

# A Study of Capability-Based Effect Systems

Fengyun Liu

*A thesis submitted for the degree of Master of Computer Science  
at École Polytechnique Fédérale de Lausanne*

*Supervisors*

Martin Odersky  
Professor

Nada Amin  
PhD Student

Sandro Stucki  
PhD Student

School of Computer and Communication Sciences



ÉCOLE POLYTECHNIQUE  
FÉDÉRALE DE LAUSANNE

Lausanne, January 2016

# Acknowledgments

I thank Professor Martin Odersky for granting me the opportunity to work on the master project to explore capability-based effect systems.

I thank Sandro Stucki and Nada Amin for their mentoring. Without their help, I could not have made such a big progress from a novice in type systems and Coq in such a short time.

Finally, I'd like to thank my wife Yinai for her long-term support and love. A special thanks goes to my little daughter Tianyao – playing with her is one of the greatest joy in my spare time.

## Abstract

The problem of *effect polymorphism* is a major obstacle to wide adoption of effect systems in the programming community. The absence of effect systems reduces compiler optimization opportunities and disables effect constraints on APIs in parallel and distributed computations.

This study shows that *capability-based* effect systems, equipped with *stoic functions* and *free functions*, can easily solve the problem of *effect polymorphism* without incurring notational burden on programmers. With this advantage, capability-based effect systems stand a better chance to be adopted by the programming community.

The central idea of *capability-based* effect system is that a capability is required in order to produce side effects. If capabilities are passed as function parameters, by tracking capabilities in the type system we can track effects in the program.

To ensure that capabilities are passed through function parameters, instead of being captured from the environment, we need to impose a *variable-capturing discipline*, stipulating that capability variables cannot be captured. Functions observing the discipline are called *stoic functions*, while functions not observing the discipline are called *free functions*.

# Contents

<b>Acknowledgments</b>	<b>1</b>
<b>Abstract</b>	<b>2</b>
<b>1 Introduction</b>	<b>6</b>
1.1 Effect Polymorphism . . . . .	6
1.2 Capability-Based Effect Systems . . . . .	7
1.3 Contributions . . . . .	8
1.4 Structure of Report . . . . .	9
<b>2 System STLC-Pure</b>	<b>10</b>
2.1 Definitions . . . . .	10
2.1.1 Variable-Capturing Discipline . . . . .	10
2.1.2 Stoic Functions and Free Functions . . . . .	12
2.1.3 Where do Effects Come From . . . . .	13
2.1.4 Where do Capabilities Come From . . . . .	13
2.1.5 Why Treat Variables as Values . . . . .	13
2.1.6 What If a Function Has More than One Side Effect . . . . .	13
2.2 Soundness . . . . .	14
2.3 Effect Safety . . . . .	15
2.3.1 Informal Formulation . . . . .	15
2.3.2 Inhabited Types and Environments . . . . .	16
2.3.3 Formalization . . . . .	17
2.3.4 Proof . . . . .	19
2.3.5 An Intuitive Proof . . . . .	19
<b>3 System STLC-Impure</b>	<b>22</b>
3.1 Definitions . . . . .	22
3.2 Soundness . . . . .	26
3.3 Effect Safety . . . . .	26
3.3.1 Formulation . . . . .	26
3.3.2 Proof . . . . .	29
3.4 Effect Polymorphism . . . . .	31
3.4.1 Axiomatic Polymorphism . . . . .	31
3.4.2 Curryng Polymorphism . . . . .	32

3.4.3	Stoic Polymorphism . . . . .	33
3.4.4	Discussion . . . . .	33
<b>4</b>	<b>System F-Pure</b>	<b>36</b>
4.1	Definitions . . . . .	36
4.2	Soundness . . . . .	38
4.3	Effect Safety . . . . .	38
4.3.1	Inhabited Types and Environments . . . . .	39
4.3.2	Formulation . . . . .	40
4.3.3	Proof . . . . .	40
<b>5</b>	<b>System F-Impure</b>	<b>43</b>
5.1	Definitions . . . . .	43
5.2	Soundness . . . . .	45
5.3	Effect Safety . . . . .	45
5.3.1	Formulation . . . . .	45
5.3.2	Proof . . . . .	47
<b>6</b>	<b>Conclusion</b>	<b>49</b>
6.1	Related Work . . . . .	49
6.2	Future Work . . . . .	50
	<b>Bibliography</b>	<b>51</b>

## List of Figures

2.1	System STLC-Pure . . . . .	11
2.2	System STLC-Pure Capsafe Environment . . . . .	18
3.1	System STLC-Impure First Formulation . . . . .	23
3.2	Subtyping: Top–Pure and Top–Impure . . . . .	24
3.3	System STLC-Impure . . . . .	25
3.4	System STLC-Impure Capsafe Environment . . . . .	28
3.5	System STLC-Impure Axioms . . . . .	30
4.1	System F-Pure . . . . .	37
4.2	System F-Pure Inhabited Environment . . . . .	39
4.3	System F-Pure Capsafe Environment . . . . .	41
5.1	System F-Impure . . . . .	44
5.2	System F-Impure Capsafe Environment . . . . .	46
5.3	System F-Impure Axioms . . . . .	48

# 1 Introduction

The main motivation of this study is to turn Scala into an effect-disciplined programming language. Currently, Scala doesn't have the ability to track effects in the type system.

The absence of an effect system disables programmers from imposing effect constraints on APIs in parallel and distributed computations. For example, in parallel computing, the functions passed to the function `pmap` should have no side effects, in order to enable parallel processing on lists; however, it is currently impossible to do so in Scala:

```
def pmap(xs: List[Int], f: Int => Int): List[Int]
```

An effect system can differentiate pure code from impure code in a program, thus enabling optimization of pure code. For example, independent pure expressions can be executed in parallel; pure code that is dead can be safely eliminated; common pure expressions can be reduced to one computation; pure operations on data structures can be fused[CLS07]. The absence of an effect system makes these optimizations impossible.

As pointed out by Lukas et al.[ROH12], a main difficulty in introducing a practical effect system is how to handle the problem of *effect polymorphism*.

## 1.1 Effect Polymorphism

*Effect polymorphism* can be illustrated by the function `map`:

```
def map[A,B](f: A => B)(l: List[A]) = l match {  
  case Nil => Nil  
  case x::xs => f(x)::map(f)(xs)  
}
```

In an effect system, the effect of `map` depends on the passed in function `f`. If `f` has IO effects, then `map` also has IO effects. If `f` is pure, then `map` is pure as well. The way to express effect polymorphism in classical type-and-effect systems is to introduce a new type variable `E` to denote the generic effect:

```
def map[A, B, E](f: A => B @E)(l: List[A]): List[B] @E
```

Java, the only industrial language integrated with an effect system for checked exceptions, exemplifies the classical approach to effect polymorphism:

```

public interface FunctionE<T, U, E extends Exception> {
    public U apply(T t) throws E;
}
public interface List<T> {
    public <U, E extends Exception> List<U>
        mapE(FunctionE<T, U, E> f) throws E;
}

```

Due to the verbosity of syntax, effect polymorphism is rarely used in Java, which reduces the effectiveness of the effect system.

The problem of effect polymorphism is how to express effect-polymorphic functions with minimal syntactical overhead. *Capability-based* effect systems provide an elegant solution to this problem.

## 1.2 Capability-Based Effect Systems

The central idea of *capability-based* effect system is that a capability is required in order to produce side effects. If capabilities are passed as function parameters, by tracking capabilities in the type system we can track effects in the program.

To ensure that capabilities are passed through function parameters, instead of being captured from the environment, we need to impose a *variable-capturing discipline*, stipulating that capability variables cannot be captured. Functions observe the discipline are called *stoic functions*, while functions don't observe the discipline are called *free functions*. We use  $\rightarrow$  to denote the type of stoic functions and  $\Rightarrow$  to denote the type of free functions. The following example shows that incorrect capturing of capability variables in a stoic function will generate a typing error:

```

def map(xs: List[Int], f: Int => Int): List[Int]
def pmap(xs: List[Int], f: Int -> Int): List[Int]
def print(x: Any, c: IO): ()

def bar(xs: List[Int])(c: IO) =
    map(xs, { x => print(x, c); x })

def foo(xs: List[Int], c: IO) =
    pmap(xs, { x => print(x, c); x }) // Error, stoic function required

```



In the code above, the function `map` can take either a stoic function or a free function, while the function `pmap` can only take a stoic function. The type signature  $\text{Int} \rightarrow \text{Int}$  in `pmap` ensures that functions passed to `pmap` must be pure. That’s why there is a typing error in the function `foo`, as the passed anonymous function is free instead of stoic.

With the combination of free functions and stoic functions, capability-based effect systems can solve the problem of *effect polymorphism* easily, while incurring no syntactical overhead. For example, the effect-polymorphic function `map` can be implemented as follows:

```
def map[A,B](f: A => B)(l: List[A]) = l match {
  case Nil => Nil
  case x::xs => f(x)::map(f)(xs)
}
def squareImpure(c: IO) = map { x => println(x)(c); x*x }
def squarePure(l: List[Int]) = map { x => x*x } l
```

In the code above, the function `map` is exactly the same as it can be written in Scala now, which would be typed as  $(A \Rightarrow B) \rightarrow \text{List}[A] \Rightarrow \text{List}[B]$ . The function `squareImpure` is stoic, which would be typed as  $\text{IO} \rightarrow \text{List}[\text{Int}] \Rightarrow \text{List}[\text{Int}]$ . The function `squarePure` is also stoic, which would be typed as  $\text{List}[\text{Int}] \rightarrow \text{List}[\text{Int}]$ . By just checking the type signature, we can conclude that `squareImpure` can only have IO effects and `squarePure` is pure. No annotation is required, effect polymorphism just works naturally. We’ll see details of effect polymorphism in the chapter STLC-Impure (Section 3.4).

### 1.3 Contributions

The main contributions of this study are as follows:

- We formulate and prove soundness and effect safety of four capability-based effect systems, which can serve as the foundation for implementing capability-based effect system in functional programming languages. The formalization is done in Coq based on the locally-nameless representation[Cha11] and hosted on Github<sup>1</sup>.
- We propose an approach to solve the problem of *effect polymorphism* with no syntactical overhead, thanks to free functions and stoic functions (Section 3.4).

<sup>1</sup><https://github.com/liufengyun/stoic>

## 1.4 Structure of Report

In following chapters, we'll introduce four systems of increasing complexity, namely STLC-Pure, STLC-Impure, F-Pure and F-Impure. The latter three are a gradual enrichment of STLC-Pure with free functions, subtyping and universal types.

- STLC-Pure is a variant of *simply typed lambda calculus* with only stoic functions.
- STLC-Impure is an extension of STLC-Pure with free functions and subtyping.
- F-Pure is an extension of STLC-Pure with universal types.
- F-Impure is an extension of F-Pure with free functions (without subtyping).

We discuss *effect polymorphism* in the chapter STLC-Impure (Section 3.4).

## 2 System STLC-Pure

This chapter describes a variant of the *simply typed lambda calculus* with the extension of capabilities. We call this system *STLC-Pure*, because in this system all functions must observe a variable-capturing discipline.

The system STLC-Pure, though conceptually simple, can quite well demonstrate the main features of capability-based effect systems. We'll first introduce the formalization, then discuss soundness and effect safety. Concepts introduced here will be a foundation for more complex systems in later chapters.

### 2.1 Definitions

Formally, STLC-Pure is obtained by introducing a capability type and imposing a variable-capturing discipline on lambda abstractions. Figure 2.1 presents the full definition of STLC-Pure.

The syntax is almost the same as standard STLC, except the addition of the capability type  $E$  and the taking of variables as values. The evaluation rules are exactly the same, with standard call-by-value small-step semantics. The typing rule T-ABS is slightly changed by performing an operation *pure* on the environment. The peculiarities in the formalization are explained below.

#### 2.1.1 Variable-Capturing Discipline

The most important change to the standard STLC lies in the following typing rule:

$$\frac{\text{pure}(\Gamma), x : S \vdash t_2 : T}{\Gamma \vdash \lambda x : S. t_2 : S \rightarrow T} \quad (\text{T-ABS})$$

This typing rule imposes a *variable-capturing discipline* on lambda abstractions. This discipline stipulates that only variables whose type is not a capability type can be captured in a lambda abstraction.

The discipline is implemented with the helper function *pure*, which removes all variable bindings of the capability type  $E$  from the typing environment. It's easy to verify that the function *pure* satisfies following properties:

<p><b>Syntax</b></p> <p><math>t ::=</math></p> <ul style="list-style-type: none"> <li><math>x</math>                      terms: variable</li> <li><math>\lambda x:T. t</math>            abstraction</li> <li><math>tt</math>                     application</li> </ul> <p><math>v ::=</math></p> <ul style="list-style-type: none"> <li><math>\lambda x:T. t</math>            values: abstraction value</li> <li><math>x</math>                     variable value</li> </ul> <p><math>T ::=</math></p> <ul style="list-style-type: none"> <li><math>B</math>                     types: basic type</li> <li><math>E</math>                     capability type</li> <li><math>T \rightarrow T</math>        type of functions</li> </ul> <p><b>Evaluation</b>                      <math>t \longrightarrow t'</math></p> $\frac{t_1 \longrightarrow t'_1}{t_1 t_2 \longrightarrow t'_1 t_2} \quad (\text{E-APP1})$ $\frac{t_2 \longrightarrow t'_2}{v_1 t_2 \longrightarrow v_1 t'_2} \quad (\text{E-APP2})$ $(\lambda x:T. t_1)v_2 \longrightarrow [x \mapsto v_2]t_1 \quad (\text{E-APPABS})$	<p><b>Typing</b>                              <math>\boxed{\Gamma \vdash x : T}</math></p> $\frac{x : T \in \Gamma}{\Gamma \vdash x : T} \quad (\text{T-VAR})$ $\frac{\text{pure}(\Gamma), x : S \vdash t_2 : T}{\Gamma \vdash \lambda x:S. t_2 : S \rightarrow T} \quad (\text{T-ABS})$ $\frac{\Gamma \vdash t_1 : S \rightarrow T \quad \Gamma \vdash t_2 : S}{\Gamma \vdash t_1 t_2 : T} \quad (\text{T-APP})$ <p><b>Pure Environment</b></p> $\begin{aligned} \text{pure}(\emptyset) &= \emptyset \\ \text{pure}(\Gamma, x : E) &= \text{pure}(\Gamma) \\ \text{pure}(\Gamma, x : T) &= \text{pure}(\Gamma), x : T \end{aligned}$
--	---

Figure 2.1: System STLC-Pure

LEMMA (Pure-Distributivity).  $\text{pure}(\Gamma, \Delta) = \text{pure} \Gamma, \text{pure} \Delta$

LEMMA (Pure-Idempotency).  $\text{pure}(\text{pure} \Gamma) = \text{pure} \Gamma$

Initially, we've tried a variable-capturing discipline where no free variables can be captured, i.e.  $\text{pure}(\Gamma) = \emptyset$ . While this definition is certainly effect safe, the system is not very expressive, as we lose the ability to create closures in the system, which is usually considered to be an essential feature of functional programming.

In addition to the capability type E, we also tried to exclude variables of the type  $T \rightarrow E$  (for any T) from the pure environment. As an over-approximation, this version is certainly effect-safe. However, it rejects more types than necessary, thus reduces expressiveness of the system.

We introduce two useful definitions with the help of the function pure.

DEFINITION (Pure-Environment). An environment  $\Gamma$  is pure if  $\text{pure } \Gamma = \Gamma$ .

DEFINITION (Pure-Type). A type is pure if the type can exist in a pure environment.

From the definition, it's obvious that the type  $B$  is pure, while  $E$  is impure. However, a pure function type doesn't imply the corresponding function is pure. For example,  $E \rightarrow B$  is a pure type, but functions that inhabit the type may have side effects, thus is impure. The meta-theory of the system ensures that if the capability type  $E$  doesn't appear in the type signature of a function, then the function must be pure. For example, all functions that can be typed as  $B \rightarrow B$  are guaranteed to be pure.

### 2.1.2 Stoic Functions and Free Functions

The variable-capturing discipline makes the functions in STLC-Pure different from functions in standard STLC. In STLC, functions can capture any variables in scope, while in STLC-Pure functions can only capture variables whose type is not the capability type  $E$ . To differentiate them (which is important as in later systems both exist), we call the more effect-disciplined functions *stoic functions* (or stoics) and the other *free functions*. We use  $\rightarrow$  to denote the type of stoic functions and  $\Rightarrow$  to denote the type of free functions.

Stoic functions are essential in capability-based effect systems. If functions are allowed to capture capability variables in scope, it will be impossible to tell whether a function has side effect or not (and what kind of effect) by just checking its type. Stoic functions are effect-disciplined in the sense that the only way for stoic functions to have side effects is to pass a capability as parameter, thus it can be captured by the type system.

Stoic functions are not necessarily pure functions. Stoic functions can have side effects, and if they do have side effects they are honest about that in their type signature. For example, the following function `hello` is a stoic function with IO effects<sup>2</sup>.

```
def hello(c:IO) = println("hello, world!", c)
```

In the following code snippet, the function `f` must be pure, as it doesn't take any capability as parameter. The type system guarantees that the function indeed cannot produce any side effects.

```
def twice(f: Int -> Int)(x: Int) = f (f x)
```

---

<sup>2</sup>For the sake of readability, we'll use a syntax similar to Scala in this report. In particular, we'll use  $\rightarrow$  for the type of stoic functions, and  $\Rightarrow$  for free functions.

### 2.1.3 Where do Effects Come From

There is no formalization of effects in the current effect system. We assume the existence of primitive functions like `println` and `readln`, which take capabilities to produce side effects.

```
def println: String -> IO -> ()
def readln: IO -> String
```

### 2.1.4 Where do Capabilities Come From

It is impossible to create capabilities in the current system. Where do capabilities come from? There are two possible answers: (1) all capabilities are from the run-time and passed to the program through the main method; (2) there are no capabilities; they can be erased before evaluation, without changing the meaning of programs.

### 2.1.5 Why Treat Variables as Values

As discussed above, there is no way to create a value of capabilities explicitly. Thus, a function taking a parameter of the capability type `E` can never be executed in the *call-by-value* semantics, unless variables are values. The same is true for the base type `B`.

Treating variables as values ensures that substitution of a term with a value of the base type `B` or the capability type `E` can actually happen in the system, thus makes the preservation proof more convincing.

### 2.1.6 What If a Function Has More than One Side Effect

There is no support for a function with more than one kind of effects in the current system. For example, in the following code snippet, `c1` cannot be used in the function body, as it's removed by `pure` in the typing of the function body.

```
def error(e:Error)(c1:IO)(c2:Throw) = {
  println("error happen!", c1) // Error, can't capture c1
  throw c2 e
}
```

It's straight-forward to extend the system with pairs or tuples to overcome this limitation. However, this is not an issue for later systems with *free functions*, thus we don't pursue the extension of pairs and tuples here.

## 2.2 Soundness

We follow the standard formulation of soundness in TAPL [Pie02], which consists of *progress* and *preservation*, defined as follows:

THEOREM (Progress). If  $\emptyset \vdash t : T$ , then either  $t$  is a value or there is some  $t'$  with  $t \longrightarrow t'$ .

THEOREM (Preservation). If  $\Gamma \vdash t : T$ , and  $t \longrightarrow t'$ , then  $\Gamma \vdash t' : T$ .

The proof of progress is the same as the proof in standard STLC. However, there is a significant difference in the proof of preservation. The classic proof of preservation for STLC (as shown in TAPL) depends on a substitution lemma, which is formulated as follows:

LEMMA (Substitution-Classic). If  $\Gamma, x : S \vdash t : T$ , and  $\Gamma \vdash s : S$ , then  $\Gamma \vdash [x \mapsto s]t : T$ .

However, this substitution lemma doesn't hold in the current system. For a counter-example, let's assume that  $\Gamma = \{f : E \rightarrow B, c : E\}$ , then it's obviously that following two typing relations hold:

$$f : E \rightarrow B, c : E, x : B \vdash \lambda z : B. x : B \rightarrow B$$

$$f : E \rightarrow B, c : E \vdash f c : B$$

However, the following typing relation doesn't hold if we replace  $x$  with  $f c$ .

$$f : E \rightarrow B, c : E \vdash \lambda z : B. f c : B \rightarrow B.$$

In fact, the substituted term  $\lambda z : B. f c$  cannot be typed, as according to the typing rule T-ABS, it cannot capture the capability variable  $c$  in the environment. To overcome this problem, we stipulate that the term  $s$  must be a value. Remember that in the current system, both lambda abstractions and variables are values, thus substitution of variables of the capability type  $E$  and the base type  $B$  can happen. The new formulation is as follows:

LEMMA (Substitution-New). If  $\Gamma, x : S \vdash t : T$ ,  $s$  is a value and  $\Gamma \vdash s : S$ , then  $\Gamma \vdash [x \mapsto s]t : T$ .

The restriction that  $s$  must be a value implies that in a function call, arguments must first be evaluated to a value before the function is applied. Therefore, capabilities can only work with strict evaluation.

Interestingly, this strict evaluation requirement contrasts capability-based effect systems with monad-based effect systems. In Haskell, if strict evaluation is adopted, it will be impossible to track effects in the type system, as demonstrated by the following code snippet:

```
inc n = (\x -> n + 1) (putStrLn (show n))
```

The function `inc` has the type  $(\text{Num } a, \text{Show } a) \Rightarrow a \rightarrow a$ . By just checking its type, we would think it has no side effects because no IO monads appear in the type signature. However, if Haskell adopts strict evaluation, the function call `putStrLn (show n)` will be executed, thus breaking the monad-based effect system.

## 2.3 Effect Safety

Does the system really work? This question prompts us to formulate and prove effect safety of the system. We start by formulating effect safety informally, then put forward a formal formulation, and finally prove effect safety of the system.

### 2.3.1 Informal Formulation

A straight-forward violation of effect safety is for functions that are taken as pure to have side effects inside the function body. Thus, a tentative formulation would be as follows:

DEFINITION (Effect-Safety-Informally-1). A function typed in a pure environment cannot have side effects inside.

However, this formulation is obviously problematic, as we know stoic functions can have side effects if it takes a capability parameter. Thus, we need to restrict the functions to those not taking capability parameters:

DEFINITION (Effect-Safety-Informally-2). A function, not taking any capability parameter and typed in a pure environment, cannot have side effects inside.

This formulation looks more satisfactory, but it's a little cumbersome. If we inspect the typing rule T-ABS closely, we can find that if  $S$  is not a capability type,  $\text{pure}(\Gamma), x : S$  is equal to  $\text{pure}(\Gamma, x : S)$ .

$$\frac{\text{pure}(\Gamma), x : S \vdash t_2 : T}{\Gamma \vdash \lambda x : S. t_2 : S \rightarrow T} \quad (\text{T-ABS})$$



Thus, instead of saying the function  $\lambda x:S. t_2$  cannot have side effects inside, we say the term  $t_2$  cannot have side effects in a pure environment. As we know, capabilities are required to produce side effects. Thus, the term  $t_2$  cannot have side effects if we cannot construct a term of the capability type  $E$  in a pure environment. This observation leads us to the following statement of effect safety:

DEFINITION (Effect-Safety-Informally-3). It's impossible to construct a term of the capability type  $E$  in a pure environment.

However, this formulation cannot be proved. For a counter-example, let's assume  $\Gamma = \{f : B \rightarrow E, x : B\}$ . It's obvious that  $\Gamma$  is pure, but we can construct the term  $f x$  of the capability type  $E$ .

The cause of the problem is that in a pure environment, there might exist uninhabited types like  $B \rightarrow E$ . Existence of uninhabited types in a pure environment doesn't pose a problem to the system; a function taking a parameter of an uninhabited type can never be actually called, thus is always effect-safe. So we only need to consider environments with only variables of inhabited types.

To convince readers that the current system is effect-safe, we need to exclude and only exclude uninhabited types from the pure environment and then prove that it is impossible to construct a term of the capability type  $E$  in this restricted environment. We arrive at the following formulation:

DEFINITION (Effect-Safety-Informally-4). It's impossible to construct a term of the capability type  $E$  in a pure environment with only variables of inhabited types.

### 2.3.2 Inhabited Types and Environments

We need to define the concept *inhabited types* precisely. What types are inhabited? Obviously, if  $\emptyset \vdash t : T$ , then  $T$  is inhabited. However, given a typing relation  $\Gamma \vdash t : T$ , we cannot immediately conclude that  $T$  is inhabited. We need to ensure that  $\Gamma$  only contains inhabited types. Otherwise, any type is inhabited if  $\Gamma$  contains a variable of the corresponding type.

An intuition is that, given  $\Gamma \vdash t : T, x : S \in \Gamma$  and  $S$  is inhabited, we can remove  $x : S$  from  $\Gamma$  and substitute  $x$  in the term  $t$  with a witness of the type  $S$  to obtain a new term  $t'$ . The substitution lemma tells us that the new term  $t'$  still has the type  $T$ . Continue this line of thought, we'll find out that all inhabited types can be inhabited in the empty environment. Thus, a tentative definition is as follows:

DEFINITION (Inhabited-Type-First-Try). A type  $T$  is inhabited if there exists a term  $t$  with  $\emptyset \vdash t : T$ .

However, this definition is not satisfactory in our case, as in STLC-Pure there doesn't exist values for the base type  $B$  and the capability type  $E$ , except variables, and we do want both types to be inhabited. A natural approach is to extend the empty environment with one variable of the base type and one of the capability type:

DEFINITION (Inhabited-Type-Second-Try). A type  $T$  is inhabited if there exists a term  $t$  with  $x : B, y : E \vdash t : T$ .

This definition indeed gives us all inhabited types in STLC-Pure. Types like  $E, B, E \rightarrow B, E \rightarrow E, (B \rightarrow E) \rightarrow E$ , etc., are all inhabited, while types like  $B \rightarrow E$  and  $E \rightarrow B \rightarrow E$  are uninhabited. The reader might want to construct a witness of the type  $B \rightarrow E$  as follows:

$$x : B, y : E \vdash \lambda x : B. y : B \rightarrow E$$

The typing relation doesn't hold, because the typing rule T-ABS would remove the binding  $y : E$  from the environment in the typing of the function body  $y$ .

The second formulation looks good, however, we can do better with the following definition:

DEFINITION (Inhabited-Type-Final). A type  $T$  is inhabited if there exists a value  $v$  with the typing  $x : B, y : E \vdash v : T$ .

What if the term  $t$  in the second definition is an application? In that case,  $t$  must be able to take a step until it becomes a value due to *progress* and *normalization* of the system<sup>3</sup>. And the *preservation* theorem tells us the type remains unchanged during evaluation. This final definition makes proofs related to inhabited types simpler. It's useful to give a definition of inhabited environments as well:

DEFINITION (Inhabited-Environment). An environment  $\Gamma$  is inhabited if it only contains variables of inhabited types.

### 2.3.3 Formalization

With the formal definition of inhabited environment, we can formalize effect safety as follows:

---

<sup>3</sup>We didn't prove normalization of STLC-Pure, but the proof should be similar to the proof in standard STLC. We only proved progress in the empty environment, and the proof can be adapted to prove progress under  $\{x : B, y : E\}$ .

DEFINITION (Effect-Safety-Inhabited). If  $\Gamma$  is a pure and inhabited environment, then there doesn't exist  $t$  with  $\Gamma \vdash t : E$ .

However, this formulation doesn't give rise to a direct proof. In fact, this statement is too strong. Some uninhabited types, such as  $E \rightarrow B \rightarrow E$ , don't enable us to create a term of the type  $E$ , thus it's safe to keep them in the environment. This implies it's possible to impose a looser restriction on  $\Gamma$ , as long as all types that can appear in a pure and inhabited environment can also appear in  $\Gamma$ .

When we examine the problem more closely, we found that through the lens of the *Curry-Howard isomorphism*, effect safety actually says that it is impossible to prove the capability type  $E$  from a group of "good" premises. Thus, we can classify all types (propositions) into two groups: in one group  $E$  cannot be proved and in the other group  $E$  can be proved. This leads us to a formulation of *capsafe environment*<sup>4</sup> given in Figure 2.2. The coined word *capsafe* is an abbreviation of capability-safe, and *caprod* an abbreviation of capability-producing.

<b>Capsafe Type</b>	<b>Caprod Type</b>
$B \text{ capsafe} \quad (\text{CS-BASE})$	$E \text{ caprod} \quad (\text{CP-EFF})$
$\frac{S \text{ caprod}}{S \rightarrow T \text{ capsafe}} \quad (\text{CS-FUN1})$	$\frac{S \text{ capsafe} \quad T \text{ caprod}}{S \rightarrow T \text{ caprod}} \quad (\text{CP-FUN})$
$\frac{T \text{ capsafe}}{S \rightarrow T \text{ capsafe}} \quad (\text{CS-FUN2})$	<b>Capsafe Environment</b>
	$\emptyset \text{ capsafe} \quad (\text{CE-EMPTY})$
	$\frac{\Gamma \text{ capsafe} \quad T \text{ capsafe}}{\Gamma, x : T \text{ capsafe}} \quad (\text{CE-VAR})$

Figure 2.2: System STLC-Pure Capsafe Environment

In the definition, types like  $B \rightarrow B$ ,  $E \rightarrow E$ ,  $E \rightarrow B$  and  $(B \rightarrow E) \rightarrow B$  are considered as *capsafe*, while types like  $B \rightarrow E$ ,  $(E \rightarrow B) \rightarrow E$  are considered as *caprod*. Only *capsafe* types can appear in a *capsafe* environment. To inspect the formalization in detail, we can ask several questions.

*Are capsafe types inhabited?* Not necessarily. The type  $E \rightarrow B \rightarrow E$  is capsafe but uninhabited.

<sup>4</sup>Sandro Stucki initially suggested the idea of using *caprod* for the definition of *pure* environments. I developed it to be a formulation of *capsafe* environments and used it in the proof of effect safety.

Allowing this type in the capsafe environment doesn't enable us to construct a term of the capability type E.

*Are inhabited types capsafe?* Yes, except the capability type E. As the capability type E cannot appear in the pure environment, this is not a problem.

*Are caprod types uninhabited?* Yes, except the capability type E. As E is also excluded in the pure environment, it's justified to remove it from the capsafe environment.

Why this formulation of capsafe environment is acceptable? In short, it is because the statement *Effect-Safety-Inhabited* is logically implied by the more general statement *Effect-Safety*:

DEFINITION (Effect-Safety). If  $\Gamma$  is capsafe, then there doesn't exist  $t$  with  $\Gamma \vdash t : E$ .

The logical implication holds because a pure and inhabited environment is also a capsafe environment. This claim has been formally proved:

LEMMA (Inhabited-Capsafe). If the type  $T$  is inhabited, then either  $T$  is capsafe or  $T = E$ .

THEOREM (Inhabited-Pure-Capsafe-Env). If  $\Gamma$  is pure and inhabited, then  $\Gamma$  is capsafe.

### 2.3.4 Proof

The proof of effect safety depends on following lemmas, most of them are straight-forward to prove. Effect safety follows immediately from the lemma *Capsafe-Env-Capsafe*.

LEMMA (Capsafe-Not-Caprod). If type  $T$  is capsafe, then  $T$  is not caprod.

LEMMA (Capsafe-Or-Caprod). For any  $T$ ,  $T$  is either capsafe or caprod.

LEMMA (Capsafe-Env-Capsafe). If  $\Gamma$  is capsafe and  $\Gamma \vdash t : T$ , then  $T$  is capsafe.

THEOREM (Effect-Safety). If  $\Gamma$  is capsafe, then there doesn't exist  $t$  with  $\Gamma \vdash t : E$ .

### 2.3.5 An Intuitive Proof

There exists an intuitive proof of effect safety without resorting to capsafe environments. The main insight is that the statement *Effect-Safety-Inhabited* is logically implied by the statement *Effect-Safety-Intuitive*:

DEFINITION (Effect-Safety-Inhabited). If  $\Gamma$  is a pure and inhabited environment, then there doesn't exist  $t$  with  $\Gamma \vdash t : E$ .

DEFINITION (Effect-Safety-Intuitive). There doesn't exist value  $v$  with  $x : B \vdash v : E$ .

The statement *Effect-Safety-Intuitive* trivially holds, because the value  $v$  can either be a variable or a function, in neither case can it be typed as the capability type  $E$ .

Why does the logical implication hold? In short, if *Effect-Safety-Inhabited* doesn't hold, then *Effect-Safety-Intuitive* doesn't hold either. Thus, the latter logically implies the former. In the following, we present an informal proof, the core idea is that if  $\Gamma$  is pure and inhabited, and  $\Gamma \vdash t : E$ , then we can “collapse”  $\Gamma$  to  $\{b : B\}$  through substitution of the witnesses of the pure and inhabited types in  $\Gamma$ .

For a pure and inhabited environment  $\Gamma = \{x : T, y : S, \dots, z : U\}$ , if there exists  $t$  with  $\Gamma \vdash t : E$ , then the typing relation still holds by extending the environment with  $b : B$ :

$$b : B, x : T, y : S, \dots, z : U \vdash t : E$$

The type  $U$  is pure and inhabited, as  $\Gamma$  is pure and inhabited. According to the definition of *inhabited type*, there exists a value  $u$  with  $b : B, e : E \vdash u : U$ . As  $u$  is a value, it can be either a variable or a function. If  $u$  is a variable, it can only be  $b$ , as  $U$  is a pure type, thus it cannot be the capability type  $E$ . If  $u$  is a function, we have  $b : B \vdash u : U$ , as in the typing of stoic functions the rule T-ABS will remove the binding  $e : E$  from the environment. In both cases, we have  $b : B \vdash u : U$ . Now using the substitution lemma, we have:

$$b : B, x : T, y : S, \dots \vdash [z \mapsto u]t : E$$

Continue this process, we can reduce the typing environment to be  $\{b : B\}$  and the term to  $t'$ :

$$b : B \vdash t' : E$$

Now combining *progress* and *normalization* of STLC-Pure<sup>5</sup>,  $t'$  can take finite evaluation steps to become a value  $v$ :

$$b : B \vdash v : E$$

To summarize, we have given an intuitive proof of the following statement:

$$\neg \text{Effect-Safety-Inhabited} \rightarrow \neg \text{Effect-Safety-Intuitive}$$

By the logical law *contraposition*, we have:

$$\text{Effect-Safety-Intuitive} \rightarrow \text{Effect-Safety-Inhabited}$$

<sup>5</sup>We only proved that progress holds in an empty environment, it's easy to prove it also holds under  $\{b : B\}$ . We didn't prove normalization of STLC-Pure, but the proof should be similar to the proof in standard STLC.

As *Effect-Safety-Intuitive* trivially holds, *Effect-Safety-Inhabited* follows by *modus ponens*.

Though conceptually simpler, the mechanized proof necessitates the proof of the normalization theorem, which is more involved than the proof based on capsafe environments. On the other hand, the approach based on capsafe environments works even if normalization doesn't hold. Therefore, we don't take the intuitive approach in the formal development.

### 3 System STLC-Impure

This chapter describes an extension of the system STLC-Pure with free functions. As stoic functions can be used as free functions, it's natural to integrate subtyping in the system.

We'll first introduce the formalization, then discuss soundness and effect safety. In the discussion, we'll focus on its difference from the system STLC-Pure.

#### 3.1 Definitions

Initially, we arrived at a formulation of the system shown in Figure 3.1. It's a straight-forward extension of STLC-Pure with subtyping and free functions.

The definition is all good, except that preservation breaks. The problem is caused by typing an impure term as  $\text{Top}$ . To see a concrete example, let's assume  $\Gamma = \{c : E\}$ . It's obvious that the following typing relation holds:

$$c : E \vdash (\lambda x:\text{Top}. \lambda y:\text{B}. x) c : \text{B} \rightarrow \text{Top}$$

However, after one evaluation step<sup>6</sup>, we get the term  $\lambda y:\text{B}. c$ , which can at best be typed as  $\text{B} \Rightarrow \text{Top}$ . Thus preservation breaks. This problem leads us to two different formulations.

The first one is to introduce two different  $\text{Top}$  types,  $\text{Top-Pure}$  and  $\text{Top-Impure}$ , the former is pure and the latter is impure<sup>7</sup>. The capability type  $E$  and free function type  $S \Rightarrow T$  are not subtypes of  $\text{Top-Pure}$ . The subtyping hierarchy is shown in Figure 3.2. This formulation works well, and we've proved soundness and effect safety for the formulation. However, we lose the simplicity of the type system. And it's counter-intuitive to forbid variables of the type  $\text{Top-Impure}$  in pure environments, as we cannot create side effects with a variable of the type  $\text{Top-Impure}$ .

The second possibility is to keep the elegance of the type system and change the evaluation rules. All terms of the type  $\text{Top}$  are equivalent, because we can do nothing with a term of the type  $\text{Top}$ . This observation inspires us to introduce a value  $\text{top}$  and replace the standard E-APPABS rule with two evaluation rules as follows:

---

<sup>6</sup>Note that variables are values, thus we can take a step here. We can also construct a counter-example by wrapping  $c$  in a free function like  $\lambda x:\text{B}. c$ .

<sup>7</sup>This idea is suggested by Sandro Stucki.

<p><b>Syntax</b></p> <p><math>t ::=</math></p> <table style="width: 100%; border: none;"> <tr> <td style="width: 30%;"><math>x</math></td> <td style="width: 30%;"></td> <td style="width: 40%;">terms:</td> </tr> <tr> <td><math>\lambda x:T. t</math></td> <td></td> <td>variable</td> </tr> <tr> <td><math>tt</math></td> <td></td> <td>abstraction</td> </tr> <tr> <td></td> <td></td> <td>application</td> </tr> </table> <p><math>v ::=</math></p> <table style="width: 100%; border: none;"> <tr> <td style="width: 30%;"><math>\lambda x:T. t</math></td> <td style="width: 30%;"></td> <td style="width: 40%;">values:</td> </tr> <tr> <td><math>x</math></td> <td></td> <td>abstraction value</td> </tr> <tr> <td></td> <td></td> <td>variable value</td> </tr> </table> <p><math>T ::=</math></p> <table style="width: 100%; border: none;"> <tr> <td style="width: 30%;"><math>\text{Top}</math></td> <td style="width: 30%;"></td> <td style="width: 40%;">types:</td> </tr> <tr> <td><math>B</math></td> <td></td> <td>top type</td> </tr> <tr> <td><math>E</math></td> <td></td> <td>basic type</td> </tr> <tr> <td><math>T \rightarrow T</math></td> <td></td> <td>capability type</td> </tr> <tr> <td><math>T \Rightarrow T</math></td> <td></td> <td>type of stocic funs</td> </tr> <tr> <td></td> <td></td> <td>type of free funs</td> </tr> </table> <p><b>Evaluation</b></p> <div style="text-align: right; border: 1px solid black; padding: 2px; display: inline-block;"><math>t \longrightarrow t'</math></div> $\frac{t_1 \longrightarrow t'_1}{t_1 t_2 \longrightarrow t'_1 t_2} \quad (\text{E-APP1})$ $\frac{t_2 \longrightarrow t'_2}{v_1 t_2 \longrightarrow v_1 t'_2} \quad (\text{E-APP2})$ $(\lambda x:T. t_1)v_2 \longrightarrow [x \mapsto v_2]t_1 \quad (\text{E-APPABS})$ <p><b>Pure Environment</b></p> <table style="width: 100%; border: none;"> <tr> <td><math>\text{pure}(\emptyset)</math></td> <td><math>=</math></td> <td><math>\emptyset</math></td> </tr> <tr> <td><math>\text{pure}(\Gamma, x : E)</math></td> <td><math>=</math></td> <td><math>\text{pure}(\Gamma)</math></td> </tr> <tr style="background-color: #e0e0e0;"> <td><math>\text{pure}(\Gamma, x : S \Rightarrow T)</math></td> <td><math>=</math></td> <td><math>\text{pure}(\Gamma)</math></td> </tr> <tr> <td><math>\text{pure}(\Gamma, x : T)</math></td> <td><math>=</math></td> <td><math>\text{pure}(\Gamma), x : T</math></td> </tr> </table>	$x$		terms:	$\lambda x:T. t$		variable	$tt$		abstraction			application	$\lambda x:T. t$		values:	$x$		abstraction value			variable value	$\text{Top}$		types:	$B$		top type	$E$		basic type	$T \rightarrow T$		capability type	$T \Rightarrow T$		type of stocic funs			type of free funs	$\text{pure}(\emptyset)$	$=$	$\emptyset$	$\text{pure}(\Gamma, x : E)$	$=$	$\text{pure}(\Gamma)$	$\text{pure}(\Gamma, x : S \Rightarrow T)$	$=$	$\text{pure}(\Gamma)$	$\text{pure}(\Gamma, x : T)$	$=$	$\text{pure}(\Gamma), x : T$	<p><b>Typing</b></p> <div style="text-align: right; border: 1px solid black; padding: 2px; display: inline-block;"><math>\Gamma \vdash x : T</math></div> $\frac{x : T \in \Gamma}{\Gamma \vdash x : T} \quad (\text{T-VAR})$ $\frac{\text{pure}(\Gamma), x : S \vdash t_2 : T}{\Gamma \vdash \lambda x:S. t_2 : S \rightarrow T} \quad (\text{T-ABS1})$ $\frac{\Gamma, x : S \vdash t_2 : T}{\Gamma \vdash \lambda x:S. t_2 : S \Rightarrow T} \quad (\text{T-ABS2})$ $\frac{\Gamma \vdash t_1 : S \Rightarrow T \quad \Gamma \vdash t_2 : S}{\Gamma \vdash t_1 t_2 : T} \quad (\text{T-APP})$ $\frac{\Gamma \vdash t : S \quad S <: T}{\Gamma \vdash t : T} \quad (\text{T-SUB})$ <p><b>Subtyping</b></p> <div style="text-align: right; border: 1px solid black; padding: 2px; display: inline-block;"><math>S &lt;: T</math></div> $T <: \text{Top} \quad (\text{S-TOP})$ $T <: T \quad (\text{S-REFL})$ $\frac{S <: U \quad U <: T}{S <: T} \quad (\text{S-TRANS})$ $S \rightarrow T <: S \Rightarrow T \quad (\text{S-DEGEN})$ $\frac{S1 <: S2 \quad T2 <: T1}{S2 \rightarrow T2 <: S1 \rightarrow T1} \quad (\text{S-FUN1})$ $\frac{S1 <: S2 \quad T2 <: T1}{S2 \Rightarrow T2 <: S1 \Rightarrow T1} \quad (\text{S-FUN2})$
$x$		terms:																																																		
$\lambda x:T. t$		variable																																																		
$tt$		abstraction																																																		
		application																																																		
$\lambda x:T. t$		values:																																																		
$x$		abstraction value																																																		
		variable value																																																		
$\text{Top}$		types:																																																		
$B$		top type																																																		
$E$		basic type																																																		
$T \rightarrow T$		capability type																																																		
$T \Rightarrow T$		type of stocic funs																																																		
		type of free funs																																																		
$\text{pure}(\emptyset)$	$=$	$\emptyset$																																																		
$\text{pure}(\Gamma, x : E)$	$=$	$\text{pure}(\Gamma)$																																																		
$\text{pure}(\Gamma, x : S \Rightarrow T)$	$=$	$\text{pure}(\Gamma)$																																																		
$\text{pure}(\Gamma, x : T)$	$=$	$\text{pure}(\Gamma), x : T$																																																		

Figure 3.1: System STLC-Impure First Formulation



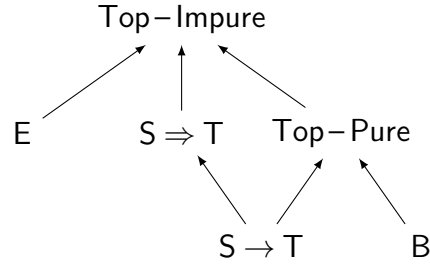


Figure 3.2: Subtyping: Top-Pure and Top-Impure

$$\frac{T \neq \text{Top}}{(\lambda x:T. t_1)v_2 \longrightarrow [x \mapsto v_2]t_1} \text{ (E-APPABS1)} \quad \frac{T = \text{Top}}{(\lambda x:T. t_1)v_2 \longrightarrow [x \mapsto \text{top}]t_1} \text{ (E-APPABS2)}$$

The two rules have the effect that if a function takes a parameter of the type Top, when called it will drop the passed parameter and replace it with the value top. There is no need to consider the case that the parameter type T is like  $B \rightarrow \text{Top}$  or  $(\text{Top} \rightarrow B) \Rightarrow B$ , as function types cannot mask impure terms as pure like Top does. We follow this approach in the formulation and the full definition is presented in Figure 3.3.

Note that in the current system, we need to change the definition of the function pure to exclude free function types from the pure environment. This restriction is important, because there's no way to know what side effects there might be inside free functions. If stoic functions have access to free functions, we'll lose the ability to track the effects of stoic functions in the type system.

However, mixed-arrow types are allowed in pure environments, as long as the first arrow is  $\rightarrow$ . For example, the type  $E \rightarrow B \Rightarrow B$  is allowed in the pure environment. As will be shown in the effect safety section, we have formally proved that having such mixed-arrow types in the pure environment doesn't enable pure stoic functions to produce side effects.

In principle, we can keep some free functions that can never be called in the pure environment, such as  $(B \rightarrow E) \Rightarrow B$ , as it's impossible to get an actual instance of the inhabited type  $B \rightarrow E$  in order to call this function. However, it's not useful to have functions in the environment if they can never be called. Thus, to simplify the system without sacrificing usability, we removed all free function types from the pure environment.

<p><b>Syntax</b></p> <p><math>t ::=</math></p> <table style="width: 100%; border: none;"> <tr> <td style="padding: 2px 10px;"><math>\text{top}</math></td> <td style="padding: 2px 10px;">terms:</td> </tr> <tr> <td style="padding: 2px 10px;"><math>x</math></td> <td style="padding: 2px 10px;">top value</td> </tr> <tr> <td style="padding: 2px 10px;"><math>\lambda x:T. t</math></td> <td style="padding: 2px 10px;">variable</td> </tr> <tr> <td style="padding: 2px 10px;"><math>tt</math></td> <td style="padding: 2px 10px;">abstraction</td> </tr> <tr> <td></td> <td style="padding: 2px 10px;">application</td> </tr> </table> <p><math>v ::=</math></p> <table style="width: 100%; border: none;"> <tr> <td style="padding: 2px 10px;"><math>\lambda x:T. t</math></td> <td style="padding: 2px 10px;">values:</td> </tr> <tr> <td style="padding: 2px 10px;"><math>x</math></td> <td style="padding: 2px 10px;">abstraction value</td> </tr> <tr> <td style="padding: 2px 10px;"><math>\text{top}</math></td> <td style="padding: 2px 10px;">variable value</td> </tr> <tr> <td></td> <td style="padding: 2px 10px;">top value</td> </tr> </table> <p><math>T ::=</math></p> <table style="width: 100%; border: none;"> <tr> <td style="padding: 2px 10px;"><math>\text{Top}</math></td> <td style="padding: 2px 10px;">types:</td> </tr> <tr> <td style="padding: 2px 10px;"><math>B</math></td> <td style="padding: 2px 10px;">top type</td> </tr> <tr> <td style="padding: 2px 10px;"><math>E</math></td> <td style="padding: 2px 10px;">basic type</td> </tr> <tr> <td style="padding: 2px 10px;"><math>T \rightarrow T</math></td> <td style="padding: 2px 10px;">capability type</td> </tr> <tr> <td style="padding: 2px 10px;"><math>T \Rightarrow T</math></td> <td style="padding: 2px 10px;">type of stoc funs</td> </tr> <tr> <td></td> <td style="padding: 2px 10px;">type of free funs</td> </tr> </table> <p><b>Evaluation</b> <span style="float: right; border: 1px solid black; padding: 2px;"><math>t \longrightarrow t'</math></span></p> $\frac{t_1 \longrightarrow t'_1}{t_1 t_2 \longrightarrow t'_1 t_2} \quad (\text{E-APP1})$ $\frac{t_2 \longrightarrow t'_2}{v_1 t_2 \longrightarrow v_1 t'_2} \quad (\text{E-APP2})$ $\frac{\text{Top} \neq \text{Top}}{(\lambda x:T. t_1) v_2 \longrightarrow [x \mapsto v_2] t_1} \quad (\text{E-APPABS1})$ $\frac{\text{Top} = \text{Top}}{(\lambda x:T. t_1) v_2 \longrightarrow [x \mapsto \text{top}] t_1} \quad (\text{E-APPABS2})$ <p><b>Pure Environment</b></p> <table style="width: 100%; border: none;"> <tr> <td style="padding: 2px 10px;"><math>\text{pure}(\emptyset)</math></td> <td style="padding: 2px 10px;"><math>=</math></td> <td style="padding: 2px 10px;"><math>\emptyset</math></td> </tr> <tr> <td style="padding: 2px 10px;"><math>\text{pure}(\Gamma, x : E)</math></td> <td style="padding: 2px 10px;"><math>=</math></td> <td style="padding: 2px 10px;"><math>\text{pure}(\Gamma)</math></td> </tr> <tr> <td style="padding: 2px 10px;"><math>\text{pure}(\Gamma, x : S \Rightarrow T)</math></td> <td style="padding: 2px 10px;"><math>=</math></td> <td style="padding: 2px 10px;"><math>\text{pure}(\Gamma)</math></td> </tr> <tr> <td style="padding: 2px 10px;"><math>\text{pure}(\Gamma, x : T)</math></td> <td style="padding: 2px 10px;"><math>=</math></td> <td style="padding: 2px 10px;"><math>\text{pure}(\Gamma), x : T</math></td> </tr> </table>	$\text{top}$	terms:	$x$	top value	$\lambda x:T. t$	variable	$tt$	abstraction		application	$\lambda x:T. t$	values:	$x$	abstraction value	$\text{top}$	variable value		top value	$\text{Top}$	types:	$B$	top type	$E$	basic type	$T \rightarrow T$	capability type	$T \Rightarrow T$	type of stoc funs		type of free funs	$\text{pure}(\emptyset)$	$=$	$\emptyset$	$\text{pure}(\Gamma, x : E)$	$=$	$\text{pure}(\Gamma)$	$\text{pure}(\Gamma, x : S \Rightarrow T)$	$=$	$\text{pure}(\Gamma)$	$\text{pure}(\Gamma, x : T)$	$=$	$\text{pure}(\Gamma), x : T$	<p><b>Typing</b> <span style="float: right; border: 1px solid black; padding: 2px;"><math>\Gamma \vdash x : T</math></span></p> $\frac{}{\Gamma \vdash \text{top} : \text{Top}} \quad (\text{T-TOP})$ $\frac{x : T \in \Gamma}{\Gamma \vdash x : T} \quad (\text{T-VAR})$ $\frac{\text{pure}(\Gamma), x : S \vdash t_2 : T}{\Gamma \vdash \lambda x:S. t_2 : S \rightarrow T} \quad (\text{T-ABS1})$ $\frac{\Gamma, x : S \vdash t_2 : T}{\Gamma \vdash \lambda x:S. t_2 : S \Rightarrow T} \quad (\text{T-ABS2})$ $\frac{\Gamma \vdash t_1 : S \Rightarrow T \quad \Gamma \vdash t_2 : S}{\Gamma \vdash t_1 t_2 : T} \quad (\text{T-APP})$ $\frac{\Gamma \vdash t : S \quad S <: T}{\Gamma \vdash t : T} \quad (\text{T-SUB})$ <p><b>Subtyping</b> <span style="float: right; border: 1px solid black; padding: 2px;"><math>S &lt;: T</math></span></p> $T <: \text{Top} \quad (\text{S-TOP})$ $T <: T \quad (\text{S-REFL})$ $\frac{S <: U \quad U <: T}{S <: T} \quad (\text{S-TRANS})$ $S \rightarrow T <: S \Rightarrow T \quad (\text{S-DEGEN})$ $\frac{S1 <: S2 \quad T2 <: T1}{S2 \rightarrow T2 <: S1 \rightarrow T1} \quad (\text{S-FUN1})$ $\frac{S1 <: S2 \quad T2 <: T1}{S2 \Rightarrow T2 <: S1 \Rightarrow T1} \quad (\text{S-FUN2})$
$\text{top}$	terms:																																										
$x$	top value																																										
$\lambda x:T. t$	variable																																										
$tt$	abstraction																																										
	application																																										
$\lambda x:T. t$	values:																																										
$x$	abstraction value																																										
$\text{top}$	variable value																																										
	top value																																										
$\text{Top}$	types:																																										
$B$	top type																																										
$E$	basic type																																										
$T \rightarrow T$	capability type																																										
$T \Rightarrow T$	type of stoc funs																																										
	type of free funs																																										
$\text{pure}(\emptyset)$	$=$	$\emptyset$																																									
$\text{pure}(\Gamma, x : E)$	$=$	$\text{pure}(\Gamma)$																																									
$\text{pure}(\Gamma, x : S \Rightarrow T)$	$=$	$\text{pure}(\Gamma)$																																									
$\text{pure}(\Gamma, x : T)$	$=$	$\text{pure}(\Gamma), x : T$																																									

Figure 3.3: System STLC-Impure

## 3.2 Soundness

We proved both progress and preservation of the system.

THEOREM (Progress). If  $\emptyset \vdash t : T$ , then either  $t$  is a value or there is some  $t'$  with  $t \longrightarrow t'$ .

THEOREM (Preservation). If  $\Gamma \vdash t : T$ , and  $t \longrightarrow t'$ , then  $\Gamma \vdash t' : T$ .

As you can imagine, now we need two different substitution lemmas in the proof of preservation, corresponding to the two reduction rules.

LEMMA (Substitution-Not-Top). If  $\Gamma, x : S \vdash t : T$ ,  $S \neq \text{Top}$ ,  $s$  is a value and  $\Gamma \vdash s : S$ , then  $\Gamma \vdash [x \mapsto s]t : T$ .

LEMMA (Substitution-Top). If  $\Gamma, x : \text{Top} \vdash t : T$ , then  $\Gamma \vdash [x \mapsto \text{top}]t : T$ .

We restrict  $s$  to be a value in the lemma *Substitution-Not-Top* for the same reason as in the system STLC-Pure.

## 3.3 Effect Safety

We follow the same approach as in the system STLC-Pure in the formulation of effect safety. The formulation is an extension of the definition of *capsafe* and *caprod* in STLC-Pure with free function types.

### 3.3.1 Formulation

The definitions of *inhabited type* and *inhabited environment* are the same as in the system STLC-Pure.

With the presence of free functions, the previous formulation of effect safety is not enough. We not only need to ensure that it's impossible to construct a term of the capability type  $E$  in pure and inhabited environments, but also need to ensure only stoic functions can be called in pure and inhabited environments. The two conditions together guarantee that there cannot be actual side effects inside a pure stoic function. Thus, we need two statements about effect safety.

DEFINITION (Effect-Safety-Inhabited-1). If  $\Gamma$  is a pure and inhabited environment, then there doesn't exist  $t$  with  $\Gamma \vdash t : E$ .

DEFINITION (Effect-Safety-Inhabited-2). If  $\Gamma$  is a pure and inhabited environment, and  $\Gamma \vdash t_1 t_2 : T$ , then there exists  $U, V$  such that  $\Gamma \vdash t_1 : U \rightarrow V$ .

A tempting formulation of the second effect safety statement would be that *in a pure and inhabited environment it's impossible to construct a term of the free function type*. However, this formulation has no hope to be proved, as  $S \rightarrow T$  is a subtype of  $S \Rightarrow T$ , any term that can be typed as the former can also be typed as the latter.

As in the system STLC-Pure, the proof of these two statements depends on two more general statements about *capsafe environments*. If we can arrive at such a definition of *capsafe environment* that a pure and inhabited environment is also *capsafe*, then it suffices to prove following two statements about *capsafe environments*:

DEFINITION (Effect-Safety-1). If  $\Gamma$  is *capsafe*, there doesn't exist  $t$  with  $\Gamma \vdash t : E$ .

DEFINITION (Effect-Safety-2). If  $\Gamma$  is *capsafe* and  $\Gamma \vdash t_1 t_2 : T$ , then there exists  $U, V$  such that  $\Gamma \vdash t_1 : U \rightarrow V$ .

Now we need to extend the definition of *capsafe* and *caprod* for free function types. Our first attempt is to add the following rule:

$$S \Rightarrow T \quad \text{caprod} \qquad \qquad \qquad (\text{CP-FUN2})$$

With this rule, the type  $(B \Rightarrow B) \rightarrow E$  would be considered *capsafe*, according to the rule CS-FUN1. However, with a variable  $f$  of this type and another variable  $g$  of the *capsafe* type  $B \rightarrow B$  in the *capsafe* environment, the constructed term  $fg$  has the capability type  $E$ ; the first statement of effect safety breaks. This formulation also breaks the connection between pure inhabited environments and *capsafe* environments. A pure and inhabited environment is no longer *capsafe*. For example, the type  $B \rightarrow B \Rightarrow B$  is pure and inhabited, but it's not *capsafe*.

On the other hand, we cannot do the opposite to take  $S \Rightarrow T$  as *capsafe*, as it would allow calling free functions in *capsafe* environments; the second statement of effect safety breaks. In the meanwhile, the pure and inhabited type  $(B \Rightarrow E) \rightarrow E$  is not *capsafe*, thus a pure and inhabited environment is no longer a *capsafe* environment.

This dilemma prompts us to reexamine the meaning of *capsafe* and *caprod*. When these facilities were first introduced in STLC-Pure, they were formulated in terms of the provability of the capability type  $E$ . From the perspective of the provability of  $E$ ,  $S \rightarrow T$  and  $S \Rightarrow T$  don't make much difference. Thus, free function types should be formulated the same way as stoic function

types. We tried this approach, and it worked<sup>8</sup>. The full formulation is presented in Figure 3.4.

<b>Capsafe Type</b>	<b>Caprod Type</b>
$B \text{ capsafe} \quad (\text{CS-BASE})$	$E \text{ caprod} \quad (\text{CP-EFF})$
$\frac{S \text{ caprod}}{S \rightarrow T \text{ capsafe}} \quad (\text{CS-FUN1})$	$\frac{S \text{ capsafe} \quad T \text{ caprod}}{S \rightarrow T \text{ caprod}} \quad (\text{CP-FUN1})$
$\frac{T \text{ capsafe}}{S \rightarrow T \text{ capsafe}} \quad (\text{CS-FUN2})$	$\frac{S \text{ capsafe} \quad T \text{ caprod}}{S \Rightarrow T \text{ caprod}} \quad (\text{CP-FUN2})$
$\frac{S \text{ caprod}}{S \Rightarrow T \text{ capsafe}} \quad (\text{CS-FUN3})$	<b>Capsafe Environment</b>
$\frac{T \text{ capsafe}}{S \Rightarrow T \text{ capsafe}} \quad (\text{CS-FUN4})$	$\emptyset \text{ capsafe} \quad (\text{CE-EMPTY})$
	$\frac{\Gamma \text{ capsafe} \quad T \text{ capsafe}}{\Gamma, x : T \text{ capsafe}} \quad (\text{CE-VAR})$

Figure 3.4: System STLC-Impure Capsafe Environment

Now a capsafe environment is no longer pure. For example, the type  $B \Rightarrow B$  is *capsafe*, but it's not pure. As a consequence, we need to slightly update the definition of the second statement of effect safety:

**DEFINITION (Effect-Safety-2')**. If  $\Gamma$  is pure and capsafe, and  $\Gamma \vdash t_1 t_2 : T$ , then there exists  $U, V$  such that  $\Gamma \vdash t_1 : U \rightarrow V$ .

Why this formulation of capsafe environment is acceptable? In short, it's because the statement *Effect-Safety-1* and *Effect-Safety-2'* logically imply the statement *Effect-Safety-Inhabited-1* and *Effect-Safety-Inhabited-2* respectively.

The logical implications hold because a pure and inhabited environment is also a capsafe (and pure) environment. This claim has been formally proved:

**LEMMA (Inhabited-Capsafe)**. If the type  $T$  is inhabited, then either  $T$  is capsafe or  $T = E$  or  $T$  is a free function type.

<sup>8</sup>Sandro Stucki suggested to take this approach.

THEOREM (Inhabited-Pure-Capsafe-Env). If  $\Gamma$  is pure and inhabited, then  $\Gamma$  is capsafe.

COROLLARY (Inhabited-Pure-Capsafe-Env'). If  $\Gamma$  is pure and inhabited, then  $\Gamma$  is pure and capsafe.

Note that the last corollary follows immediately from the second theorem, as we already know from the premise that  $\Gamma$  is pure.

### 3.3.2 Proof

The proof of the first statement of effect safety is almost the same as in STLC-Pure. The only change worth mention is that, in the presence of subtyping, an additional lemma *Capsafe-Sub* is required. The first effect safety statement follows immediately from the lemma *Capsafe-Env-Capsafe*.

LEMMA (Capsafe-Not-Caprod). If type  $T$  is capsafe, then  $T$  is not caprod.

LEMMA (Capsafe-Or-Caprod). For any  $T$ ,  $T$  is either capsafe or caprod.

LEMMA (Capsafe-Sub). If  $S$  is capsafe and  $S <: T$ , then  $T$  is capsafe.

LEMMA (Capsafe-Env-Capsafe). If  $\Gamma$  is capsafe and  $\Gamma \vdash t : T$ , then  $T$  is capsafe.

THEOREM (Effect-Safety-1). If  $\Gamma$  is capsafe, then there doesn't exist term  $t$  with  $\Gamma \vdash t : E$ .

However, it's impossible to prove the second statement of effect safety, which states that if we can construct an application  $t_1 t_2$  in a pure and capsafe environment, then  $t_1$  can be typed as  $S \rightarrow T$  for some  $S$  and  $T$ . To see an example why it cannot be proved, let's assume  $\Gamma = \{f : B \rightarrow B \Rightarrow B, x : B\}$ . It's obvious  $\Gamma$  is capsafe and pure. However,  $f x$  in the term  $(f x) x$  has the type  $B \Rightarrow B$ . It's impossible to prove that  $f x$  has the type  $B \rightarrow B$ . If  $f$  is not a variable, but a fully defined function, then we can prove that  $f$  also has the type  $B \rightarrow B \rightarrow B$ . This is because the inner function of the type  $B \Rightarrow B$  cannot capture any capabilities or free functions. Otherwise, the outer function cannot be typed as stoic.

This observation leads us to assume four axioms listed in Figure 3.5. These axioms can only be proved if  $\Gamma$  is empty<sup>9</sup>. Otherwise, if  $t$  is a variable, we can do nothing.

The justification for the axiom AX-BASE is as follows. Suppose  $t = \lambda x:B. \lambda y:S. t_1$  and  $\Gamma \vdash t : B \rightarrow S \Rightarrow T$ . The typing rule for  $t$  should be the typing rule for stoic functions:

<sup>9</sup>Except the axiom AX-POLY, which cannot be proved even if  $\Gamma$  is empty.

$\frac{\Gamma \vdash t : B \rightarrow S \Rightarrow T}{\Gamma \vdash t : B \rightarrow S \rightarrow T} \text{ (AX-BASE)}$	$\frac{\Gamma \vdash t : (U \rightarrow V) \rightarrow S \Rightarrow T}{\Gamma \vdash t : (U \rightarrow V) \rightarrow S \rightarrow T} \text{ (AX-STOIC)}$
$\frac{\Gamma \vdash t : \text{Top} \rightarrow S \Rightarrow T}{\Gamma \vdash t : \text{Top} \rightarrow S \rightarrow T} \text{ (AX-TOP)}$	$\frac{\Gamma \vdash t_2 : U \rightarrow V \quad \Gamma \vdash t_1 : (U \Rightarrow V) \rightarrow S \Rightarrow T}{\Gamma \vdash t_1 t_2 : S \rightarrow T} \text{ (AX-POLY)}$

Figure 3.5: System STLC-Impure Axioms

$$\frac{\text{pure}(\Gamma), x : B \vdash \lambda y : S. t_1 : S \Rightarrow T}{\Gamma \vdash \lambda x : B. \lambda y : S. t_1 : B \rightarrow S \Rightarrow T} \text{ (STEP-1)}$$

Then what's the rule used in the typing of  $\lambda y : S. t_1$ ? If it's first typed as  $S \rightarrow T$  and then subsumed as  $S \Rightarrow T$ , we are done. Otherwise,  $\lambda y : S. t_1$  is typed using the rule of free functions:

$$\frac{\text{pure}(\Gamma), x : B, y : S \vdash t_1 : T}{\text{pure}(\Gamma), x : B \vdash \lambda y : S. t_1 : S \Rightarrow T} \text{ (STEP-2)}$$

According to the definition of *pure*, we know that following two equations hold:

$$\begin{aligned} \text{pure}(\text{pure}(\Gamma)) &= \text{pure}(\Gamma) \\ \text{pure}(\Gamma, x : B) &= \text{pure}(\Gamma), x : B \end{aligned}$$

Then we can obtain the following equation:

$$\text{pure}(\text{pure}(\Gamma), x : B) = \text{pure}(\Gamma), x : B$$

If we substitute the equation in Step-2, we get exactly the precondition for typing stoic functions. Thus  $\lambda y : S. t_2$  can be typed as stoic function:

$$\frac{\text{pure}(\text{pure}(\Gamma), x : B), y : S \vdash t_1 : T}{\text{pure}(\Gamma), x : B \vdash \lambda y : S. t_1 : S \rightarrow T} \text{ (STEP-2')}$$

Given that  $\lambda y : S. t_1$  can be typed as  $S \rightarrow T$ , Step-1 can be updated as follows:

$$\frac{\text{pure}(\Gamma), x : B \vdash \lambda y : S. t_1 : S \rightarrow T}{\Gamma \vdash \lambda x : B. \lambda y : S. t_1 : B \rightarrow S \rightarrow T} \text{ (STEP-1')}$$

Put all the steps together, we obtained  $\Gamma \vdash t : B \rightarrow S \rightarrow T$  from the fact  $\Gamma \vdash t : B \rightarrow S \Rightarrow T$ . That's the justification for the axiom AX-BASE. The justifications for AX-TOP and AX-STOIC are similar.

The axiom AX-POLY can be justified as follows. Suppose the following typing relation holds, what impure variables can be captured by  $t$ ?

$$\Gamma \vdash \lambda f:U \Rightarrow V. \lambda y:S. t : (U \Rightarrow V) \rightarrow S \Rightarrow T.$$

As the function is stoic, the only impure variable can be captured is  $f$ . Now, if we supply a stoic function as parameter to the function, it has the same effect as saying that  $f$  has the type  $U \rightarrow V$ . Thus, in this context, the function can be typed as  $(U \rightarrow V) \rightarrow S \Rightarrow T$ . Now according to the axiom AX-STOIC, the term can also be typed as  $(U \rightarrow V) \rightarrow S \rightarrow T$ . Finally by the rule T-APP, the type of the application is  $S \rightarrow T$ . That's the justification for the axiom AX-POLY.

Assuming the axioms above, it's straight-forward to prove a lemma *Capsafe-Pure-Stoic*, and the second statement of effect safety follows immediately from the lemma.

LEMMA (Capsafe-Pure-Stoic). If  $\Gamma$  is pure and capsafe, and  $\Gamma \vdash t : S \Rightarrow T$ , then  $\Gamma \vdash t : S \rightarrow T$ .

THEOREM (Effect-Safety-2'). If  $\Gamma$  is pure and capsafe, and  $\Gamma \vdash t_1 t_2 : T$ , then there exists  $U, V$  such that  $\Gamma \vdash t_1 : U \rightarrow V$ .

## 3.4 Effect Polymorphism

There are three kinds of effect polymorphism in capability-based effect systems, namely axiomatic polymorphism, curry polymorphism and stoic polymorphism. The first one depends on the axiom AX-POLY, while the other two are inherent in capability-based effect systems.

### 3.4.1 Axiomatic Polymorphism

Axiomatic polymorphism depends on the axiom AX-POLY. It can be illustrated by the following example:

```
def map[A,B](f: A => B)(l: List[A]) = l match {
  case Nil => Nil
  case x::xs => f(x)::map(f)(xs)
}
def squareImpure(c: IO) = map { x => println(x)(c); x*x }
```



```
def squarePure = map { x => x*x }
```

In the code above, the function `map` has the type  $(A \Rightarrow B) \rightarrow \text{List}[A] \Rightarrow \text{List}[B]$ . Without the axiom AX-POLY, the function `squareImpure` would be typed as  $\text{IO} \rightarrow \text{List}[\text{Int}] \Rightarrow \text{List}[\text{Int}]$ , while the function `squarePure` would be typed as  $\text{List}[\text{Int}] \Rightarrow \text{List}[\text{Int}]$ . Effect polymorphism doesn't work in this case, as `squarePure` can actually be typed as a stoic function.

In the presence of the axiom AX-POLY, `squarePure` can be typed as  $\text{List}[\text{Int}] \rightarrow \text{List}[\text{Int}]$ , as the function passed to `map` has the stoic function type  $\text{Int} \rightarrow \text{Int}$ . Now both `squarePure` and `squareImpure` are typed as stoic functions, all effects are captured in the type system even in the presence of (local) free functions.

The reader might be wondering, what if we want to stipulate that the passed function `f` can only have IO effects? How does effect polymorphism work in such cases? The answer is as follows:

```
private def mapImpl[A,B](f: A => B)(l: List[A]) = l match {  
  case Nil => Nil  
  case x::xs => f(x)::map(f)(xs)  
}
```

```
def map[A,B](f: A -> B) = mapImpl(f)
```

```
def map[A,B](c: IO)(f: IO -> A => B) = mapImpl(f c)
```

In the code above, we protect the implementation function `mapImpl` with the keyword `private`. All callings of the `map` function has to go through the two exposed interfaces. The two interfaces impose that the passed function `f` must either be pure or only have IO side effects. It's straightforward to add a new interface to allow exception effects without changing the implementation or affecting existing code.

### 3.4.2 Currying Polymorphism

There is another kind of effect polymorphism inherent in capability-based effect systems without resorting to the axiom AX-POLY. This kind of polymorphism is related to *currying*, thus is called *currying polymorphism*. It can be demonstrated by the following example:

```
def map[A,B](f: A => B)(l: List[A]) = l match {  
  case Nil => Nil  
  case x::xs => f(x)::map(f)(xs)
```

```

}
def squareImpure(c: IO) = map { x => println(x)(c); x*x }
def squarePure(l: List[Int]) = map { x => x*x } l

```

Note that the only change we made compared to the example in *axiomatic polymorphism* is to add the parameter `l : List[Int]` explicitly to `squarePure`, instead of resorting to currying. This seemingly trivial change makes a big difference in capability-based effect systems. Now `squarePure` can be typed as `List[Int] → List[Int]` without using the axiom AX-POLY.

### 3.4.3 Stoic Polymorphism

Stoic functions that take a free function and return a value of a pure type are inherently effect-polymorphic. This kind of effect polymorphism is called *stoic polymorphism*. Stoic polymorphism can be illustrated by the following example:

```

def twice(f: Int => Int) = f (f 0)
def pure(x: Int) = twice { n => n + x }
def impure(x: Int)(c: IO) = twice { n => println(n)(c); n + x }

```

In the code above, the function `twice` is typed as `(Int ⇒ Int) → Int`. By just checking the type signature of `twice`, we know it might have side effects, if the passed function `f` has side effects. Without adding any annotation or resorting to any axiom, the type system automatically types the function `pure` as `Int → Int` and `impure` as `Int → IO → Int`. Effect polymorphism works naturally.

### 3.4.4 Discussion

As reported in section 1.6 of the thesis[Lip09a], Haskell has fractured into monadic and non-monadic sub-languages. In Haskell, almost every general purpose higher-order function needs both a monadic version and a non-monadic version. For example, in Haskell it requires two versions of the function `map`, but in capability-based effect systems, only one effect-polymorphic `map` is required, as the following code shows:

```

-- monad-based effect system
map :: (a -> b) -> List a -> List b
mapM :: Monad m => (a -> m b) -> List a -> m (List b)
-- capability-based effect system

```

```
map :: (a => b) -> List a => List b
```

In Haskell, it's possible to implement the non-monadic version based on the monadic version by using the identity monad:

```
map :: (a -> b) -> List a -> List b
map f xs = runIdentity (mapM (\x -> return (f x)) xs)
```

However, in practice programmers usually first come up with the non-monadic version, and later realize the need for a monadic version. Turning the non-monadic code into the monadic version requires almost a rewrite of the function, as demonstrated by the following code:

```
map :: (a -> b) -> List a -> List b
map f xs
  = case xs of
      Nil -> Nil
      Cons x xs -> Cons (f x) (map f xs)

mapM :: Monad m => (a -> m b) -> List a -> m (List b)
mapM f xs
  = case xs of
      Nil -> return Nil
      Cons x xs -> do x' <- f x
                    xs' <- mapM f xs
                    return (Cons x' xs')
```

In capability-based effect systems, if programmers first come up with the pure version, turning it to the effect-polymorphic version only requires changing  $\rightarrow$  to  $\Rightarrow$ , as demonstrated by the following code:

```
// pure version
def map[A,B](f: A -> B)(l: List[A]) = l match {
  case Nil => Nil
  case x::xs => f(x)::map(f)(xs)
}

// effect-polymorphic version
def map[A,B](f: A => B)(l: List[A]) = l match {
  case Nil => Nil
  case x::xs => f(x)::map(f)(xs)
}
```

Monads are not only heavier in handling effect polymorphism, but also more difficult to learn and use than capabilities. The concept *capability* has an intuitive meaning as it has in daily life, while the abstruse concept *monad* comes from category theory, which is the most abstract field of mathematics. Therefore, we expect capability-based effect systems to be welcomed by a larger audience than monad-based effect systems.

Type-and-effect systems based on type-and-effect inference[TJ92, TJ94] can greatly reduce the syntactical overhead in effect systems, but it still complicates the type signature of effect-polymorphic functions by introducing generic effect type variables. For example, the type signature inferred for the effect-polymorphic function `map` would look like  $\forall E. (A \rightarrow B@E) \rightarrow \text{List}[A] \rightarrow \text{List}[B]@E$ , which is a little daunting to programmers, compared to  $(A \Rightarrow B) \rightarrow \text{List}[A] \Rightarrow \text{List}[B]$  in our case.

Effect polymorphism is an advantage of capability-based effect systems over other kinds of effect systems. In capability-based effect systems, effect polymorphism can be easily achieved without introducing generic effect type variables in the type signature or depending on additional annotations.

## 4 System F-Pure

In this chapter, we present the system F-Pure, which is an extension of the system STLC-Pure with universal types. In this system, not only functions need to observe a variable-capturing discipline, type abstractions also need to observe the same variable-capturing discipline.

We'll first introduce the formalization, then discuss soundness and effect safety. In the discussion, we'll focus on its difference from the system STLC-Pure.

### 4.1 Definitions

The system F-Pure extends STLC-Pure with universal types. Figure 4.1 presents the full definition of F-Pure, with the difference from the system STLC-Pure highlighted.

The extension of syntax and evaluation rules are exactly the same as the extension of standard STLC with universal types. The essential difference lies in the two new typing rules T-TABS and T-TAPP. The typing rule T-TABS stipulates that type abstraction must observe the variable-capturing discipline.

$$\frac{\text{pure}(\Gamma), X \vdash t_2 : T}{\Gamma \vdash \lambda X. t_2 : \forall X. T} \quad (\text{T-TABS})$$

We made this design choice in order to allow universal types to be present in pure environments. Otherwise, if type abstractions can capture capability variables, application of a type abstraction could generate a term of the capability type or have side effects. This makes it incorrect to have universal types in pure environments, thus renders universal types useless in the system.

The typing rule T-TAPP requires that the type parameter cannot be the capability type E. However, it's allowed to supply uninhabited types like  $B \rightarrow E$  as parameter to type abstraction.

$$\frac{\Gamma \vdash t_1 : \forall X. T \quad T_2 \neq E}{\Gamma \vdash t_1 [T_2] : [X \mapsto T_2]T} \quad (\text{T-TAPP})$$

Without the restriction, preservation of the system breaks <sup>10</sup>, as can be seen from the following term  $t$ , which has the type  $\forall X. X \rightarrow B \rightarrow X$ :

<sup>10</sup>It's also possible to allow E as argument to type application, and restore preservation by treating variables with the type of a type variable as impure. We'll follow this approach in the system F-Impure.



$$t = \lambda X. \lambda x : X. \lambda y : B. x$$

If we allow  $E$  as parameter to type application, the term  $t [E]$  would have the type  $E \rightarrow B \rightarrow E$ . However, after one evaluation step, the term  $\lambda x : E. \lambda y : B. x$  cannot be typed anymore, as the capability variable  $x$  cannot be captured in the inner-most lambda; thus preservation breaks.

The definition of the function `pure` is changed slightly by allowing type variables to be in the pure environment. Type variables themselves are harmless to effect safety, their presence in the pure environment ensures the well-formedness of the environment.

Note that a hidden change to the function `pure` is that the binding  $x : X$ , where  $X$  is a type variable, may appear in pure environments. This is natural, as we know in the typing rule T-TAPP that the type variable  $X$  cannot be the capability type  $E$ .

## 4.2 Soundness

We proved both progress and preservation of the system.

THEOREM (Progress). If  $\emptyset \vdash t : T$ , then either  $t$  is a value or there is some  $t'$  with  $t \longrightarrow t'$ .

THEOREM (Preservation). If  $\Gamma \vdash t : T$ , and  $t \longrightarrow t'$ , then  $\Gamma \vdash t' : T$ .

The proof of progress is the same as in System F. In the proof of preservation, we need to make small changes to the standard substitution lemmas in System F.

LEMMA (Substitution-Term). If  $\Gamma, x : S \vdash t : T$ ,  $s$  is a value and  $\Gamma \vdash s : S$ , then  $\Gamma \vdash [x \mapsto s]t : T$ .

LEMMA (Substitution-Type). If  $\Gamma, X \vdash t : T$  and  $P \neq E$ , then  $\Gamma \vdash [X \mapsto P]t : [X \mapsto P]T$ .

We restrict  $s$  to be a value in the lemma *Substitution-Term* for the same reason as in the system STLC-Pure. In the lemma *Substitution-Type*, we restrict that  $P$  is not the capability type  $E$ . Otherwise, the lemma cannot be proved as explained in the previous section.

## 4.3 Effect Safety

We follow the same approach as in the system STLC-Pure in the formulation of effect safety.

### 4.3.1 Inhabited Types and Environments

In the presence of universal types, we need to adapt the definition of inhabited type and inhabited environment.

The definition of inhabited environment in STLC-Pure would reject a well-formed environment with type variables like  $\{T, S, f : T \rightarrow S, x : T\}$ , which could be the environment for the following well-defined function:

```
def apply[T, S](f: T -> S)(x: T) = f x
```

To handle this problem, we propose the definition of *inhabited type* and *inhabited environment* as shown in Figure 4.2.

<b>Primitive Environment</b>	<b>Inhabited Environment</b>
$x : B, y : E$ primitive (P-BASE)	$\emptyset$ inhabited (IE-EMPTY)
$\frac{\Sigma \text{ primitive}}{\Sigma, X \text{ primitive}} \quad (\text{P-TVAR})$	$\frac{\Gamma \text{ inhabited}}{\Gamma, X \text{ inhabited}} \quad (\text{IE-TVAR})$
$\frac{\Sigma \text{ primitive}}{\Sigma, x : X \text{ primitive}} \quad (\text{P-TYPE})$	$\frac{\Gamma \text{ inhabited} \quad T \text{ inhabited}}{\Gamma, x : T \text{ inhabited}} \quad (\text{IE-TYPE})$
<b>Inhabited Type</b>	
$\frac{\Sigma \text{ primitive} \quad \Sigma \vdash v : T}{T \text{ inhabited}} \quad (\text{IT})$	

Figure 4.2: System F-Pure Inhabited Environment

The definition of inhabited type depends on *primitive environments*. A primitive environment is an extension of the environment  $\{x : B, y : E\}$  with any type variables and type variable bindings. This definition would take types like  $X, X \rightarrow Y$  as inhabited, as expected.

The definition of inhabited environment not only allow bindings of inhabited types, but also allow any type variables to be present in an inhabited environment. This definition ensures that only uninhabited types, such as  $B \rightarrow E, \forall X.X$  and  $\forall X.\forall Y.X \rightarrow Y$  are rejected from pure environments.



### 4.3.2 Formulation

As in STLC, the standard formulation is given based on inhabited environments:

DEFINITION (Effect-Safety-Inhabited). If  $\Gamma$  is a pure and inhabited environment, then there doesn't exist  $t$  with  $\Gamma \vdash t : E$ .

The proof of the statement depends on a more general statement about *capsafe environments*. If we can arrive at such a definition of *capsafe environment* that a pure and inhabited environment is also *capsafe*, then it suffices to prove the following statement about *capsafe environments*:

DEFINITION (Effect-Safety). If  $\Gamma$  is *capsafe*, there doesn't exist  $t$  with  $\Gamma \vdash t : E$ .

What *capsafe* and *caprod* rules we need for universal types? Obviously, we need to take the uninhabited type  $\forall X.X$  as *caprod*, as with a variable of this type, it's possible to create a term of the capability type  $E$ . For example, if  $x$  is of the type  $\forall X.X$  and  $b$  is of the type  $B$ , then  $x [B \rightarrow E] b$  has the type  $E$ . We also need to take the uninhabited type  $\forall X.\forall Y.X \rightarrow Y$  as *caprod*. Otherwise, if  $x$  is of the type  $\forall X.\forall Y.X \rightarrow Y$  and  $b$  is of the type  $B$ , then  $x [B] [B \rightarrow E] b$  has the type  $E$ . This observation leads us to the definition of *capsafe environment* presented in Figure 4.3, with differences from STLC-Pure highlighted.

Why this formulation of *capsafe environment* is acceptable? In short, it's because the statement *Effect-Safety* logically implies the statement *Effect-Safety-Inhabited*.

The logical implication holds because a pure and inhabited environment is also a *capsafe environment*. This claim has been formally proved:

LEMMA (Inhabited-Capsafe). If the type  $T$  is inhabited, then either  $T$  is *capsafe* or  $T = E$ .

THEOREM (Inhabited-Pure-Capsafe-Env). If  $\Gamma$  is pure and inhabited, then  $\Gamma$  is also *capsafe*.

### 4.3.3 Proof

The proof of effect safety is more involved than in STLC-Pure. We need to introduce the *degree* of types in the proof of relevant lemmas about types.

<b>Capsafe Type</b>	<b>Caprod Type</b>
$B \text{ capsafe} \quad (\text{CS-BASE})$	$E \text{ caprod} \quad (\text{CP-EFF})$
$X \text{ capsafe} \quad (\text{CS-VAR})$	$\frac{S \text{ capsafe} \quad T \text{ caprod}}{S \rightarrow T \text{ caprod}} \quad (\text{CP-FUN})$
$\frac{S \text{ caprod}}{S \rightarrow T \text{ capsafe}} \quad (\text{CS-FUN1})$	$\frac{[X \mapsto B]T \text{ caprod}}{\forall X.T \text{ caprod}} \quad (\text{CP-ALL1})$
$\frac{T \text{ capsafe}}{S \rightarrow T \text{ capsafe}} \quad (\text{CS-FUN2})$	$\frac{[X \mapsto E]T \text{ caprod}}{\forall X.T \text{ caprod}} \quad (\text{CP-ALL2})$
$\frac{[X \mapsto B]T \text{ capsafe} \quad [X \mapsto E]T \text{ capsafe}}{\forall X.T \text{ capsafe}} \quad (\text{CS-ALL})$	<b>Capsafe Environment</b>
	$\emptyset \text{ capsafe} \quad (\text{CE-EMPTY})$
	$\frac{\Gamma \text{ capsafe} \quad T \text{ capsafe}}{\Gamma, x : T \text{ capsafe}} \quad (\text{CE-VAR})$
	$\frac{\Gamma \text{ capsafe}}{\Gamma, X \text{ capsafe}} \quad (\text{CE-TVAR})$

Figure 4.3: System F-Pure Capsafe Environment

DEFINITION (Degree of Type). The degree of a type  $T$  is defined as follows:

$$\text{degree}(T) = \begin{cases} \max(\text{degree}(t_1), \text{degree}(t_2)) & \text{if } T = T_1 \rightarrow T_2, \\ \text{degree}(T_1) + 1 & \text{if } T = \forall X.T_1, \\ 0 & \text{others} \end{cases}$$

With the help of the definition above, it's possible to prove following lemmas based on a nested induction on the degree of types and the type  $T$ .

LEMMA (Capsafe-Not-Caprod). If type  $T$  is capsafe, then  $T$  is not caprod.

LEMMA (Capsafe-Or-Caprod). For any type  $T$ ,  $T$  is either capsafe or caprod.

LEMMA (Capsafe-All-Subst). If  $\forall X. T$  is capsafe, then for any type  $U$ ,  $[X \mapsto U]T$  is capsafe.

To prove the lemma *Capsafe-Env-Capsafe*, we need a similar definition on terms, and then do a nested induction on the degree of terms and the typing relation.

DEFINITION (Degree of Term). The degree of a term  $t$  is defined as follows:

$$\text{degree}(t) = \begin{cases} \text{degree}(t_1) & \text{if } t = \lambda x:T. t_1, \\ \max(\text{degree}(t_1), \text{degree}(t_2)) & \text{if } t = t_1 t_2, \\ \text{degree}(t_1) + 1 & \text{if } t = \lambda X. t_1, \\ \text{degree}(t_1) & \text{if } t = t_1 [T], \\ 0 & \text{others} \end{cases}$$

Effect safety follows immediately from the lemma *Capsafe-Env-Capsafe*.

LEMMA (Capsafe-Env-Capsafe). If  $\Gamma$  is capsafe and  $\Gamma \vdash t : T$ , then  $T$  is capsafe.

THEOREM (Effect-Safety). If  $\Gamma$  is capsafe, then there doesn't exist  $t$  with  $\Gamma \vdash t : E$ .

## 5 System F-Impure

In this chapter, we present the system F-Impure, which is an extension of the system F-Pure with free functions. It can also be seen as an extension of the system STLC-Impure with universal types, but without subtyping. Extending the system with subtyping would lead us to bounded quantification, which we are still working on. Given the importance of parametric polymorphism and the fact that subtyping is not a necessary add-on of functional programming, the system F-Impure deserves a separate presentation here.

We'll first introduce the formalization, then discuss soundness and effect safety. In the discussion, we'll focus on its difference from the system STLC-Impure and F-Pure.

### 5.1 Definitions

Figure 5.1 presents the full definition of F-Impure, with the difference from the system F-Impure highlighted. As can be seen from the figure, we introduced free function types and added a typing rule for free functions. We have to add a typing rule T-DEGEN to restore the subtyping relation between stoic function types and free function types, as there is no subtyping in the current system.

As in STLC-Impure, we adapted the definition of pure to exclude free function types from pure environments. If stoic functions have access to free functions, we'll lose the ability to track the effects of stoic functions in the type system.

Different from F-Pure, in the rule T-TAPP we allow both capability types and free function types as type parameter in type application. We know in F-pure, without such restriction preservation of the system breaks. To restore preservation of the system, in the definition of pure we remove  $x : X$  from the pure environment. That is, we treat variables typed with type variable as impure, which cannot be captured by a stoic function.

We could also introduce free type abstractions in the system. Such an extension will not be very useful in real-world programming, as in practice polymorphic functions rarely capture free variables, not mention capability variables. Thus for the sake of simplicity, we don't pursue the extension.

<p><b>Syntax</b></p> <p><math>t ::=</math></p> <table style="width: 100%; border: none;"> <tr> <td style="width: 30%;"><math>x</math></td> <td style="width: 30%;"></td> <td style="width: 40%;">terms: variable</td> </tr> <tr> <td><math>\lambda x:T. t</math></td> <td></td> <td>abstraction</td> </tr> <tr> <td><math>tt</math></td> <td></td> <td>application</td> </tr> <tr> <td><math>\lambda X. t</math></td> <td></td> <td>type abstraction</td> </tr> <tr> <td><math>t[T]</math></td> <td></td> <td>type application</td> </tr> </table> <p><math>v ::=</math></p> <table style="width: 100%; border: none;"> <tr> <td style="width: 30%;"><math>\lambda x:T. t</math></td> <td style="width: 30%;"></td> <td style="width: 40%;">values: abstraction value</td> </tr> <tr> <td><math>x</math></td> <td></td> <td>variable value</td> </tr> <tr> <td><math>\lambda X. t</math></td> <td></td> <td>type abstraction value</td> </tr> </table> <p><math>T ::=</math></p> <table style="width: 100%; border: none;"> <tr> <td style="width: 30%;"><math>X</math></td> <td style="width: 30%;"></td> <td style="width: 40%;">types: type variable</td> </tr> <tr> <td><math>B</math></td> <td></td> <td>basic type</td> </tr> <tr> <td><math>E</math></td> <td></td> <td>capability type</td> </tr> <tr> <td><math>T \rightarrow T</math></td> <td></td> <td>type of stoic funs</td> </tr> <tr> <td><math>T \Rightarrow T</math></td> <td></td> <td>type of free funs</td> </tr> <tr> <td><math>\forall X. T</math></td> <td></td> <td>universal type</td> </tr> </table> <p><b>Evaluation</b></p> <div style="text-align: right; border: 1px solid black; padding: 2px; display: inline-block;"><math>t \longrightarrow t'</math></div> $\frac{t_1 \longrightarrow t'_1}{t_1 t_2 \longrightarrow t'_1 t_2} \quad (\text{E-APP1})$ $\frac{t_2 \longrightarrow t'_2}{v_1 t_2 \longrightarrow v_1 t'_2} \quad (\text{E-APP2})$ $(\lambda x:T. t_1)v_2 \longrightarrow [x \mapsto v_2]t_1 \quad (\text{E-APPABS})$ $\frac{t_1 \longrightarrow t'_1}{t_1 [T_2] \longrightarrow t'_1 [T_2]} \quad (\text{E-TAPP})$ $(\lambda X. t_1)[T_2] \longrightarrow [X \mapsto T_2]t_1 \quad (\text{E-TAPPTABS})$	$x$		terms: variable	$\lambda x:T. t$		abstraction	$tt$		application	$\lambda X. t$		type abstraction	$t[T]$		type application	$\lambda x:T. t$		values: abstraction value	$x$		variable value	$\lambda X. t$		type abstraction value	$X$		types: type variable	$B$		basic type	$E$		capability type	$T \rightarrow T$		type of stoic funs	$T \Rightarrow T$		type of free funs	$\forall X. T$		universal type	<p><b>Typing</b></p> <div style="text-align: right; border: 1px solid black; padding: 2px; display: inline-block;"><math>\Gamma \vdash x : T</math></div> $\frac{x : T \in \Gamma}{\Gamma \vdash x : T} \quad (\text{T-VAR})$ $\frac{\text{pure}(\Gamma), x : S \vdash t_2 : T}{\Gamma \vdash \lambda x:S. t_2 : S \rightarrow T} \quad (\text{T-ABS1})$ $\frac{\Gamma, x : S \vdash t_2 : T}{\Gamma \vdash \lambda x:S. t_2 : S \Rightarrow T} \quad (\text{T-ABS2})$ $\frac{\Gamma \vdash t : S \rightarrow T}{\Gamma \vdash t : S \Rightarrow T} \quad (\text{T-DEGEN})$ $\frac{\Gamma \vdash t_1 : S \Rightarrow T \quad \Gamma \vdash t_2 : S}{\Gamma \vdash t_1 t_2 : T} \quad (\text{T-APP})$ $\frac{\text{pure}(\Gamma), X \vdash t_2 : T}{\Gamma \vdash \lambda X. t_2 : \forall X. T} \quad (\text{T-TABS})$ $\frac{\Gamma \vdash t_1 : \forall X. T}{\Gamma \vdash t_1 [T_2] : [X \mapsto T_2]T} \quad (\text{T-TAPP})$ <p><b>Pure Environment</b></p> <table style="width: 100%; border: none;"> <tr> <td><math>\text{pure}(\emptyset)</math></td> <td><math>=</math></td> <td><math>\emptyset</math></td> </tr> <tr> <td><math>\text{pure}(\Gamma, x : E)</math></td> <td><math>=</math></td> <td><math>\text{pure}(\Gamma)</math></td> </tr> <tr> <td><math>\text{pure}(\Gamma, x : X)</math></td> <td><math>=</math></td> <td><math>\text{pure}(\Gamma)</math></td> </tr> <tr> <td><math>\text{pure}(\Gamma, x : S \Rightarrow T)</math></td> <td><math>=</math></td> <td><math>\text{pure}(\Gamma)</math></td> </tr> <tr> <td><math>\text{pure}(\Gamma, x : T)</math></td> <td><math>=</math></td> <td><math>\text{pure}(\Gamma), x : T</math></td> </tr> <tr> <td><math>\text{pure}(\Gamma, X)</math></td> <td><math>=</math></td> <td><math>\text{pure}(\Gamma), X</math></td> </tr> </table>	$\text{pure}(\emptyset)$	$=$	$\emptyset$	$\text{pure}(\Gamma, x : E)$	$=$	$\text{pure}(\Gamma)$	$\text{pure}(\Gamma, x : X)$	$=$	$\text{pure}(\Gamma)$	$\text{pure}(\Gamma, x : S \Rightarrow T)$	$=$	$\text{pure}(\Gamma)$	$\text{pure}(\Gamma, x : T)$	$=$	$\text{pure}(\Gamma), x : T$	$\text{pure}(\Gamma, X)$	$=$	$\text{pure}(\Gamma), X$
$x$		terms: variable																																																											
$\lambda x:T. t$		abstraction																																																											
$tt$		application																																																											
$\lambda X. t$		type abstraction																																																											
$t[T]$		type application																																																											
$\lambda x:T. t$		values: abstraction value																																																											
$x$		variable value																																																											
$\lambda X. t$		type abstraction value																																																											
$X$		types: type variable																																																											
$B$		basic type																																																											
$E$		capability type																																																											
$T \rightarrow T$		type of stoic funs																																																											
$T \Rightarrow T$		type of free funs																																																											
$\forall X. T$		universal type																																																											
$\text{pure}(\emptyset)$	$=$	$\emptyset$																																																											
$\text{pure}(\Gamma, x : E)$	$=$	$\text{pure}(\Gamma)$																																																											
$\text{pure}(\Gamma, x : X)$	$=$	$\text{pure}(\Gamma)$																																																											
$\text{pure}(\Gamma, x : S \Rightarrow T)$	$=$	$\text{pure}(\Gamma)$																																																											
$\text{pure}(\Gamma, x : T)$	$=$	$\text{pure}(\Gamma), x : T$																																																											
$\text{pure}(\Gamma, X)$	$=$	$\text{pure}(\Gamma), X$																																																											

Figure 5.1: System F-Impure

## 5.2 Soundness

We proved both progress and preservation of the system. The subtleties in the proof are the same as stated in the system F-Pure.

THEOREM (Progress). If  $\emptyset \vdash t : T$ , then either  $t$  is a value or there is some  $t'$  with  $t \longrightarrow t'$ .

THEOREM (Preservation). If  $\Gamma \vdash t : T$ , and  $t \longrightarrow t'$ , then  $\Gamma \vdash t' : T$ .

## 5.3 Effect Safety

We first introduce the formulation, which is a combination of the formulation in STLC-Impure and F-Pure, then discuss the proof of effect safety.

### 5.3.1 Formulation

The definitions of *inhabited type* and *inhabited environment* are the same as in the system F-Pure.

As in the system STLC-Impure, in the presence of free functions, we need two statements of effect safety:

DEFINITION (Effect-Safety-Inhabited-1). If  $\Gamma$  is a pure and inhabited environment, then there doesn't exist  $t$  with  $\Gamma \vdash t : E$ .

DEFINITION (Effect-Safety-Inhabited-2). If  $\Gamma$  is a pure and inhabited environment, and  $\Gamma \vdash t_1 t_2 : T$ , then there exists  $U, V$  such that  $\Gamma \vdash t_1 : U \rightarrow V$ .

As in the system STLC-Impure, the proof of these two statements depends on two more general statements about *capsafe environments*. Given that we've seen how universal types and free function types are extended in the formulation of *capsafe environment*, we can easily combine them to arrive at the formulation shown in Figure 5.2, with the changes from F-Pure highlighted. Notice that type variables are now treated as caprod (CP-VAR).

Why this formulation of capsafe environment is acceptable? In short, it's because the statement *Effect-Safety-Inhabited-1* and *Effect-Safety-Inhabited-2* are logically implied by the statement *Effect-Safety-1* and *Effect-Safety-2* respectively.

DEFINITION (Effect-Safety-1). If  $\Gamma$  is capsafe, there doesn't exist  $t$  with  $\Gamma \vdash t : E$ .

Capsafe Type	Caprod Type
$B \text{ capsafe} \quad (\text{CS-BASE})$	$E \text{ caprod} \quad (\text{CP-EFF})$
$\frac{S \text{ caprod}}{S \rightarrow T \text{ capsafe}} \quad (\text{CS-FUN1})$	$\frac{S \text{ capsafe} \quad T \text{ caprod}}{S \rightarrow T \text{ caprod}} \quad (\text{CP-FUN1})$
$\frac{T \text{ capsafe}}{S \rightarrow T \text{ capsafe}} \quad (\text{CS-FUN2})$	$\frac{S \text{ capsafe} \quad T \text{ caprod}}{S \Rightarrow T \text{ caprod}} \quad (\text{CP-FUN2})$
$\frac{S \text{ caprod}}{S \Rightarrow T \text{ capsafe}} \quad (\text{CS-FUN3})$	$\frac{[X \mapsto B]T \text{ caprod}}{\forall X.T \text{ caprod}} \quad (\text{CP-ALL1})$
$\frac{T \text{ capsafe}}{S \Rightarrow T \text{ capsafe}} \quad (\text{CS-FUN4})$	$\frac{[X \mapsto E]T \text{ caprod}}{\forall X.T \text{ caprod}} \quad (\text{CP-ALL2})$
$\frac{[X \mapsto B]T \text{ capsafe} \quad [X \mapsto E]T \text{ capsafe}}{\forall X.T \text{ capsafe}} \quad (\text{CS-ALL})$	<b>Capsafe Environment</b>
	$\emptyset \text{ capsafe} \quad (\text{CE-EMPTY})$
	$\frac{\Gamma \text{ capsafe} \quad T \text{ capsafe}}{\Gamma, x : T \text{ capsafe}} \quad (\text{CE-VAR})$
	$\frac{\Gamma \text{ capsafe}}{\Gamma, X \text{ capsafe}} \quad (\text{CE-TVAR})$

Figure 5.2: System F-Impure Capsafe Environment

DEFINITION (Effect-Safety-2). If  $\Gamma$  is pure and capsafe, and  $\Gamma \vdash t_1 t_2 : T$ , then there exists  $U, V$  such that  $\Gamma \vdash t_1 : U \rightarrow V$ .

The logical implications hold because a pure and inhabited environment is also a capsafe (and pure) environment. This claim has been formally proved:

LEMMA (Inhabited-Capsafe). If the type  $T$  is inhabited, then either  $T$  is capsafe or  $T = E$  or  $T$  is a free function type.

THEOREM (Inhabited-Pure-Capsafe-Env). If  $\Gamma$  is pure and inhabited, then  $\Gamma$  is also capsafe.

COROLLARY (Inhabited-Pure-Capsafe-Env'). If  $\Gamma$  is pure and inhabited, then  $\Gamma$  is pure and capsafe.

Note that the last corollary follows immediately from the second theorem, as we already know from the premise that  $\Gamma$  is pure.

### 5.3.2 Proof

The proof of the first effect safety statement is almost the same as in the system F-Pure, thus we omit here.

The proof of the second statement of effect safety faces the same problem as in the system STLC-Impure. We need to assume a set of axioms, as shown in Figure 5.3, with the newly added axioms highlighted. The justification for the axiom AX-ALL is the same as the justification for the axiom AX-BASE in STLC-Impure. In short, because the outer function is stoic and the type of the first parameter is pure, the inner function cannot capture variables of capabilities or free functions, thus it's fair enough to type the inner function as stoic.

The justification for the axiom AX-TABS is similar. If a term  $t$  can be typed as  $\forall X. T_1 \Rightarrow T_2$  under  $\Gamma$ , according to the typing rule T-ALL, the whole term can be typed under  $\text{pure}(\Gamma)$ . Then the inner function can be typed under  $\text{pure}(\Gamma), X$ , which is equal to  $\text{pure}(\Gamma, X)$ . Thus, it's fair enough to type the inner function as stoic.

The justification for the axiom AX-POLY is the same as given in STLC-Impure, thus we omit here.

Assuming these axioms, it's straight-forward to prove a lemma *Capsafe-Pure-Stoic*, and the second statement of effect safety follows immediately from the lemma.

LEMMA (Capsafe-Pure-Stoic). If  $\Gamma$  is pure and capsafe, and  $\Gamma \vdash t : S \Rightarrow T$ , then  $\Gamma \vdash t : S \rightarrow T$ .



$$\begin{array}{c}
\frac{\Gamma \vdash t : B \rightarrow S \Rightarrow T}{\Gamma \vdash t : B \rightarrow S \rightarrow T} \text{ (AX-BASE)} \qquad \frac{\Gamma \vdash t : \forall X. (T_1 \Rightarrow T_2)}{\Gamma \vdash t : \forall X. (T_1 \rightarrow T_2)} \text{ (AX-TABS)} \\
\frac{\Gamma \vdash t : (\forall X. T) \rightarrow S \Rightarrow T}{\Gamma \vdash t : (\forall X. T) \rightarrow S \rightarrow T} \text{ (AX-ALL)} \qquad \frac{\Gamma \vdash t : (U \rightarrow V) \rightarrow S \Rightarrow T}{\Gamma \vdash t : (U \rightarrow V) \rightarrow S \rightarrow T} \text{ (AX-STOIC)} \\
\frac{\Gamma \vdash t_2 : U \rightarrow V \quad \Gamma \vdash t_1 : (U \Rightarrow V) \rightarrow S \Rightarrow T}{\Gamma \vdash t_1 t_2 : S \rightarrow T} \text{ (AX-POLY)}
\end{array}$$

Figure 5.3: System F-Impure Axioms

**THEOREM (Effect-Safety-2).** If  $\Gamma$  is pure and capsafe, and  $\Gamma \vdash t_1 t_2 : T$ , then there exists  $U, V$  such that  $\Gamma \vdash t_1 : U \rightarrow V$ .

## 6 Conclusion

We formalized four capability-based effect systems and proved soundness and effect safety for each system. The four systems can serve as the theoretical foundation for implementing capability-based effect systems in functional languages.

The existence of *stoic functions* is the main trait of capability-based effect systems. The interplay between *stoic functions* and *free functions* enables flexible programming patterns that trivially solve the problem of *effect polymorphism*.

Capability-based effect systems have to be paired with strict evaluation, just like monad-based effect systems have to be paired with lazy evaluation.

### 6.1 Related Work

Lucassen and Gifford first introduced type-and-effect systems[GL86] and effect polymorphism using effect type parameterization [LG88], which is further developed by Talpin and Jouvelot to provide type-and-effect inference [TJ92, TJ94]. Type-and-effect inference can greatly reduce verbosity in syntax, but it only works in languages with global type inference, while Scala is based on local type inference. Even in those languages with global type-and-effect inference, the type signature for effect-polymorphic functions are much more complex than in capability-based effect systems, which is an obstacle to programmers.

Moggi introduced the usage of monads for computation effects[Mog91]. Wadler popularized the usage of monads[Wad92, Wad95] and proved that it's possible to transpose any type-and-effect system into a corresponding monad system[WT03]. As reported in section 1.6 of the thesis[Lip09a], almost every general purpose higher-order function in Haskell needs both a monadic version and non-monadic version. Lippmeier proposed the usage of region variables and dependently kinded witness to encode mutability polymorphism[Lip09b].

Lukas *et al.* studied type-and-effect systems for Scala[ROH12, RAO13, Ryt14]. In lightweight polymorphic effects[ROH12], the dichotomy between *effect-polymorphic function type* and *monomorphic function type* resembles the dichotomy between *stoic function type* and *free function type*. However, in the system effect polymorphism doesn't work for nested functions. To overcome this problem, they unified the two function types in a framework called *relative effect annotation* based on dependent types. In the framework, the function map can be marked as effect-polymorphic without introducing a generic effect variable:

```
def map[A, B](f: A => B)(l: List[A]): List[B] @pure(f)
```

The annotation `@pure(f)` says that the effect of `map` depends on the function `f`. Nested functions are handled specially in the system to avoid duplicate annotations.

Heather *et al.* proposed a variant of functions called *spores* for Scala[MHO14]. Compared to normal functions, spores observe a variable capturing discipline. The set of types that can or cannot be captured is part of the type signature of spores, thus can be used by library authors to impose constraints on parameters of function types.

Crary *et al.* proposed a capability calculus for typed memory management[CWM99], which improves the LIFO-style (last-in, first-out) region-based memory management introduced by Tofte and Talpin[TT97] by allowing arbitrary allocation and deallocation order. The safety of deallocation of memory is guaranteed by the type system. In the capability calculus, a capability is a set of regions that are presently valid to access. It tracks uniqueness of capabilities in the type system to control aliasing of capabilities. Functions may be polymorphic over types, regions or capabilities. Capability variables cannot be captured in closures and have to be passed around in the program.

## 6.2 Future Work

Recursive types, objects, mutation and general recursion are predominant features of industrial languages. To provide a more practical model for industrial languages, it's useful to extend existing systems with these features.

Effect masking can be a useful feature in real-world programming, especially when dealing with exception effects and local mutations. It's worthy to explore effect masking in existing systems.

Another direction of work is to extend the proposed systems with some implicit calculus, in order to avoid explicitly passing capabilities around in the code. This would result in significant savings in boilerplate code, thus make the system more friendly to programmers.

## References

- [Cha11] Arthur Charguéraud. The locally nameless representation. *Journal of Automated Reasoning*, pages 1–46, 2011. 10.1007/s10817-011-9225-2.
- [CLS07] Duncan Coutts, Roman Leshchinskiy, and Don Stewart. Stream fusion: From lists to streams to nothing at all. In *ACM SIGPLAN Notices*, volume 42, pages 315–326. ACM, 2007.
- [CWM99] Karl Crary, David Walker, and Greg Morrisett. Typed memory management in a calculus of capabilities. In *Proceedings of the 26th ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, pages 262–275. ACM, 1999.
- [GL86] David K Gifford and John M Lucassen. Integrating functional and imperative programming. In *Proceedings of the 1986 ACM conference on LISP and functional programming*, pages 28–38. ACM, 1986.
- [LG88] John M Lucassen and David K Gifford. Polymorphic effect systems. In *Proceedings of the 15th ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, pages 47–57. ACM, 1988.
- [Lip09a] Ben Lippmeier. *Type inference and optimisation for an impure world*. Australian National University, 2009.
- [Lip09b] Ben Lippmeier. Witnessing purity, constancy and mutability. In *Programming Languages and Systems*, pages 95–110. Springer, 2009.
- [MHO14] Heather Miller, Philipp Haller, and Martin Odersky. Spores: a type-based foundation for closures in the age of concurrency and distribution. In *ECOOP 2014–Object-Oriented Programming*, pages 308–333. Springer, 2014.
- [Mog91] Eugenio Moggi. Notions of computation and monads. *Information and computation*, 93(1):55–92, 1991.
- [Pie02] Benjamin C Pierce. *Types and programming languages*. MIT press, 2002.
- [RAO13] Lukas Rytz, Nada Amin, and Martin Odersky. A flow-insensitive, modular effect system for purity. In *Proceedings of the 15th Workshop on Formal Techniques for Java-like Programs*, page 4. ACM, 2013.
- [ROH12] Lukas Rytz, Martin Odersky, and Philipp Haller. Lightweight polymorphic effects. In *ECOOP 2012–Object-Oriented Programming*, pages 258–282. Springer, 2012.

- [Ryt14] Lukas Rytz. *A Practical Effect System for Scala*. EPFL, 2014.
- [TJ92] Jean-Pierre Talpin and Pierre Jouvelot. Polymorphic type, region and effect inference. *Journal of functional programming*, 2(03):245–271, 1992.
- [TJ94] Jean-Pierre Talpin and Pierre Jouvelot. The type and effect discipline. *Information and computation*, 111(2):245–296, 1994.
- [TT97] Mads Tofte and Jean-Pierre Talpin. Region-based memory management. *Information and computation*, 132(2):109–176, 1997.
- [Wad92] Philip Wadler. Comprehending monads. *Mathematical structures in computer science*, 2(04):461–493, 1992.
- [Wad95] Philip Wadler. Monads for functional programming. In *Advanced Functional Programming*, pages 24–52. Springer, 1995.
- [WT03] Philip Wadler and Peter Thiemann. The marriage of effects and monads. *ACM Transactions on Computational Logic (TOCL)*, 4(1):1–32, 2003.