepsilon$ -IMPORT is the typing rule used and E-IMPORT2 is the reduction rule used. To show the reduction  $\mathrm{import}(\varepsilon)\,x=\hat{v}\,\mathrm{in}\,e\,\longrightarrow\,[\hat{v}/x]\mathrm{annot}(e,\varepsilon)\mid\varnothing$  preserves soundness requires a few things. First, if  $\hat{\Gamma}\vdash\mathrm{import}(\varepsilon)\,x=\hat{v}\,\mathrm{in}\,e:\hat{\tau}_A\,\mathrm{with}\,\varepsilon_A$ , then we need to be able to type the reduced expression in the same context:  $\hat{\Gamma}\vdash[\hat{v}/x]\mathrm{annot}(e,\varepsilon):\hat{\tau}_B\,\mathrm{with}\,\varepsilon_B$ , where the type and effects are preserved. Our proof strategy for this case is to do this in two parts. First we show taht the typing judgement  $\hat{\Gamma}\vdash\mathrm{annot}(e,\varepsilon):\hat{\tau}_B\,\mathrm{with}\,\varepsilon_B$  can be made; then we the same judgement will hold of  $[\hat{v}/x]\mathrm{annot}(e,\varepsilon)$  in the same context by the substitution lemma. To prove the first part can be done, we introduce a new lemma.

**Lemma 7** (Annotation). *If the following are true:* 

```
\begin{array}{ll} \bullet & \hat{\Gamma} \vdash \hat{v} : \hat{\tau} \; \mathtt{with} \; \varnothing \\ \bullet & \Gamma, y : \mathtt{erase}(\hat{\tau}) \vdash e : \tau \\ \bullet & \varepsilon = \mathtt{effects}(\hat{\tau}) \\ \bullet & \mathtt{ho\text{-safe}}(\hat{\tau}, \varepsilon) \end{array}
```

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

*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.  $\Gamma$  encapsulates the ambeint authority exercised by e. At the top-level of any program, we will always have  $\Gamma = \varnothing$ ; compare this with how the premise of  $\varepsilon$ -IMPORT types the body of an import expression with only a single binding for the import. 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, this is not ambient authority at the top-level.

Proof of the Annotation lemma is long but routine, save for the use of an additional pair of lemmas. These lemmas relate  $\hat{\tau}$  and annot(erase( $\hat{\tau}$ ),  $\varepsilon$ ).

**Lemma 8.** If effects( $\hat{\tau}$ )  $\subseteq \varepsilon$  and ho-safe( $\hat{\tau}, \varepsilon$ ) then  $\hat{\tau} <:$  annot(erase( $\hat{\tau}$ ),  $\varepsilon$ ).

**Lemma 9.** 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 13** (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:  $\varepsilon$ -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,  $\tau_A = \tau_B = \mathtt{annot}(\tau, \varepsilon)$ .

We can now combine Progress and Preservation into the Soundness theorem for  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ . The proof of multi-step soundness in  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  is identical to the proof in  $\lambda_{\pi}^{\rightarrow}$ .

**Theorem 14** (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  $\varepsilon_B \cup \varepsilon \subseteq \varepsilon_A$ .

**Theorem 15** (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  $\varepsilon_B \cup \varepsilon \subseteq \varepsilon_A$ .

## Chapter 4

## **Applications**

In this chapter we show how  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  can be used in practice, and show how its rules can enable effect reasoning in existing capability-safe languages. This will take the form of writing a program in a high-level, capability-safe language, translating it to an equivalent  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  program, and demonstrating how the rules of  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  enable reasoning about the use of effects.

In this section the high-level programs will be written in a version of Wyvern. Wyvern is a pure, object-oriented, capability-safe language. It has a first-class module system, in which modules and objects are treated uniformly. Although  $\lambda_{\pi,\varepsilon}^{\to}$  does not have objects, the example Wyvern programs can be expressed using functions. This does not mean the examples given aren't demonstrating useful and realistic situations — we simply do not need the added expressiveness given by self-referential objects.

In section 4.1. we discuss how the translation from Wyvern to  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  will work, and what simplifying assumptions are made in our examples. This also serves as a gentle introduction to Wyvern's syntax. A variety of scenarios are then explored in 4.2. to show how the rules of  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  can help developers in practice.

### 4.1 Translations and Encodings

Our aim is to develop some notation to help us translate Wyvern programs into  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ . Our approach will be to encode these additional rules and forms into the base language of  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ ; essentially, to give common patterns and forms a short-hand, so they can be easily named and recalled. This is called *sugaring*. When these derived forms are collapsed into their underlying representation, it is called *desugaring*. We are going to introduce several rules to show a Wyvern program might be considered syntactic sugar for an  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  program, and translate examples by desugaring according to our rules.

#### 4.1.1 Unit

Unit is a type inhabited by exactly one value. It conveys the absence of information; in  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  an operation call on a resource literal reduces to unit for this reason. We define unit  $\stackrel{\text{def}}{=} \lambda x : \varnothing.x$ . The unit literal is the same in both annotated and naked code. In annotated code, it has the type Unit  $\stackrel{\text{def}}{=} \varnothing \rightarrow_{\varnothing} \varnothing$ , while in naked code it has the type Unit  $\stackrel{\text{def}}{=} \varnothing \rightarrow \varnothing$ . While these are technically two seperate types, we will not distinguish between the annotated and naked versions, simply referring to them both as Unit.

Note that unit is a value, and because  $\varnothing$  is uninhabited (there is no empty resource literal), unit cannot be applied to anything. Furthermore,  $\vdash$  unit : Unit with  $\varnothing$  by  $\varepsilon$ -ABS, and  $\vdash$  unit : Unit by T-ABS. This leads to the derived rules in 4.1.

```
\begin{array}{c} \hline \Gamma \vdash e : \tau \\ \\ \hline \hat{\Gamma} \vdash \hat{e} : \hat{\tau} \; \text{with} \; \varepsilon \\ \\ \hline \hline \Gamma \vdash \text{unit} : \text{Unit} \; \; (\text{T-UNIT}) \quad \overline{\hat{\Gamma} \vdash \text{unit} : \text{Unit} \; \text{with} \; \varnothing} \; \left( \varepsilon \text{-UNIT} \right) \end{array}
```

Figure 4.1: Derived Unit rules.

Since unit represents the absence of information, we also use it as the type when a function either takes no argument, or returns nothing. 4.2 shows the definition of a Wyvern function which takes no argument and returns nothing, and its corresponding representation in  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ .

```
def method():Unit unit \lambda x:Unit. unit
```

Figure 4.2: Desugaring of functions which take no arguments or return nothing.

#### 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). The new rules in 4.4 capture these derivations.

let expressions can be used to sequence computations. Intuitively, the let expression simply names the results of the intemediate steps and then ignores them in its body.

Figure 4.3: Derived 1et rules.

When we ignore the result of a computation we shall bind it to \_ instead of a real name, to suggest the result isn't important and prevent the naming of unused variables. 4.4 shows how this is done.

```
def method(f: {File}):Unit with {File.open, File.write, File.close}

f.open
f.write(''hello, world!'')
f.close

\( \lambda f: \{\text{File}\}. \)
let _ = f.open in
let _ = f.write in
f.close
```

Figure 4.4: Desugaring of a sequence of computations.

### 4.1.3 Modules and Objects

Wyvern's modules are first-class and desugar into objects; invoking a method inside a module is no different from invoking an object's method. There are two kinds of modules: pure and resourceful. For our purposes, a pure module is one with no (transitive) authority over any resources, while a resource module has (transitive) authority over some resource. A pure module may still be given a capability, for example by requesting it in a function signature, but it may not possess or capture the capability for longer than the duration of the method call. 4.5 shows an example of two modules, one pure and one

resourceful, each declared in a seperate file. Note how pure modules are declared with the module keyword, while resource modules are declared with the resource module keywords.

```
module PureMod

def tick(f: {File}):Unit
   f.append

resource module ResourceMod
require File

def tick():Unit with {File.append}
   File.append
```

Figure 4.5: Definition of two modules, one pure and the other resourceful.

Wyvern is capability-safe, so resource modules must be instantiated with the capabilities they require. In 4.5, ResourceMod requests the use of a File capability, which must be supplied to it from someone already possessing it. Modules are behaving like objects in this way, because they require explicit instantiation. 4.6 demonstrates how the two modules above would be instantiated and used.

To prevent infinite regress the File must, at some point, be introduced into the program. This happens in a special main module. When the program begins execution, the File capability is passed into the program from the system environment. All these initial capabilities are modelled in  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  as resource literals. They are then propagated by the top-level entry point.

```
require File
instantiate PureMod
instantiate ResourceMod(File)

def main():Unit
PureMod.tick(File)
ResourceMod.tick()
```

Figure 4.6: Definition of two modules, one pure and the other resourceful.

Before explaining our translation of Wyvern programs into  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ , we must explain several simplifications made in all of our examples which enable our particular desugaring.

Objects are only ever used in the form of modules. Modules only ever contain functions and other modules, and have no mutable fields. The examples contain no recursion

or self-reference, including a module invoking its own functions. Modules will not reference each other cyclically. Lastly, modules only contain one function definition. Despite these simplifications, the chosen examples will highlight the essential aspects of  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ .

Because modules do not exercise self-reference and only contain one function definition, they will be modelled as functions in  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ . Applying this function will be equivalent to applying the single function definition in the module.

A collection of modules is desugared into  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  as follows. First, a sequence of letbindings are used to name constructor functions which, when given the capabilities requested by a module, will return an instance of the module. If the module does not require any capabilities then it will take Unit as its argument. The constructor function for M is called MakeM. A function is then defined which represents the main function, which is the entry point into the program. This main function will instantiate all the modules by invoking the constructor functions, and then execute the body of code in main. Finally, the main function is invoked with the primitive capabilities it needs.

To demonstrate this process, 4.7 shows how the examples above desugar. Lines 1-3 define the constructor for PureMod; since PureMod requires no capabilities, the constructor takes Unit as an argument on line 2. Lines 6-8 define the constructor for ResourceMod; it requires a File capability, so the constructor takes {File} as its input type on line 7. The entry point to the program is defiend on lines 11-15, which invokes the constructors and then runs the body of the main method. Lastly, line 17 starts everything off by invoking Main with the initial set of capabilities, which in this case is just File.

### 4.2 Examples

In this section we present several scenarios where a developer may be forced to reason about the use of effects, and show how the capability-based reasoning of effects can assist them. In some scenarios, a program exhibits a certain nefarious behaviour, in which case capability-based reasoning can automatically detect this behaviour and reject it. Other scenarios are more qualitative; perhaps a developer must make a design choice and none of the alternatives *prima facie* stand out. In such cases, capability-based reasoning might supply them with useful information, enabling tehm to make more informed design choices. We also hope to convince the reader that the rules of  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  have practical worth, and could be used to enrich existing capability-safe languages.

The format of each section is as follows. A program is introduced which exhibits some bad behaviour or demonstrates a particular story about software development. The language used is *Wyvern*; a pure, object-oriented, capability-safe language with first-class modules-as-objects. We show how the Wyvern program can be written as a corresponding  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  program and sketch a derivation showing how the rules of  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  and a sketch a derivation showing how the relevant problem.

```
let MakePureMod =
      \lambda x: Unit.
         \lambda f:\{File\}. f.append
3
   in
4
5
   let MakeResourceMod =
      \lambda f: \{File\}.
         \lambdax:Unit. f.append
8
   in
9
10
   let MakeMain =
11
      \lambda f:\{File\}.
12
         \lambda x: Unit.
13
            let PureMod = (MakePureMod unit) in
14
           let ResourceMod = (MakeResourceMod f) in
15
            let _ = (PureMod f) in (ResourceMod unit) in
16
17
   (MakeMain File) unit
18
```

Figure 4.7: Desugaring of PureMod and ResourceMod into  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ .

We take some shortcuts with the translation of Wyvern into  $\lambda_{\pi,\varepsilon}^{\rightarrow}$ . Our "objects" are really records of functions; the difference between the two is self-reference. The particular examples chosen do not require self-reference, so no important properties are lost by treating Wyvern objects as records.

#### 4.2.1 Unannotated Client

In Figure 4.8 an annotated Logger module provides its client the ability to append to a particular File resource. File is a primitive capability passed into the program when it begins execution, perhaps from the system environment or a virtual machine. The Logger module presents a controlled subset of operations on the File viz. File.append. The program consists of an unannotated client which instantiates the Logger module with the capability it selects (File) and then attempts to log.

If the client code is executed, what effects will it have? The answer is not immediately clear from the client's source-code, but a capability-based argument goes as follows: because the client code can typecheck needing only Logger, then whatever effects presented by Logger are an upper-bound on the effects of the client.

The desugaring first creates two functions, MakeLogger and MakeClient, which instantiate the Logger and Client modules; the client code is treated as an implicit module. Lines 1-4 define a function which, given a File, returns a record containing a single

```
resource module Logger
require File

def log(): Unit with File.append =
    File.append('`message logged'')

module Client
require Logger

def run(): Unit =
    Logger.log()

resource module Main
require File
instantiate Logger(File)
instantiate Client(Logger)

Client.run()
```

Figure 4.8: A logger client doesn't need to add effect labels; these can be inferred.

log function. Lines 6-8 define a function which, given a Logger, returns the unannotated client code, wrapped inside an import expression selecting its needed authority. Lines 10-14 are the meat of the program; this function, when given a File capability, creates the modules and then runs the client code. Program execution begins on line 16, where Main is given its initial set of capabilities — which, in this case, is just File.

The interesting part is on lines 7-8, where the unannotated code selects File.append as its authority. This is exactly the effects of the logger, i.e. effects(Unit  $\rightarrow_{\tt File.append}$  Unit) = {File.append}. The code also satisfies the higher-order safety predicates, and the body of the import expression typechecks in the empty context. Therefore, the unannotated code typechecks by  $\varepsilon$ -IMPORT.

In such a small example the client could simply inspect the source code of Logger to determine what effects it might have. Several situations can make this impossible or tedious. First, the manual approach loses efficiency when the system involves many modules of large size across code-ownership boundaries; capability-based reasoning tells you automatically. Second, the source code of Logger might be obfuscated or unavailable, and the only useful information is that given by its signature. Lastly, the client may not care about effects in this situation; the program may be a quick-and-dirty throwaway, in which case it is nice that the capability-based reasoning still accepts the client code without annotations..

```
let MakeLogger =
      (\lambdaf: File.
        \lambda x: Unit. f.append) in
3
   let MakeClient =
5
      (\lambdalogger: Logger.
        import(File.append) logger = logger in
           \lambda x: Unit. logger unit) in
8
   let MakeMain =
10
      (\lambdaf: File.
11
        \lambda x: Unit.
12
           let LoggerModule = MakeLogger f in
13
           let ClientModule = MakeClient LoggerModule in
14
           ClientModule unit) in
15
16
   (MakeMain File) unit
17
```

Figure 4.9: Desugared version of Figure 4.8.

#### 4.2.2 API Violation

Figure 4.10 inverts the roles of the last scenario: now, the annotated Client wants to use the unannotated Logger. The Logger module captures the File capability, and exposes a single function log with the File.append effect. However, the Client has a function run which executes Logger.log, incurring the effect of File.append, but declares its set of effects as  $\varnothing$ . The implementation and the signature of Client.run are inconsistent — does the type system recognise this?

A desugaring is given in Figure 4.11. Lines 1-3 define the function which instantiates the Logger module. Lines 5-8 define the function which instantiates the Client module. Lines 10-15 define the function which instantiates the Main module. Line 17 initiates the program, supplying File to the Main module and invoking its main method. On lines 3-4, the unannotated code is modelled using an import expression which selects  $\varnothing$  as its authority. So far this coheres to the expectations of Client. However,  $\varepsilon$ -IMPORT cannot be applied because the name being bound, f, has the type  $\{\text{File}\}$ , and  $\text{effects}(\{\text{File}\}) = \{\text{File.*}\}$ , which is inconsistent with the declared effects  $\varnothing$ .

The only way for this to typecheck would be to annotate Client.run as having every effect on File. This demonstrates how the effect-system of  $\lambda_{\pi,\varepsilon}^{\rightarrow}$  approximates unlabelled code: it simply considers it as having every effect which could be incurred on those resources in scope, which here is File.\*.

```
resource module Logger
require File

def log(): Unit =
   File.append('`message logged'')

resource module Client
require Logger

def run(): Unit with Ø =
   Logger.log()

resource module Main
require File
instantiate Logger(File)
instantiate Client (Logger)

Client.run()
```

Figure 4.10: The unlabelled code in Logger exercises authority exceeding that selected by Client.

#### 4.2.3 API Violation

Figure 4.12 is a variation of the last example, but now Logger.log is passed the File capability, rather than possessing it. Logger.log still incurs File.append inside unannotated code, which causes the implementation of Client.run to violate its signature.

- 4.2.4 API Violation
- 4.2.5 Hidden Authority
- 4.2.6 Hidden Authority 2
- 4.2.7 Hidden Authority 2
- 4.2.8 Resource Leak

```
let MakeLogger =
      (\lambdaf: File.
        import(\emptyset) f = f in
           \lambda x: Unit. f.append) in
   let MakeClient =
      (\lambdalogger: Logger.
        \lambda x: Unit. logger unit) in
   let MakeMain =
10
      (\lambdaf: File.
11
        let LoggerModule = MakeLogger f in
12
        let ClientModule = MakeClient LoggerModule in
13
        ClientModule unit) in
14
15
   (MakeMain File) unit
16
```

Figure 4.11: Desugared version of Figure 4.10.

```
module Logger

def log(f: {File}): Unit
    f.append('`message logged'')

module Client
instantiate Logger(File)

def run(f: {File}): Unit with Ø
    Logger.log(File)
```

A desugared version is given in Figure 4.12, which is largely the same as in the previous example, except

```
resource module Main
require File
instantiate Client

Client.run(File)
```

Figure 4.12: The unlabelled code in Logger exercises authority exceeding that selected by Client.

```
let MakeClient =
      (\lambda x: Unit.
        let MakeLogger =
3
            (\lambda x: Unit.
              import (\emptyset) x=x in
5
                 \lambda f: \{File\}. f.append) in
         let LoggerModule = MakeLogger unit in
        \lambda f: {File}. LoggerModule f) in
   let MakeMain =
10
      (\lambda f: {File}.
11
        \lambda x: Unit.
12
           let ClientModule = MakeClient unit in
13
           ClientModule f) in
14
   (MakeMain File) unit
16
```

Figure 4.13: Desugared version of 4.12.

```
resource module Logger
require File

def log(): Unit with {File.append, File.write} =
   File.append('`message logged'')
   File.write('`message written'')

module Client

def run(1: Logger): Unit with {File.append} =
   1.log()

resource module Main
require File
instantiate Logger(File)

Client.run(Logger)
```

Figure 4.14: This won't type because of a mismatch between the effects of Client and the effects of Logger.

```
let MakeLogger =
     (\lambdaf: File.
        \lambda x: Unit. let _ = f.append in f.write) in
  let MakeClient =
     (\lambda x: Unit.
        \lambdalogger: Logger unit) in
  let MakeMain =
     (\lambdaf: File.
10
        \lambda x: Unit.
11
          let LoggerModule = MakeLogger f in
12
          let ClientModule = MakeClient unit in
13
           ClientModule.run LoggerModule) in
14
15
   (MakeMain File) unit
16
```

Figure 4.15: Desugared version of Figure 4.14.

```
module Malicious

def stealData(f: {File}):Unit with {File.read} =
    f.read

module Plugin
instantiate Malicious

def run(f: {File}): Unit with Ø =
    Malicious.stealData(f)

resource module Main
require File
instantiate Plugin

Plugin.run(File)
```

Figure 4.16: The Main module transitively invokes a File.read effect, violating its selected authority.

```
let MakePlugin =
      (\lambda x: Unit.
        let MakeMalicious =
           (\lambda x: Unit. \lambda f: \{File\}. f.read) in
        let MaliciousModule = (MakeMalicious unit) in
        \lambda f: {File}. MaliciousModule f) in
   let MakeMain =
      (\lambdaf: File.
        \lambda x: Unit.
10
           let PluginModule = MakePlugin unit in
11
           PluginModule.run f) in
12
13
   (MakeMain File) unit
14
```

Figure 4.17: Desugared version of Figure 4.16.

```
module Malicious

def stealData(f: {File}):Unit =
    f.read

module Plugin
instantiate Malicious

def run(f: {File}): Unit with Ø =
    Malicious.stealData(f)

resource module Main
require File
instantiate Plugin

Plugin.run(File)
```

Figure 4.18: The transitive invocation of File.read now happens inside unannotated code, but the type system will still reject this program.

```
let MakePlugin =
      (\lambda x: Unit.
        let MakeMalicious =
            (\lambda x: Unit.
              import(\emptyset) x=x in
                 \lambda f: \{File\}. f.read) in
        let MaliciousModule = (MakeMalicious unit) in
        \lambda f: {File}. MaliciousModule f) in
   let MakeMain =
10
      (\lambdaf: File.
11
        \lambda x: Unit.
12
           let PluginModule = MakePlugin unit in
13
           PluginModule.run f) in
14
15
   (MakeMain File) unit
16
```

Figure 4.19: Desugared version of Figure 4.18.

```
module Malicious

def log(f: Unit → Unit):Unit
f()

module Plugin
instantiate Malicious

def run(f: {File}): Unit with Ø

Malicious.log(λx:Unit. f.read)

resource module Main
require File
instantiate Plugin

Plugin.run(File)
```

Figure 4.20: The transitive invocation of File.read happens when the unannotated code executes the function given to it.

```
module Malicious

def log(f: Unit → File):Unit
   f().read

module Plugin
instantiate Malicious

def run(f: {File}): Unit with ∅
   Malicious.log(λx:Unit. f)

resource module Main
require File
instantiate Plugin

Plugin.run(File)
```

Figure 4.21: A resource leak allows Malicious to gain access to the File capability directly.

## 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 [8]. 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

## $\lambda_{\pi}^{\rightarrow}$ **Proofs**

**Lemma 10** (Canonical Forms). *Unless the rule used is*  $\varepsilon$ -SUBSUME, *the following are true:* 

```
1. If \Gamma \vdash v : \tau with \varepsilon then \varepsilon = \emptyset.
```

- 2. If  $\Gamma \vdash v : \{\bar{r}\}$  with  $\varepsilon$  then v = r for some  $r \in R$  and  $\{\bar{r}\} = \{r\}$ .
- 3. If  $\Gamma \vdash v : \tau_1 \rightarrow_{\varepsilon'} \tau_2$  with  $\varepsilon$  then  $v = \lambda x : \tau.e.$

Proof.

- 1. A value is either a resource literal or a lambda. The only rules which can type values are  $\varepsilon$ -RESOURCE and  $\varepsilon$ -ABS. In the conclusions of both,  $\varepsilon = \emptyset$ .
- 2. The only rule ascribing the type  $\{\bar{r}\}$  is  $\varepsilon$ -RESOURCE. Its premises imply the result.

3. The only rule ascribing the type  $\tau_1 \to_{\varepsilon'} \tau_2$  is  $\varepsilon$ -ABS. Its premises imply the result.

**Theorem 16** (Progress). *If*  $\Gamma \vdash e : \tau$  with  $\varepsilon$  and e is not a value or variable, then  $e \longrightarrow e' \mid \varepsilon$ .

*Proof.* By induction on  $\Gamma \vdash e : \tau$  with  $\varepsilon$ .

Case:  $\varepsilon$ -VAR,  $\varepsilon$ -RESOURCE, or  $\varepsilon$ -ABS. Then e is a value or variable, and the theorem statement holds vacuously.

Case:  $\varepsilon$ -APP. Then  $e=e_1\ e_2$ . If  $e_1$  is a non-value it can be reduced by inductive assumption, and then  $e_1\ e_2 \longrightarrow e_1'\ e_2$  by E-APP1. If  $e_1=v_1$  is a value and  $e_2$  a non-value, then  $e_2$  can be reduced by inductive assumption and then  $e_1\ e_2 \longrightarrow v_1\ e_2'$  by E-APP2. Otherwise  $e_1=v_1$  and  $e_2=v_2$  are both values. By canonical forms,  $v_1=\lambda x:\tau.e$ . Then  $(\lambda x:\tau.e)v_2 \longrightarrow [v_2/x] \mid \varnothing$  by E-APP3.

Case:  $\varepsilon$ -OPER. Then  $e=e_1.\pi$ . If  $e_1$  is a non-value it can be reduced by inductive assumption, and then  $e_1.\pi \longrightarrow e_1'.\pi \mid \varepsilon_1$  by E-OPERCALL1. Otherwise  $e_1=v_1$  is a value. By

canonical forms,  $v_1 = r$ . Then  $r.\pi \longrightarrow \text{unit} \mid \{r.\pi\}$  by E-OPERCALL2.

Case:  $\varepsilon$ -SUBSUME. Then  $\Gamma \vdash e : \tau'$  with  $\varepsilon'$ . By inversion,  $\Gamma \vdash e : \tau$  with  $\varepsilon$ , where  $\tau' <: \tau$  and  $\varepsilon' \subseteq \varepsilon$ . If e is a value or variable, the theorem holds vacuously. Otherwise the reduction exists by applying the inductive assumption to the sub-derivation.

**Lemma 11** (Substitution). *If*  $\Gamma, x : \tau' \vdash e : \tau$  with  $\varepsilon$  and  $\Gamma \vdash v : \tau'$  with  $\varnothing$  then  $\Gamma \vdash [v/x]e : \tau$  with  $\varepsilon$ .

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

Case:  $\varepsilon$ -VAR. Then e = y and either y = x or  $y \neq x$ .

**Subcase 1:**  $y \neq x$ . By theorem assumption,  $\Gamma \vdash y : \tau$  with  $\varepsilon$ , where  $\varepsilon = \emptyset$  by inversion on  $\varepsilon$ -VAR. Since [v/x]y = y, then  $\Gamma \vdash [v/x]y : \tau$  with  $\emptyset$ .

**Subcase 2:** y=x. By inversion on  $\varepsilon$ -VAR, the typing judgement from the theorem assumption is  $\Gamma, x: \tau' \vdash x: \tau'$  with  $\varnothing$ . Since [v/x]y=v, and by assumption  $\Gamma \vdash v: \tau'$  with  $\varnothing$ , then  $\Gamma \vdash [v/x]x: \tau'$  with  $\varnothing$ .

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

Case:  $\varepsilon$ -APP. By inversion,  $\Gamma, x: \tau' \vdash e_1: \tau_2 \to_{\varepsilon_3} \tau_3$  with  $\varepsilon_A$  and  $\Gamma, x: \tau' \vdash e_2: \tau_2$  with  $\varepsilon_B$ , where  $\varepsilon = \varepsilon_A \cup \varepsilon_B \cup \varepsilon_3$  and  $\tau = \tau_3$ . By inductive assumption,  $\Gamma \vdash [v/x]e_1: \tau_2 \to_{\varepsilon_3} \tau_3$  with  $\varepsilon_A$  and  $\Gamma \vdash [v/x]e_2: \tau_2$  with  $\varepsilon_B$ . By  $\varepsilon$ -APP we have  $\Gamma \vdash ([v/x]e_1)([v/x]e_2): \tau_3$  with  $\varepsilon_A \cup \varepsilon_B \cup \varepsilon_3$ . By simplifying and applying the definition of substitution, this is the same as  $\Gamma \vdash [v/x](e_1e_2): \tau$  with  $\varepsilon$ .

Case:  $\varepsilon$ -OPERCALL. By inversion,  $\Gamma, x: \tau' \vdash e_1: \{\bar{r}\}$  with  $\varepsilon_1$ , where  $\tau = \{\bar{r}\}$  and  $\varepsilon = \varepsilon_1 \cup \{r.\pi \mid r \in \bar{r}, \pi \in \Pi\}$ . By applying the inductive assumption,  $\Gamma \vdash [v/x]e_1: \{\bar{r}\}$  with  $\varepsilon_1$ . Then by  $\varepsilon$ -OPERCALL,  $\Gamma \vdash ([v/x]e_1).\pi: \{\bar{r}\}$  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  $\Gamma \vdash [v/x](e_1.\pi): \tau$  with  $\varepsilon$ .

Case:  $\varepsilon$ -SUBSUME. By inversion,  $\Gamma, x : \tau' \vdash e : \tau_2$  with  $\varepsilon_2$ , where  $\tau_2 <: \tau$  and  $\varepsilon_2 \subseteq \varepsilon$ . By inductive hypothesis,  $\Gamma \vdash [v/x]e : \tau_2$  with  $\varepsilon_2$ . Then  $\Gamma \vdash [v/x]e : \tau$  with  $\varepsilon$  by  $\varepsilon$ -SUBSUME.

**Theorem 17** (Preservation). *If*  $\Gamma \vdash e_A : \tau_A$  with  $\varepsilon_A$  and  $e_A \longrightarrow e_B \mid \varepsilon_C$ , then  $\Gamma \vdash e_B : \tau_B$  with  $\varepsilon_B$ , where  $e_B <: e_A$  and  $\varepsilon \cup \varepsilon_B \subseteq \varepsilon_A$ .

*Proof.* By induction on  $\Gamma \vdash e_A : \tau_A$  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  $e_A=e_1\ e_2$  and  $e_1:\tau_2\longrightarrow_\varepsilon \tau_3$  with  $\varepsilon_1$  and  $\Gamma\vdash e_2:\tau_2$  with  $\varepsilon_2$  and  $\tau_B=\tau_3$  and  $\varepsilon_A=\varepsilon_1\cup\varepsilon_2\cup\varepsilon$ .

**Subcase:** E-APP1. Then  $e_1 e_2 \longrightarrow e'_1 e_2 \mid \varepsilon_C$ . From inversion on E-APP1,  $e_1 \longrightarrow e'_1 \mid \varepsilon_C$ . By inductive hypothesis and application of  $\varepsilon$ -SUBSUME,  $\Gamma \vdash e'_1 : \tau_2 \longrightarrow_{\varepsilon} \tau_3$  with  $\varepsilon_1$ . Then  $\Gamma \vdash e'_1 e_2 : \tau_3$  with  $\varepsilon_1 \cup \varepsilon_2 \cup \varepsilon$  by  $\varepsilon$ -APP.

**Subcase:** E-APP2. Then  $e_1=v_1$  is a value and  $e_2\longrightarrow e_2'\mid \varepsilon_C$ . By inversion on E-APP2,  $e_2\longrightarrow e_2'\mid \varepsilon_C$ . By inductive hypothesis and application of  $\varepsilon$ -SUBSUME,  $\Gamma\vdash e_2':\tau_2$  with  $\varepsilon_2$ . Then  $\Gamma\vdash v_1\ e_2':\tau_3$  with  $\varepsilon_1\cup\varepsilon_2\cup\varepsilon$  by  $\varepsilon$ -APP.

**Subcase:** E-APP3. Then  $e_1 = \lambda x : \tau_2.e$  and  $e_2 = v_2$  are values, and  $(\lambda x : \tau_2.e) \ v_2 \longrightarrow [v_2/x]e \mid \varnothing$ . By inversion on the typing rule for  $\lambda x : \tau_2.e$ , we know  $\Gamma, x : \tau_2 \vdash e : \tau_3$  with  $\varepsilon_3$ . By canonical forms,  $\varepsilon_2 = \varnothing$  because  $e_2 = v_2$  is a value. Then by the substitution lemma,  $\Gamma \vdash [v_2/x]e : \tau_3$  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. Then  $e_A=e_1.\pi$  and  $e_1:\{\bar{r}\}$  with  $\varepsilon_1$  and  $\tau_A=$  Unit and  $\varepsilon_A=\varepsilon_1\cup\{r.\pi\mid r\in\bar{r},\pi\in\Pi\}$ .

**Subcase:** E-OPERCALL1. Then  $e_1.\pi \longrightarrow e'_1.\pi \mid \varepsilon_C$ . By inversion on E-OPERCALL1,  $e_1 \longrightarrow e'_1 \mid \varepsilon_C$ . By inductive hypothesis and application of  $\varepsilon$ -Subsume,  $\Gamma \vdash e_1 : \{\bar{r}\}$  with  $\varepsilon_1$ . Then  $\Gamma \vdash e'_1.\pi : \{\bar{r}\}$  with  $\varepsilon_1 \cup \{r.\pi \mid r \in \bar{r}, \pi \in \Pi\}$ .

**Subcase:** E-OPERCALL2. Then  $e_1 = r$  is a resource literal and  $r.\pi \longrightarrow \text{unit} \mid \{r.\pi\}$ . By canonical forms,  $\Gamma \vdash r$ : unit with  $\{r.\pi\}$ . Trivially,  $\Gamma \vdash \text{unit}$ : Unit with  $\varnothing$ . Therefore  $\tau_B = \tau_A$  and  $\varepsilon_C \cup \varepsilon_B = \{r.\pi\} = \varepsilon_A$ .

**Theorem 18** (Soundness). *If*  $\Gamma \vdash e_A : \tau_A$  with  $\varepsilon_A$  and  $e_A$  is not a value, then  $e_A \longrightarrow e_B \mid \varepsilon$ , where  $\Gamma \vdash e_B : \tau_B$  with  $\varepsilon_B$  and  $\tau_B <: \tau_A$  and  $\varepsilon_B \cup \varepsilon \subseteq \varepsilon_A$ .

*Proof.* If  $e_A$  is not a value or variable then the reduction exists by the progress theorem. The rest of the theorem statement follows by the preservation theorem.

**Theorem 19** (Multi-step Soundness). *If*  $\Gamma \vdash e_A : \tau_A \text{ with } \varepsilon_A \text{ and } e_A \longrightarrow^* e_B \mid \varepsilon$ , where  $\Gamma \vdash e_B : \tau_B \text{ with } \varepsilon_B \text{ and } \tau_B <: \tau_A \text{ and } \varepsilon_B \cup \varepsilon \subseteq \varepsilon_A$ .

*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  $e_A \longrightarrow^* e_C \mid \varepsilon'$ , 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,  $\Gamma \vdash e_C : \tau_C$  with  $\varepsilon_C$  and  $\Gamma \vdash e_B : \tau_B$  with  $\varepsilon_B$ . By inductive assumption,  $\tau_C <: \tau_A$  and  $\varepsilon_C \cup \varepsilon' \subseteq \varepsilon_A$ . By single-step soundness,  $\tau_B <: \tau_C$  and  $\varepsilon_B \cup \varepsilon'' \subseteq \varepsilon_C$ . Then by transitivity,  $\tau_B <: \tau$  and  $\varepsilon_B \cup \varepsilon' \cup \varepsilon'' = \varepsilon_B \cup \varepsilon \subseteq \varepsilon_A$ .

## Appendix B

# $\lambda_{\pi,\varepsilon}^{\rightarrow}$ Proofs

**Theorem 20** (Progress). *If*  $\hat{\Gamma} \vdash \hat{e}_A : \hat{\tau}$  with  $\varepsilon$  and  $\hat{e}_A$  is not a value or variable, then  $\hat{e}_A \longrightarrow \hat{e}_B \mid \varepsilon$ .

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

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

**Lemma 12** (Substitution). *If*  $\hat{\Gamma}, x : \hat{\tau}' \vdash \hat{e}_A : \hat{\tau}_A \text{ with } \varepsilon_A \text{ and } \hat{\Gamma} \vdash \hat{v} : \hat{\tau}' \text{ with } \varnothing \text{ then } \hat{\Gamma} \vdash [\hat{v}/x]e : \hat{\tau}_A \text{ with } \varepsilon_A.$ 

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

*Case:*  $\varepsilon$ -MODULE. Then the following are true.

- 1.  $\hat{\Gamma}, x : \hat{\tau}' \vdash \mathtt{import}(\varepsilon) \ x = \hat{e} \ \mathtt{in} \ e : \mathtt{annot}(\tau, \varepsilon) \ \mathtt{with} \ \varepsilon \cup \varepsilon_1$
- 2.  $\hat{\Gamma}, x : \hat{\tau}' \vdash \hat{e} : \hat{\tau} \text{ with } \varepsilon_1$
- 3.  $\varepsilon = \text{effects}(\hat{\tau}) \cup \text{ho-effects}(\text{annot}(\tau, \varepsilon))$
- 4.  $x : erase(\hat{\tau}) \vdash e : \tau$
- 5.  $\hat{\tau}_A = \mathtt{annot}(\tau, \varepsilon)$
- 6.  $\hat{\varepsilon}_A = \varepsilon \cup \varepsilon_1$

By applying inductive assumption to (2),  $\hat{\Gamma} \vdash [\hat{v}/x]\hat{e} : \hat{\tau}$  with  $\varepsilon_1$ . Then by  $\varepsilon$ -Module,  $\hat{\Gamma} \vdash \text{import}(\varepsilon) \ x = [\hat{v}/x]\hat{e} \ \text{in} \ e : \text{annot}(\tau, \varepsilon) \ \text{with} \ \varepsilon \cup \varepsilon_1$ .

**Lemma 13.** If  $effects(\hat{\tau}) \subseteq \varepsilon$  and  $ho-safe(\hat{\tau}, \varepsilon)$  then  $\hat{\tau} <: annot(erase(\hat{\tau}), \varepsilon)$ .

**Lemma 14.** 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 derivations of safe and ho-safe.

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 effects( $\hat{\tau}$ )  $\subseteq \varepsilon$  we have ho-effects( $\hat{\tau}_1$ )  $\cup \varepsilon' \cup$  effects( $\hat{\tau}_2$ )  $\subseteq \varepsilon$  and therefore effects( $\hat{\tau}_2$ )  $\subseteq \varepsilon$ . From ho-safe( $\hat{\tau}, \varepsilon$ ) we have 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 15** (Annotation). *If the following are true:* 

```
 \begin{array}{l} 1. \ \hat{\Gamma} \vdash \hat{v} : \hat{\tau} \ \text{with} \varnothing \\ 2. \ \Gamma, y : \texttt{erase}(\hat{\tau}) \vdash e : \tau \\ 3. \ \varepsilon = \texttt{effects}(\hat{\tau}) \cup \texttt{ho-effects}(\texttt{annot}(\tau, \varepsilon)) \\ 4. \ \texttt{ho-safe}(\hat{\tau}, \varepsilon) \end{array}
```

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

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

*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 assumption,  $y: \text{erase}(\hat{\tau}) \vdash x:\tau$ , so  $\tau=\text{erase}(\hat{\tau})$ . By  $\varepsilon\text{-VAR}$ ,  $y:\hat{\tau} \vdash x:\hat{\tau}$  with  $\varnothing$ . By definition,  $\text{annot}(x,\varepsilon')=x$ , so  $y:\hat{\tau} \vdash \text{annot}(x,\varepsilon'):$ 

 $\hat{\tau}$  with  $\varnothing$ . By assumptions (3) and (4) we know effects $(\hat{\tau}) \subseteq \varepsilon$  and ho-safe $(\hat{\tau}, \varepsilon)$ , so applying Lemma 2 gives  $\hat{\tau} <: \operatorname{annot}(\operatorname{erase}(\hat{\tau}), \varepsilon)$ . Then applying  $\varepsilon$ -SUBSUME to the type of  $y: \hat{\tau} \vdash \operatorname{annot}(x, \varepsilon): \hat{\tau}$  with  $\varnothing$  gives  $y: \hat{\tau} \vdash \operatorname{annot}(x, \varepsilon): \operatorname{annot}(\operatorname{erase}(\hat{\tau}), \varepsilon)$  with  $\varnothing$ . We already know  $\operatorname{erase}(\hat{\tau}) = \tau$ , so this judgement is the same as  $y: \hat{\tau} \vdash \operatorname{annot}(x, \varepsilon): \operatorname{annot}(x, \varepsilon):$ 

Subcase 2:  $x \neq y$ . Because  $\Gamma, y : \operatorname{erase}(\hat{\tau}) \vdash x : \tau$ , and  $x \neq y$ , then  $x : \tau \in \Gamma$ . Then  $x : \operatorname{annot}(\tau, \varepsilon) \in \operatorname{annot}(\Gamma, \varepsilon)$ , so  $\operatorname{annot}(\Gamma, \varepsilon) \vdash x : \operatorname{annot}(\tau, \varepsilon)$  with  $\varnothing$  by  $\varepsilon$ -VAR. By definition,  $\operatorname{annot}(x, \varepsilon) = x$ , so  $\operatorname{annot}(\Gamma, \varepsilon) \vdash \operatorname{annot}(x, \varepsilon) : \operatorname{annot}(\tau, \varepsilon)$  with  $\varnothing$ . Applying  $\varepsilon$ -SUBSUME gives  $\operatorname{annot}(\Gamma, \varepsilon) \vdash \operatorname{annot}(x, \varepsilon) : \operatorname{annot}(\tau, \varepsilon)$  with  $\varepsilon \cup \operatorname{effects}(\operatorname{annot}(\Gamma, \varepsilon))$ . Lastly, by widening the context,  $\hat{\Gamma}$ ,  $\operatorname{annot}(\Gamma, \varepsilon)$ ,  $y : \hat{\tau} \vdash \operatorname{annot}(\tau, \varepsilon)$  with  $\varepsilon \cup \operatorname{effects}(\operatorname{annot}(\Gamma, \varepsilon))$ .

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

Case: T-ABS. Then  $\Gamma, y: \operatorname{erase}(\hat{\tau}) \vdash \lambda x: \tau_1.e_{body}: \tau_1 \to \tau_2$ . Applying definitions,  $\operatorname{annot}(e, \varepsilon') = \operatorname{annot}(\lambda x: \tau_1.e_2, \varepsilon') = \lambda x: \operatorname{annot}(\tau_1, \varepsilon').\operatorname{annot}(e_2, \varepsilon')$  and  $\operatorname{annot}(\tau, \varepsilon') = \operatorname{annot}(\tau_1 \to \tau_2, \varepsilon') = \operatorname{annot}(\tau_1, \varepsilon') \to_{\varepsilon'} \operatorname{annot}(\tau_2, \varepsilon')$ . By inversion on T-ABS, we get the sub-derivation  $\Gamma, y: \operatorname{erase}(\hat{\tau}), x: \tau_1 \vdash e_{body}: \tau_2$ . We shall apply the inductive assumption to this judgement with an unlabelled context  $\Gamma, x: \tau_1$ . Define  $\varepsilon'' = \varepsilon' \cup \operatorname{effects}(\operatorname{annot}(\tau_1, \varepsilon))$ .  $\varepsilon''$  is the inductive assumption's version of  $\varepsilon'$ . Applying the inductive assumption,  $\hat{\Gamma}$ ,  $\operatorname{annot}(\Gamma, \varepsilon''), y: \hat{\tau}, x: \operatorname{annot}(\tau_1, \varepsilon'') \vdash \operatorname{annot}(e_{body}, \varepsilon''): \operatorname{annot}(\tau_2, \varepsilon'')$  with  $\varepsilon''$ .

Ostensibly,  $\varepsilon''$  seems to differ from  $\varepsilon'$  by effects  $(\operatorname{annot}(\tau_1, \varepsilon'))$ , but we shall show  $\varepsilon' = \varepsilon''$ . By assumption (3),  $\varepsilon' = \operatorname{effects}(\hat{\tau}) \cup \operatorname{ho-effects}(\operatorname{annot}(\tau, \varepsilon))$ . In this case,  $\tau = \tau_1 \to \tau_2$ , so  $\operatorname{annot}(\tau, \varepsilon) = \operatorname{annot}(\tau_1, \varepsilon) \to_{\varepsilon} \operatorname{annot}(\tau_2, \varepsilon)$ . By definition, effects  $(\operatorname{annot}(\tau_1, \varepsilon)) \subseteq \operatorname{ho-effects}(\operatorname{annot}(\tau, \varepsilon))$ , so effects  $(\operatorname{annot}(\tau_1, \varepsilon)) \subseteq \varepsilon$ . Therefore,  $\varepsilon'' = \varepsilon \cup \operatorname{effects}(\operatorname{annot}(\Gamma, \varepsilon))$ . But this is the definition of  $\varepsilon'$ , so  $\varepsilon'' = \varepsilon'$ .

We can therefore rewrite the judgement obtained from the inductive hypothesis as  $\hat{\Gamma}$ , annot $(\Gamma, \varepsilon'), y: \hat{\tau}, x: \mathrm{annot}(\tau_1, \varepsilon') \vdash \mathrm{annot}(e_{body}, \varepsilon'): \mathrm{annot}(\tau_2, \varepsilon')$  with  $\varepsilon'$ . Applying  $\varepsilon$ -ABS,  $\hat{\Gamma}$ , annot $(\Gamma, \varepsilon'), y: \hat{\tau} \vdash \lambda x: \mathrm{annot}(\tau_1, \varepsilon').\mathrm{annot}(e_{body}, \varepsilon')$  with  $\varnothing$ . Lastly, by  $\varepsilon$ -Subsume,  $\hat{\Gamma}$ , annot $(\Gamma, \varepsilon), y: \hat{\tau} \vdash \mathrm{annot}(e, \varepsilon): \mathrm{annot}(\hat{\tau}_1) \rightarrow_{\varepsilon} \mathrm{annot}(\hat{\tau}_2)$  with  $\varepsilon$ -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  $\Gamma, y: \operatorname{erase}(\hat{\tau}) \vdash e_1.\pi: \operatorname{Unit}$ . By inversion we get the subderivation  $\Gamma, y: \operatorname{erase}(\hat{\tau}) \vdash e_1: \{\bar{r}\}$ . By definition,  $\operatorname{annot}(\{\bar{r}\}, \varepsilon) = \{\bar{r}\}$ . By inductive assumption,  $\hat{\Gamma}$ ,  $\operatorname{annot}(\Gamma, \varepsilon), y: \hat{\tau} \vdash e_1: \{\bar{r}\}$  with  $\varepsilon \cup \operatorname{effects}(\operatorname{annot}(\Gamma, \varepsilon))$ . By  $\varepsilon$ -OPERCALL,  $\hat{\Gamma}$ ,  $\operatorname{annot}(\Gamma, \varepsilon), y: \hat{\tau} \vdash e_1.\pi: \{\bar{r}\}$  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 \text{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  $r=\operatorname{erase}(\hat{\tau})=\hat{\tau}$  and one of the above subcases must also hold.

**Theorem 21** (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:  $\varepsilon$ -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.  $\Gamma \vdash e_A : \mathtt{annot}(\tau, \varepsilon) \mathtt{ with } \varepsilon \cup \varepsilon_1$
- 3.  $import(\varepsilon) x = \hat{v} in e \longrightarrow [\hat{v}/x] annot(e, \varepsilon) \mid \varnothing$
- 4.  $\hat{\Gamma} \vdash \hat{v} : \hat{\tau} \text{ with } \varnothing$
- 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 22** (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  $\varepsilon_B \cup \varepsilon \subseteq \varepsilon_A$ .

**Theorem 23** (Multi-step Soundness). If  $\hat{\Gamma} \vdash \hat{e}_A : \hat{\tau}_A$  with  $\varepsilon_A$  and  $\hat{e}_A \longrightarrow^* e_B \mid \varepsilon$ , 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.* The proofs are the same as for  $\lambda_{\pi}^{\rightarrow}$ .

## 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] KLEENE, S. Recursive predicates and quantifiers. *Journal of Symbolic Logic 8*, 1 (1943), 32–34.
- [6] 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.
- [7] 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.
- [8] LIU, F. A study of capability-based effect systems. Master's thesis, École Polytechnique Fédérale de Lausanne, 2016.
- [9] MAFFEIS, S., MITCHELL, J. C., AND TALY, A. Object Capabilities and Isolation of Untrusted Web Applications. In *IEEE Symposium on Security and Privacy* (2010).
- [10] MILLER, M., YEE, K.-P., AND SHAPIRO, J. Capability myths demolished. Tech. rep., 2003.
- [11] MILLER, M. S. Robust Composition: Towards a Unified Approach to Access Control and Concurrency Control. PhD thesis, Johns Hopkins University, 2006.
- [12] NIELSON, F., AND NELSON, H. R. Type and Effect Systems. pp. 114–136.

68 BIBLIOGRAPHY

[13] ODERSKY, M., ALTHERR, P., CREMET, V., DUBOCHET, G., EMIR, B., HALLER, P., MICHELOUD, S., MIHAYLOV, N., MOORS, A., RYTZ, L., SCHINZ, M., STENMAN, E., AND ZENGER, M. Scala Language Specification. http://scala-lang.org/files/archive/spec/2.11/. Last accessed: Nov 2016.

- [14] PIERCE, B. C. *Types and Programming Languages*. The MIT Press, Cambridge, MA, USA, 2002.
- [15] SALTZER, J. H. Protection and the Control of Information Sharing in Multics. *Communications of the ACM 17*, 7 (1974), 388–402.