

Modular Monadic Semantics and Compilation

Sheng Liang

Research Report YALEU/DCS/RR-1140
December 1997

Abstract

Modular Monadic Semantics and Compilation

Sheng Liang

Yale University

1998

Modular monadic semantics is a high-level and modular form of denotational semantics. It is capable of capturing individual programming language features and their interactions. This thesis explores the theory and applications of modular monadic semantics, including: building blocks for individual programming features, equational reasoning with laws and axioms, modular proofs, program transformation, modular interpreters, and compiler construction. We will demonstrate that the modular monadic semantics framework makes programming languages easy to specify, reason about, and implement.

Modular Monadic Semantics and Compilation

A Dissertation
Presented to the Faculty of the Graduate School
of
Yale University
in Candidacy for the Degree of
Doctor of Philosophy

by
Sheng Liang
May 1998

© Copyright by Sheng Liang 1997

All Rights Reserved

Acknowledgments

Most of all, I would like to thank my advisor, Paul Hudak, for his support and guidance. The Yale Haskell Project was the source of many of the ideas presented in this thesis. In particular, Mark Jones led me to the initial idea of modular monadic interpreters.

The work benefited from discussions and exchanges with Rajiv Mirani, Ross Paterson, Zhong Shao, Tim Sheard, and Dan Rabin.

Members of my reading committee, including Paul Hudak, Mark Jones, John Peterson, and Zhong Shao, helped me a great deal to improve the presentation of this thesis.

Finally, I would like to thank my family and friends for their support during the years at Yale.

Contents

List of Tables	iv
List of Figures	v
1 Introduction	1
1.1 Overview	1
1.2 The Lack of Modularity in Denotational Semantics	2
1.3 Monad — An Abstraction Mechanism	4
1.4 Background and Organization of the Thesis	6
1.5 The Source Language	7
1.6 A Notation	8
2 Modular Monadic Semantics	9
2.1 Modular Semantic Building Blocks	11
2.1.1 The Arithmetic Building Block	11
2.1.2 The Function Building Block	12
2.1.3 The References and Assignment Building Block	13
2.1.4 The Lazy Evaluation Building Block	14
2.1.5 The Program Tracing Building Block	15
2.1.6 The Continuation Building Block	15
2.1.7 The Nondeterminism Building Block	16
2.1.8 Alternative Definitions of the Environment and Store	17

2.2	Monads With Operations	18
2.3	Monad Transformers	19
2.3.1	The State Monad Transformer	21
2.3.2	The Environment Monad Transformer	22
2.3.3	The Error Monad Transformer	23
2.3.4	The Continuation Monad Transformer	23
2.3.5	The List Monad	24
2.4	Liftings	25
2.4.1	The Easy Cases	26
2.4.2	Lifting <i>Callcc</i>	27
2.4.3	Lifting <i>InEnv</i>	28
2.5	Summary	28
3	A Theory of Monads and Monad Transformers	33
3.1	Monad and Monad Transformers	34
3.1.1	Monads	34
3.1.2	Monad Transformers	35
3.2	Environment Axioms	36
3.3	Natural Liftings	37
3.3.1	Lifting Types	37
3.3.2	Natural Lifting Condition	38
3.3.3	Examples	40
3.4	Ordering of Monad Transformers	45
4	Modular Monadic Interpreters	47
4.1	Extensible Union Types	49
4.2	Interpreter Building Blocks	50
4.3	Constructor Classes	52
4.4	Monads	53

4.5	Monad Transformers	54
4.6	Summary	56
5	Compilation	57
5.1	Using Monad Laws to Transform Programs	58
5.2	A Natural Semantics	60
5.2.1	Definition of a Natural Semantics	60
5.2.2	Correspondence between Natural and Standard Semantics	61
5.2.3	Benefits of Reasoning in Monadic Style	62
5.3	Targeting Monadic Code	63
5.3.1	The Target Language Monad	63
5.3.2	Utilizing Target Language Features	64
5.3.3	Limitations of This Approach	65
5.3.4	Implications for a Common Back-end	65
5.4	An Experiment: Retargeting a Haskell Compiler	66
5.4.1	The STG Language	66
5.4.2	Compiling the STG Language	68
5.4.3	Monadic Semantics of the STG Language	70
5.4.4	Implementing Standard Entries	76
5.4.5	Connecting to the SML/NJ Back-end	78
6	Related Work, Future Work and Conclusion	83
6.1	Related Work	84
6.1.1	Modular Semantics	84
6.1.2	Reasoning with Monads	85
6.1.3	Semantics-directed Compilation	86
6.2	Future Work	87
6.2.1	Theory of Programming Language Features	87
6.2.2	Monadic Program Transformation	87

6.2.3	A Common Back-end for Modern Languages	87
6.2.4	Concurrency	88
6.3	Conclusion	88
	Bibliography	89
	A Proofs	95
	B Gofer Code for Monad Transformers	113
	C Monadic Semantics of the STG Language	117
C.1	Literals, Variables and Atoms	117
C.2	Function Applications	118
C.3	Constructors and Primitives	119
C.4	Case Expressions	120
C.5	Let Bindings	121
C.6	Right-hand Sides	123

List of Tables

2.1	Monad operations used in the semantics	18
5.1	Sample programs	80
5.2	Timing results of sample programs (in seconds)	80

List of Figures

1.1	A parameterized arithmetic semantics	5
2.1	The organization of modular monadic semantics	10
2.2	Monad transformers	29
2.3	Liftings	30
4.1	Gofer specification of a modular interpreter	48
5.1	The STG language	67
5.2	How an expression <code>(fac 3)</code> is evaluated in cell mode	70

Chapter 1

Introduction

1.1 Overview

Denotational semantics [Stoy, 1977] is among the most important developments in programming language theory. It gives a precise mathematical description of programming languages, useful in designing and implementing languages as well as reasoning about programs. For example, advances in denotational semantics have led to clarifications of features, to more consistent programming language design, and to new programming constructs.

It has long been recognized, however, that traditional denotational semantics lacks modularity and extensibility [Mosses, 1984] [Lee, 1989]. This is regarded as a major obstacle in applying denotational semantics to realistic programming languages.

In this thesis, we take advantage of a new development in programming language theory—a monadic approach [Moggi, 1990] to structured denotational semantics. The resulting *modular monadic semantics* achieves a high level of modularity and extensibility. It is able to capture individual programming language features in reusable building blocks, and to specify programming languages by composing the necessary features.

Because modular monadic semantics is no more than a structured denotational semantics, all the equational reasoning methods apply. In addition, we will show that modular monadic semantics further facilitates reasoning by allowing us to specify axioms of programming language features and to construct reusable modular proofs.

Modular monadic semantics can be implemented using modern programming languages such as Haskell [Hudak *et al.*, 1992], ML [Milner *et al.*, 1990], or Scheme [Clinger and Rees, 1991]. The result is a modular interpreter. We have discovered that the relatively new idea of *constructor classes* in Gofer (and Haskell 1.3) are particularly suitable for representing some rather complex typing relationships in modular interpreters.

Our work is directly applicable to semantics-directed compiler construction. We will present a compilation method based on monadic semantics and monadic program transformations, and test our ideas by retargeting a Haskell compiler.

Before introducing modular monadic semantics, in the next section we will give an example to demonstrate that traditional denotational semantics lacks modularity. The presentation follows the traditional denotational semantics style, augmented with a types declaration syntax similar to that of Haskell or ML. We assume the reader has basic knowledge of denotational semantics and functional programming.

1.2 The Lack of Modularity in Denotational Semantics

Let us first look at the semantics for a simple arithmetic language:

$$E : \text{Term} \rightarrow \text{Value}$$
$$E[[n]] = n$$
$$E[[e_1 + e_2]] = E[[e_1]] + E[[e_2]]$$

Denotational semantics maps *terms* in the source language into *values* in the meta language. The source language terms are enclosed in “ $\llbracket \ \rrbracket$ ”. The n and $+$ symbols on the right hand side correspond to the meta language concepts of a *number* and the *add* arithmetic operation.

An important measure of modularity is how a semantic description can be extended to incorporate new programming language features. For example, if we extend the source language with variables and functions, we need to introduce an environment—a mapping from variable names to values:

$$\begin{aligned}
 E & : \text{Term} \rightarrow \text{Env} \rightarrow \text{Value} \\
 E\llbracket n \rrbracket & = \lambda\rho.n \\
 E\llbracket e_1 + e_2 \rrbracket & = \lambda\rho.E\llbracket e_1 \rrbracket\rho + E\llbracket e_2 \rrbracket\rho \\
 E\llbracket v \rrbracket & = \lambda\rho.\rho\llbracket v \rrbracket
 \end{aligned}$$

Even though numbers are independent of the environment, we must change the semantics of numbers to accommodate the newly introduced environment argument. Similarly, the environment argument must be passed recursively to the subexpressions of $e_1 + e_2$, even though the arithmetic operation itself is independent of the environment.

If we further add continuations to our semantics (for supporting, for example, the sequencing operator “;”), we must change the semantics of numbers once again:

$$\begin{aligned}
 E & : \text{Term} \rightarrow \text{Env} \rightarrow (\text{Value} \rightarrow \text{Ans}) \rightarrow \text{Ans} \\
 E\llbracket n \rrbracket & = \lambda\rho.\lambda k.kn \\
 E\llbracket e_1 + e_2 \rrbracket & = \lambda\rho.\lambda k.E\llbracket e_1 \rrbracket\rho(\lambda i.E\llbracket e_2 \rrbracket\rho(\lambda j.k(i + j))) \\
 E\llbracket e_1; e_2 \rrbracket & = \lambda\rho.\lambda k.E\llbracket e_1 \rrbracket\rho(\lambda x.E\llbracket e_2 \rrbracket\rho k)
 \end{aligned}$$

In summary, we must make global changes to the traditional denotational semantics in order to add new features into the source language. The lack of modularity of traditional denotational semantics has long been recognized [Mosses, 1984]

[Lee, 1989]. It is regarded as the most significant obstacle in applying denotational semantics to realistic languages.

In the next section we will show how we can enhance the modularity by using an abstraction mechanism called *monads*.

1.3 Monad — An Abstraction Mechanism

We will use a type constructor M and two functions:

$$\begin{aligned} \text{return} & : a \rightarrow M a \\ \text{bind} & : M a \rightarrow (a \rightarrow M b) \rightarrow M b \end{aligned}$$

The intuitive meanings of these constructs are as follows:

- $M a$: a *computation* returning a value of type a .
- $\text{bind } c_1 (\lambda v. c_2)$: a computation that first computes c_1 , binds the result to v , and then computes c_2 .
- $\text{return } v$: a trivial computation that simply returns v as result.

Next, we write a *parameterized semantics* for arithmetic expressions using M , bind , and return :

$$\begin{aligned} E & : \text{Term} \rightarrow M \text{Value} \\ E[[n]] & = \text{return } n \\ E[[e_1 + e_2]] & = \text{bind } (E[[e_1]]) \\ & \quad (\lambda i. \text{bind } (E[[e_2]]) \\ & \quad \quad (\lambda j. \text{return } (i + j))) \end{aligned}$$

The semantic function E maps terms to computations (of type $M \text{Value}$). The semantics of $E[[n]]$ is a trivial computation that returns n as result. The semantics of $E[[e_1 + e_2]]$ is a computation that computes $E[[e_1]]$, binds the result to i , computes $E[[e_2]]$, binds the result to j , and finally returns $i + j$.

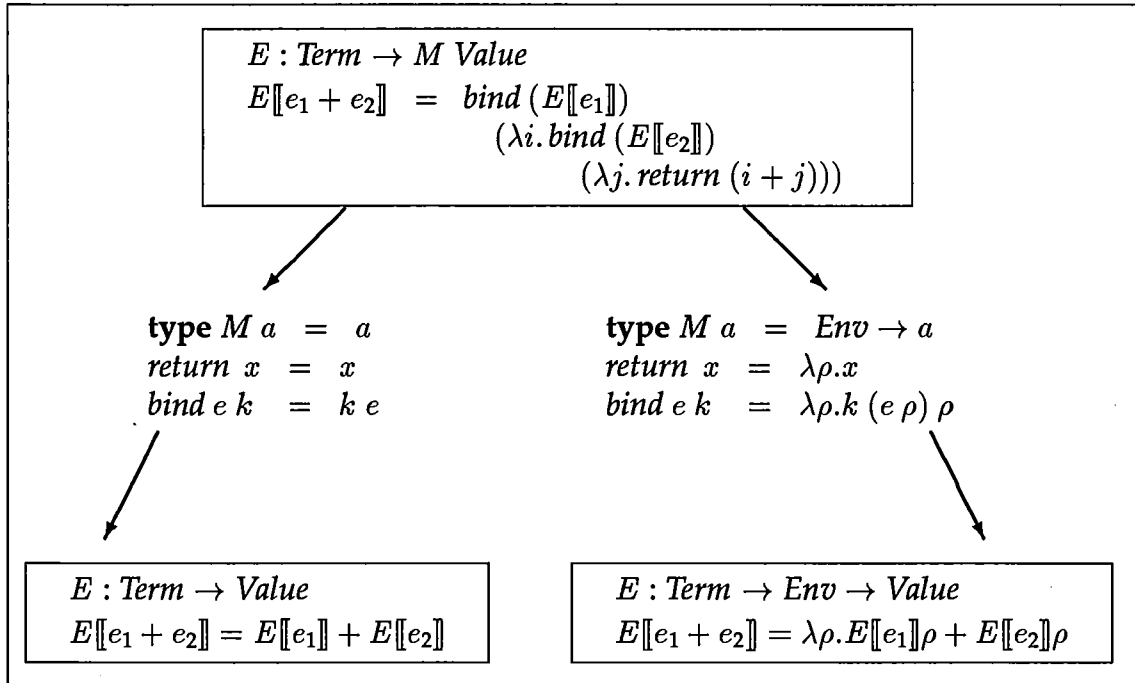


Figure 1.1: A parameterized arithmetic semantics

The advantage of the parameterized semantics is that, depending on how we instantiate M , return and bind , we get different semantics of our choice. Figure 1.1 shows how the arithmetic semantics can be instantiated to the trivial and environment-based semantics described in the last section. Later in this thesis, we will see that appropriate definitions of M , return and bind can also lead to other (for example, continuation-based) semantics.

The type constructor M , together with the two functions return and bind , are called a *monad*. The parameterized semantics defined above for arithmetic expressions is called *monadic semantics*. Monadic semantics can be instantiated using different *underlying monads*. To add a new feature to a monadic semantics, we only need to add a semantic description of the new feature, and change the underlying monad, but not the semantic descriptions of the existing features.

1.4 Background and Organization of the Thesis

This thesis explores the theory and practical applications of monads and monadic semantics, building on the previous work in this area. The concept of monads originates from category theory [Mac Lane, 1971]. The formulation of monads using a triple (*bind*, *return*, and the type constructor) is due to Kleisli. Moggi [Moggi, 1990] first proposed that monads provided a useful tool for structuring denotational semantics. Early work by Wadler [Wadler, 1990] showed the relationships between monads and functional programming. Recently, there has been a great deal of interest in using monads to construct modular semantics and modular interpreters [Wadler, 1992] [Jones and Duponcheel, 1993] [Espinosa, 1993] [Steele Jr., 1994].

In Chapter 2, we will present the modular monadic semantics for a wide range of programming language features. We will demonstrate how *monad transformers* capture individual features, and how *liftings* capture the interactions between different features.

In Chapter 3, we will investigate the theory of monads and monad transformers. This includes, for example, the formal properties of monad transformers and liftings. We will use monad laws and axioms to perform equational reasoning at a higher level than in traditional denotational semantics.

In Chapter 4, we will demonstrate how the formal concepts of monads and monad transformers fit nicely into the Gofer [Jones, 1991] type system. By implementing modular monadic semantics in Gofer, we obtain a modular interpreter.

In Chapter 5, we will apply monadic semantics to semantics-directed compilation. We will show how an effective and provably correct compilation scheme can be derived, taking advantage of the modularity and reasoning power of the monadic framework. We will put some of our ideas to test by building a retargeted Haskell compiler.

Throughout the thesis, we will use a common source language to address various issues in monadic semantics, modular interpreters, and compilation. The

source language is introduced in the next section.

1.5 The Source Language

The source language we consider in this thesis has a variety of features, including different function call mechanisms, imperative features, first-class continuation, tracing (for debugging), and nondeterminism.

$e ::= n$		$e_1 + e_2$	(arithmetic operations)				
		v		$\lambda v.e$	(variables and functions)		
		$(e_1 e_2)_n$	(call-by-name)				
		$(e_1 e_2)_v$	(call-by-value)				
		$(e_1 e_2)_l$	(lazy evaluation)				
		callcc	(first-class continuations)				
		$e_1 := e_2$		ref e		deref e	(imperative features)
		label @ e	(trace labels)				
		$\{e_0, e_1, \dots\}$	(nondeterminism)				

To simplify the presentation, we use one form of function abstraction that can be applied using any of the three function application mechanisms. We can observe the differences in various function call mechanisms with the help of trace messages. For example, evaluating:

$$((\lambda x.x + x)(l @ 1))_n$$

results in 2 after printing the trace message “ l ” twice, whereas

$$((\lambda x.x + x)(l @ 1))_v$$

prints “ l ” only once. Nondeterminism is captured by returning all possible results.

For example:

$$\{1, 3\} + \{2, 5\}$$

results in $\{3, 6, 5, 8\}$.

1.6 A Notation

For clarity, we adopt the following short-hand:

$$\begin{array}{l}
 E : \text{Term} \rightarrow M \text{ Value} \\
 E[e_1 + e_2] = \text{bind} (E[e_1]) \\
 \quad (\lambda i. \text{bind} (E[e_2]) \\
 \quad \quad (\lambda j. \text{return} (i + j)))
 \end{array}$$

↓

$$\begin{array}{l}
 E : \text{Term} \rightarrow M \text{ Value} \\
 E[e_1 + e_2] = \{ i \leftarrow E[e_1]; \\
 \quad \quad \quad j \leftarrow E[e_2]; \\
 \quad \quad \quad \text{return} (i + j) \}
 \end{array}$$

This notation is similar to the “do” syntax in Haskell [Peterson and Hammond, 1996], and is also somewhat similar to monad comprehensions [Wadler, 1990]. It is important to remember that, despite the imperative feel, the monadic semantics is still made up of lambdas and applications. We will use *bind* and its short-hand notation interchangeably, depending on whichever is more convenient.

Chapter 2

Modular Monadic Semantics

In this chapter, we present the modular monadic semantics of our source language. Compared with traditional denotational semantics, our approach captures individual programming language features using modular building blocks.

Figure 2.1 shows how our modular monadic semantics is organized. High-level features are defined based on a set of “kernel-level” operations. The expression $e_1 := e_2$, for example, is interpreted by the low-level primitive operation *update*.

While it is a well-known practice to base programming language semantics on a kernel language, the novelty of our approach lies in how the kernel-level primitive operations are organized. In our framework, depending on how much support the upper layers need, any set of primitive operations can be put together in a modular way using an abstraction mechanism called *monad transformers* [Moggi, 1990] [Liang *et al.*, 1995]. Monad transformers provide the power needed to represent the abstract notion of programming language features, but still allow us to access low-level semantic details. However, monad transformers are defined as higher-order functions and our monadic semantics is no more than a structured version of denotational semantics, so conventional reasoning methods (such as β conversion) apply.

The modular monadic semantics is composed of two parts:

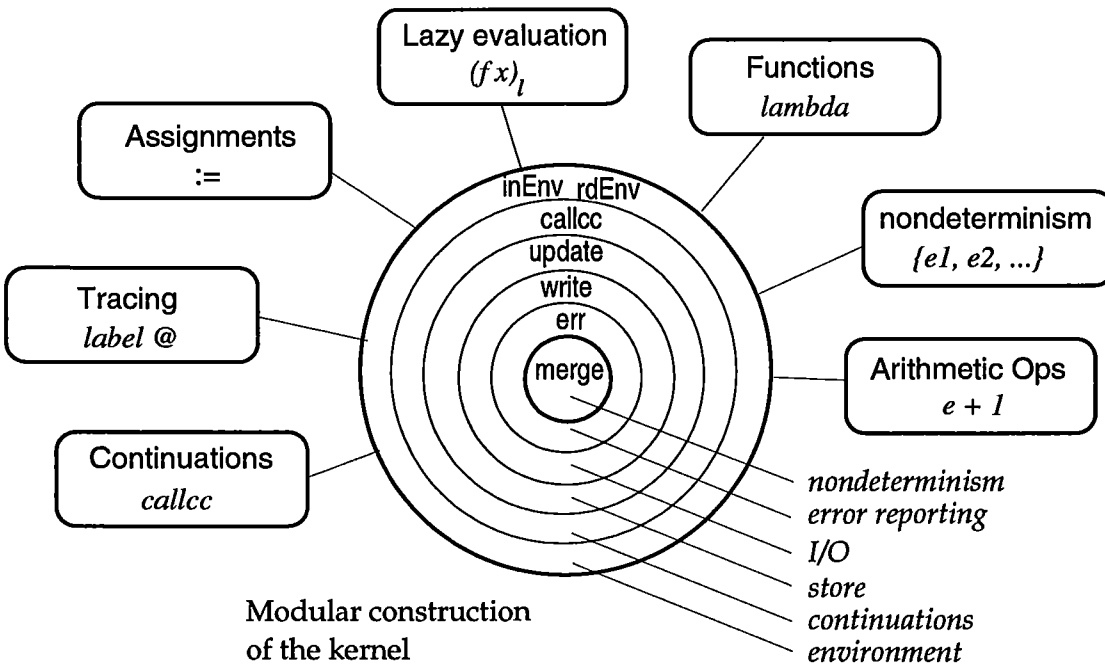


Figure 2.1: The organization of modular monadic semantics

Modular Semantic Building Blocks Semantic building blocks (represented by rectangular blocks in Figure 2.1) define the monadic semantics of individual source language features. Semantic building blocks are independent of each other, although they are based on a common set of kernel-level operations. Two building blocks may be supported by the same kernel-level operation. For example, both assignments and lazy evaluation may use the same store.

Monad Transformers Monad transformers define the kernel-level operations in a modular way. Multiple monad transformers can be composed to form the underlying monad used by all semantic building blocks. In Figure 2.1, monad transformers that support environment, continuations, store, etc. are used to construct the underlying monad.

Modular semantic building blocks and monad transformers are the topics of the following sections.

2.1 Modular Semantic Building Blocks

Each modular semantic building block defines the monadic semantics for a particular source language feature. Traditional denotational semantics maps, say, a term, an environment and a continuation to an answer. In contrast, monadic semantics maps terms to computations, where the details of the environment, store, etc. are *hidden*. Specifically, our semantic evaluation function E has type:

$$E : \text{Term} \rightarrow M \text{Value}$$

where M is a monad equipped with two basic operations:

$$\text{bind} : M a \rightarrow (a \rightarrow M b) \rightarrow M b$$

$$\text{return} : a \rightarrow M a$$

Value is the domain sum of basic values and functions; and $M \text{Value}$ represents computations that return *Value* as result. Functions map computations to computations:

$$\text{type Fun} = M \text{Value} \rightarrow M \text{Value}$$

$$\text{type Value} = \text{Int} + \text{Bool} + \text{Addr} + \text{Fun} + \dots$$

As will be seen, this generality allows us to model call-by-name, call-by-value and lazy evaluation with only one kind of lambda abstraction (but 3 kinds of function application) in the source language.

In this section, we will present the semantic building blocks needed for our source language. The monad operations *return* and *bind* are the basic operations used by every building block. In addition, each semantic building block depends on several other kernel-level operations that are specific to its purpose.

2.1.1 The Arithmetic Building Block

The semantics for arithmetic expressions is as follows:

$$\begin{aligned}
E[[n]] &= \text{return } (\mathbf{in}_{Int} n) \\
E[[e_1 + e_2]] &= \{ v_1 \leftarrow E[[e_1]]; \\
&\quad v_2 \leftarrow E[[e_2]]; \\
&\quad \text{if } (\mathbf{is}_{Int} v_1 \text{ and } \mathbf{is}_{Int} v_2) \text{ then} \\
&\quad\quad \text{return } (\mathbf{in}_{Int} (\mathbf{out}_{Int} v_1 + \mathbf{out}_{Int} v_2)) \\
&\quad \text{else} \\
&\quad\quad \text{err "type error" } \}
\end{aligned}$$

\mathbf{in}_{Int} is the injection function from Int to the $Value$ domain, whereas \mathbf{out}_{Int} is the projection function from the $Value$ domain to Int . The kernel-level function (to be defined later):

$$err : String \rightarrow M a$$

reports error conditions (which, in this case, are type errors). In other words, err is an operation supported by the underlying monad M . For clarity, from now on we will omit domain injection/projection and type checking.

$E[[n]]$ returns the number n (injected into the $Value$ domain) as the result of a trivial computation. To evaluate $e_1 + e_2$, we evaluate e_1 and e_2 in turn, and then sum the results.

2.1.2 The Function Building Block

In denotational semantics, functions are supported using an environment—a mapping from variable names to their denotation. We introduce an environment Env which maps variable names to computations,¹ and two kernel-level operations that retrieve the current environment and perform a computation in a given environment, respectively:

¹We do not need an environment that maps names to computations in order to support call-by-value. However, we need such an environment to support call-by-name and lazy evaluation. We will discuss this issue in more details in Section 2.1.8.

type *Env* = *Name* → *M Value*
rdEnv : *M Env*
inEnv : *Env* → *M Value* → *M Value*

The definition of *rdEnv* and *inEnv* will be given later. The semantics for variables, function abstraction, call-by-name and call-by-value are as follows:

$$\begin{aligned}
 E[[v]] &= \{\rho \leftarrow rdEnv; \rho [[v]]\} \\
 E[[\lambda v.e]] &= \{\rho \leftarrow rdEnv; return(\lambda c.inEnv \rho [c/[[v]]] E[[e]])\} \\
 E[[e_1 e_2]_n] &= \{f \leftarrow E[[e_1]]; \rho \leftarrow rdEnv; f(inEnv \rho E[[e_2]])\} \\
 E[[e_1 e_2]_v] &= \{f \leftarrow E[[e_1]]; v \leftarrow E[[e_2]]; f(return v)\}
 \end{aligned}$$

Because there is no risk of confusion, we drop the parentheses around $\rho[c/v]$ and $E[[e]]$ in the application of *inEnv*.

The difference between call-by-value and call-by-name is clear: the former reduces the argument before invoking the function,² whereas the latter packages the argument with the current environment to form a closure.

2.1.3 The References and Assignment Building Block

Imperative features can be supported using a store—a mapping from locations (of type *Loc*) to computations. Three functions allocate, read from and write to the memory cells in the store:

alloc : *M Loc*
read : *Loc* → *M Value*
write : (*Loc*, *M Value*) → *M ()*

The monadic semantics for references and assignment is as follows:

²To be precise, the call-by-value semantics is only preserved when the underlying monad enforces an evaluation order dependency. This is true of the continuation, state, and error monads. However, the identity and environment monads do not actually force the evaluation of c_1 before c_2 in $\{x \leftarrow c_1; c_2\}$.

$$\begin{aligned}
E[\text{ref } e] &= \{v \leftarrow E[e]; l \leftarrow \text{alloc}; \text{write } (l, \text{return } v)\} \\
E[\text{deref } e] &= \{l \leftarrow E[e]; \text{read } l\} \\
E[e_1 := e_2] &= \{l \leftarrow E[e_1]; v \leftarrow E[e_2]; \text{write } (l, \text{return } v)\}
\end{aligned}$$

To create a reference, we evaluate the expression, allocate a new memory cell, and store in the location of the memory cell a trivial computation that returns the value of the expression. The argument of `deref` evaluates to a location, at which the stored value can be read. To assign an expression to a location, we evaluate the expression, and update the location with a trivial computation that returns the value of the expression.

Note that we only store trivial computations. We could alternatively give the semantics for references and assignment using a store that maps locations to values, rather than locations to computations. The reason we store computations is to simplify the overall presentation, so that we will not need to introduce a separate kernel-level store operation for our next feature—lazy evaluation.

2.1.4 The Lazy Evaluation Building Block

Using the same store for references and assignments, we can implement lazy evaluation whose operational semantics implies *caching* of results.

$$\begin{aligned}
E[(e_1 e_2)_l] &= \{ f \leftarrow E[e_1]; \\
&\quad l \leftarrow \text{alloc}; \\
&\quad \rho \leftarrow \text{rdEnv}; \\
&\quad \text{let } \text{thunk} = \{ v \leftarrow \text{inEnv } \rho \ E[e_2]; \\
&\quad\quad\quad _ \leftarrow \text{write } (l, \text{return } v); \\
&\quad\quad\quad \text{return } v \} \\
&\quad \text{in } \{ _ \leftarrow \text{write } (l, \text{thunk}); \\
&\quad\quad f (\text{read } l) \} \}
\end{aligned}$$

Before entering the function, we allocate a memory cell and store a thunk (a computation that updates itself) in it. After the argument is first evaluated, the result is stored back to the memory cell, overwriting the thunk itself.

2.1.5 The Program Tracing Building Block

Given a kernel-level function:

$$\text{output} : \text{String} \rightarrow M ()$$

that prints out a string, we can support tracing. Labels attached to expressions cause a "trace record" to be invoked whenever that expression is evaluated:

$$\begin{aligned} E[l @ e] = \{ & _ \leftarrow \text{output} ("enter " ++ l); \\ & v \leftarrow E[e]; \\ & _ \leftarrow \text{output} ("leave " ++ l); \\ & \text{return } v \} \end{aligned}$$

Here we see that some of the features of monitoring semantics [Kishon *et al.*, 1991] are easily incorporated into our framework.

2.1.6 The Continuation Building Block

Continuation is a powerful mechanism for modeling control flow in denotational semantics. In particular, *callcc* (call with current continuation) is a useful language construct. Here is a simple example to show how *callcc* works:

$$\text{callcc } (\lambda k. (k \ 100)_v) \implies 100$$

When applied to a function, *callcc* captures the current continuation, and passes the continuation as the argument *k*. The continuation itself is captured as a function. When captured continuation is later applied to the value 100, the control flow is

transferred back to the point where the continuation was initially captured. The value (100) passed to the continuation is the result returned from *callcc*.

The power of *callcc* lies in that the captured continuation does not have to be invoked immediately. We may store the continuation into data structures, perform other computations, and then invoke the stored continuation to transfer the control back to where we issued *callcc*. For this reason, *callcc* can be used to model a wide variety of non-local control flow, including, for example, catch/throw, error handling, coroutines, and thread context switches. Scheme [Clinger and Rees, 1991] and SML [Milner *et al.*, 1990] incorporate *callcc* as a language feature.

As expected, the kernel-level operation *callcc* takes a function argument that in turn takes a continuation:

$$\text{callcc} : ((\text{Value} \rightarrow M \text{Value}) \rightarrow M \text{Value}) \rightarrow M \text{Value}$$

We define the semantics of source-level *callcc* as a function expecting another function as an argument, to which the current continuation will be passed:

$$E[\text{callcc}] = \text{return } (\lambda f. \{f' \leftarrow f; \text{callcc}(\lambda k. f'(\lambda a. \{x \leftarrow a; kx\}))\})$$

The result of $E[\text{callcc}]$ is a trivial computation that returns a function. The argument of the function, f , evaluates to the current continuation (f').

2.1.7 The Nondeterminism Building Block

Given a kernel-level function:

$$\text{merge} : \text{List } (M a) \rightarrow M a$$

that merges a list of computations into a single (nondeterministic) computation, nondeterminism semantics can be expressed as:

$$E[\{e_0, e_1, \dots\}] = \text{merge } [E[e_0], E[e_1], \dots]$$

$E[e_0]$, $E[e_1]$, etc. are a list of computations denoting the nondeterministic behavior.

2.1.8 Alternative Definitions of the Environment and Store

In the building blocks presented so far, we used one environment that mapped variable names to computations, and used one store that mapped locations to computations. As we have pointed out, this generality is not necessary for some of the building blocks. For example, the call-by-value semantics only needs an environment that maps variable names to values, whereas the reference and assignment semantics only needs a store for values.

Modular monadic semantics is flexible enough that we can easily introduce multiple environments and stores, so that each building block is supported by exactly the right set of operations. To specify call-by-value functions, for example, we can use an environment that maps variable names to values. If we later add call-by-name functions, we simply add a new environment that maps variable names to computations. Similarly for the reference and assignment building block, we can introduce a store that maps locations to values, separate from the requirements of lazy evaluation.

If we store variables in two separate environments, we will need to be able to distinguish, at the source language level, call-by-value functions from call-by-name functions. Thus instead of using one syntax for all three kinds of function abstractions (as in Section 1.5), we will need to have two separate syntactic constructs: one for call-by-value, the other for call-by-name and lazy evaluation. Variables will then be stored in either of the two environments, depending on what kind of function abstraction the variable is introduced in.

We will not present the details of designing a modular semantics with multiple environments and stores. Instead, we emphasize that the simplifications we made in previous sections to ease presentation do not fundamentally limit the modularity of our approach.

Feature	Function
Error reporting	$err : String \rightarrow M a$
Environment	$rdEnv : M Env$ $inEnv : Env \rightarrow M a \rightarrow M a$
Store	$alloc : M Loc$ $read : Loc \rightarrow M Value$ $write : (Loc, M Value) \rightarrow M ()$
Output	$output : String \rightarrow M ()$
Continuations	$callcc : ((a \rightarrow M b) \rightarrow M a) \rightarrow M a$
Nondeterminism	$merge : List (M a) \rightarrow M a$

Table 2.1: Monad operations used in the semantics

2.2 Monads With Operations

Semantic building blocks depend on other kernel-level operations in addition to *unit* and *bind*. From the last section, it is clear that the operations listed in Table 2.1 must be supported.

If we were writing the semantics in the traditional way, now would be the time to set up the domains and define the functions listed in the table. The major drawback of such a monolithic approach is that we have to take into account all other features when we define an operation for one specific feature. When we define *callcc*, for example, we have to decide how it interacts with the store and environment etc. And, if we later want to add more features, the semantic domains and all kernel-level functions may have to be redefined.

Monad transformers, on the other hand, allow us to capture individual language features. Furthermore, the concept of *lifting* allows us to account for the interactions between various features. Monad transformers and lifting are the topics of the next two sections.

To simplify the set of operations, we note that both the store and output (used by the program tracing building block) have to do with some notion of *state*. Thus

we could define *alloc*, *read*, *write*, and *output* in terms of the function:

$$\text{update} : (s \rightarrow s) \rightarrow M s$$

for some suitably chosen state type *s*. We can read the state by passing *update* the identity function, and update the state by passing it a state transformer. For example, we can model output by using *String* as the state type:

$$\begin{aligned} \text{output} & : \text{String} \rightarrow m () \\ \text{output msg} & = \{ _ \leftarrow \text{update} (\lambda \text{sofar}.\text{sofar} ++ \text{msg}); \\ & \quad \text{return} () \} \end{aligned}$$

The underscore (*_*) indicates that the return value of *update* is ignored.

2.3 Monad Transformers

To get an intuitive understanding of monad transformers, consider the merging of a state monad with an arbitrary monad, an example originally appeared in Moggi's note [Moggi, 1990]:

$$\text{type StateT } s \ m \ a = s \rightarrow m (s, a)$$

The type variable *m* represents a type constructor. We will later show that, if *m* is a monad, then so is *StateT s m*. Therefore *StateT s* is a monad transformer. For example, if we substitute the identity monad:

$$\text{type Id } a = a$$

for *m* in the above monad transformer, then we arrive at:

$$\begin{aligned} \text{StateT } s \ \text{Id } a & = s \rightarrow \text{Id} (s, a) \\ & = s \rightarrow (s, a) \end{aligned}$$

which is the standard state monad found, for example, in Wadler's work [Wadler, 1992].

We will formally define monad transformers in Section 3.1.2. For now we note that a monad transformer t has a number of capabilities:

First, it transforms any monad m to monad $t\ m$. Functions $return_{t\ m}$ and $bind_{t\ m}$ are naturally defined in terms of $return_m$ and $bind_m$.

Second, it can embed any computation in monad m as a computation in monad $t\ m$. Every monad transformer is equipped with a function:

$$lift_t : m\ a \rightarrow t\ m\ a$$

which maps any computation in monad m to a computation in monad $t\ m$.

Third, it adds operations (i.e., introduces new features) to a monad. The *StateT* monad transformer, for example, adds state s to the monad it is applied to, and the resulting monad accepts *update* as a legitimate operation.

Lastly, monad transformers compose easily. For example, applying both *StateT* s_1 and *StateT* s_2 to the identity monad, we get:

$$\begin{aligned} StateT\ s_1\ (StateT\ s_2\ Id)\ a &= s_1 \rightarrow (StateT\ s_2\ Id)\ (s_1, a) \\ &= s_1 \rightarrow s_2 \rightarrow (s_2, (s_1, a)), \end{aligned}$$

which is the expected type signature for transforming both states s_1 and s_2 . The observant reader will note, however, an immediate problem: in the resulting monad, which state does *update* act upon? In general, this is the problem of *lifting* monad operations through transformers, and will be addressed in the next section.

The remainder of this section introduces the monad transformers that cover all the features listed in Table 2.1. Some of these (*StateT*, *ContT*, and *ErrorT*) appear in an abstract form in Moggi's note [Moggi, 1990]. The *environment* monad is similar to the *state reader* by Wadler [Wadler, 1990]. The *state* and *environment* monad transformers are related to ideas found in Jones and Duponcheel's work [Jones, 1993] [Jones and Duponcheel, 1993].

We will attach subscripts to monadic operations to distinguish between the different monads they operate on. Some monad transformers use two additional functions: *map* and *join*. These functions, which can be used in any monad, are easily defined in terms of *return* and *bind*:

$$\begin{aligned} \text{map}_m & : (a \rightarrow b) \rightarrow m\ a \rightarrow m\ b \\ \text{map}_m\ f\ e & = \text{bind}_m\ e\ (\lambda a.\text{return}_m\ (f\ a)) \\ \text{join}_m\ z & : m\ (m\ a) \rightarrow m\ a \\ \text{join}_m\ z & = \text{bind}_m\ z\ (\lambda a.a) \end{aligned}$$

2.3.1 The State Monad Transformer

The state monad transformer introduces an updatable state into an existing monad. The resulting monad accepts an additional operation *update*, and is called a *state monad*.

Previously, we described the state monad transformer with a type definition:

$$\text{type StateT } s\ m\ a = s \rightarrow m\ (s, a)$$

To complete the definition, we must also provide the *return* and *bind* functions for *StateT* *s* *m*:

$$\begin{aligned} \text{return}_{\text{StateT } s\ m} & = \lambda s.\text{return}_m\ (s, x) \\ \text{bind}_{\text{StateT } s\ m}\ m\ k & = \lambda s_0.\text{bind}_m\ (m\ s_0)\ (\lambda (s_1, a).k\ a\ s_1) \end{aligned}$$

Given these definitions, if *return*_{*m*}, *bind*_{*m*}, and *m* form a monad, then so do *return*_{*StateT* *s* *m*}, *bind*_{*StateT* *s* *m*} and *StateT* *s* *m*. A more formal characterization of the relationships between *m* and *StateT* *s* *m* will be given in Chapter 3.

Next, we define the *lift* function, which simply performs the computation in the new context and preserves the state.

$$\begin{aligned} \text{lift}_{\text{StateT } s} & : m\ a \rightarrow \text{StateT } s\ m\ a \\ \text{lift}_{\text{StateT } s}\ c & = \lambda s.\{x \leftarrow c; \text{return}_m\ (s, x)\}_m \end{aligned}$$

Finally, a state monad must support the *update* operation, which transforms the state using the given f , and returns the old state:

$$\begin{aligned} \text{update}_{\text{StateT } s \ m} & : (s \rightarrow s) \rightarrow \text{StateT } s \ m \ s \\ \text{update}_{\text{StateT } s \ m} f & = \lambda s. \text{return}_m (f \ s, s) \end{aligned}$$

2.3.2 The Environment Monad Transformer

$\text{EnvT } r$ transforms any monad into an *environment monad* that supports inEnv and rdEnv . The definition of bind shows that two subsequent computation steps run under the same environment ρ (of type r). (Compare this with the state monad, where the second computation is run in the state returned by the first computation.)

$$\begin{aligned} \text{type EnvT } r \ m \ a & = r \rightarrow m \ a \\ \text{return}_{\text{EnvT } r \ m} a & = \lambda \rho. \text{return}_m a \\ \text{bind}_{\text{EnvT } r \ m} m \ k & = \lambda \rho. \text{bind}_m (m \ \rho) (\lambda a. k \ a \ \rho) \end{aligned}$$

The result of lifting a computation through the environment monad is a computation that ignores its environment.

$$\begin{aligned} \text{lift}_{\text{EnvT } r} & : m \ a \rightarrow \text{EnvT } r \ m \ a \\ \text{lift}_{\text{EnvT } r} c & = \lambda \rho. c \end{aligned}$$

inEnv ignores the environment carried inside the monad, and performs the computation in the given environment.

$$\begin{aligned} \text{inEnv}_{\text{EnvT } r \ m} & : r \rightarrow \text{EnvT } r \ m \ a \rightarrow \text{EnvT } r \ m \ a \\ \text{inEnv}_{\text{EnvT } r \ m} \rho \ m & = \lambda \rho'. m \ \rho \\ \text{rdEnv}_{\text{EnvT } r \ m} & : \text{EnvT } r \ m \ r \\ \text{rdEnv}_{\text{EnvT } r \ m} & = \lambda \rho. \text{return}_m \ \rho \end{aligned}$$

2.3.3 The Error Monad Transformer

Monad *Err* completes a series of computations if all succeed, or aborts as soon as an error occurs. The monad transformer *ErrT* transforms a monad into an *error monad* that supports *err* as a valid operation.

```

data Err a      = Ok a | Err String
type ErrT m a  = m (Err a)

returnErrT m a = returnm (Ok a)
bindErrT m m k = bindm m (λa.case a of
                                (Ok x)    → k x
                                (Err msg) → returnm (Err msg))

```

To lift a computation across *ErrT*, we tag the result with *Ok*:

```

liftErrT : m a → ErrT m a
liftErrT = mapm Ok

```

The semantic function *err* throws away any intermediate result, and returns the error value *Err*.

```

err : String → ErrT m a
err = returnm · Err

```

2.3.4 The Continuation Monad Transformer

We define the continuation monad transformer as:

```

type ContT c m a = (a → m c) → m c

returnContT c m x = λk.k x
bindContT c m m f = λk.m (λa.f a k)

```

$ContT$ introduces an additional continuation argument (of type $a \rightarrow m\ c$), where c is the answer type. By the above definitions of *return* and *bind*, all computations in monad $ContT\ c\ m$ are carried out in the continuation passing style.

Lift for $ContT\ c\ m$ turns out to be the same as $bind_m$. (Indeed they have the same type signature.)

$$lift_{ContT\ c} : m\ a \rightarrow ContT\ c\ m\ a$$

$$lift_{ContT\ c} = bind_m$$

$ContT$ transforms any monads to a *continuation monad*, which supports an additional operation *callcc*. $callcc\ f$ invokes the computation in f , passing it a continuation that, once applied, throws away the current continuation k' and invokes the captured continuation k .

$$callcc_{ContT\ c\ m} : ((a \rightarrow ContT\ c\ m\ b) \rightarrow ContT\ c\ m\ a) \rightarrow ContT\ c\ m\ a$$

$$callcc_{ContT\ c\ m}\ f = \lambda k.f\ (\lambda a.\lambda k'.k\ a)\ k$$

2.3.5 The List Monad

In denotational semantics, nondeterminism is usually captured by a list of all possible results. It is known that lists compose with a special kind of monads called *commutative monads* [Jones and Duponcheel, 1993]. It is not clear, however, if lists compose with arbitrary monads. Since many useful monads (e.g. state, error and continuation monads) are not commutative, we cannot define a list monad transformer—one which adds the operation *merge* to any monad.

Fortunately, every other monad transformer we have considered in this thesis properly transforms arbitrary monads. We thus can use lists as the *base* monad, to which other transformers can be applied. We recall the definition of the well-known list type and its monadic operations:

```
data List a = a : List a  -- Cons cell
           | []          -- Nil
```

```

returnList x = [x]
bindList m k = case m of
    []      → []
    (x : xs) → k x ++ (bindList xs k)

```

The *merge* function of the *List* monad is the well-known list concatenation operation:

```

mergeList : List (List a) → List a
mergeList []      = []
mergeList (x : xs) = x ++ mergeList xs

```

2.4 Liftings

We have introduced monad transformers that add useful operations to a given monad, but we have not addressed how these operations can be carried through other layers of monad transformers. This process is called the *lifting* of operations.

Lifting an operation f in monad m through a monad transformer t results in an operation whose type signature can be derived by substituting all occurrences of m in the type of f with $t\ m$. For example, lifting:

$$\text{inEnv} : r \rightarrow m\ a \rightarrow m\ a$$

through t results in an operation with type:

$$\text{inEnv} : r \rightarrow t\ m\ a \rightarrow t\ m\ a$$

Moggi [Moggi, 1990] studied the problem of lifting under a categorical context. The objective was to identify liftable operations from their type signatures. Unfortunately, many useful operations such as *merge*, *inEnv* and *callcc* failed to meet Moggi's criteria, and were left unsolved.

Instead, we consider how to lift these difficult cases individually. This allows us to make use of their definitions (rather than just their types), and to find ways to lift them through all of the monad transformers studied so far.

This is exactly where monad transformers provide us with an opportunity to study how various programming language features interact. The easy-to-lift cases correspond to features that are independent in nature, while the more involved cases require a deeper analysis of monad structures to clarify the semantics.

An unfortunate consequence of our approach is that, as we consider more monad transformers, the number of possible liftings grows quadratically. It seems, however, that there are not too many different kinds of monad transformers (although there may be many *instances* of the same monad transformer such as *StateT*). The monad transformers that we have introduced so far are able to model almost all commonly known features of sequential languages.³

Some operations are more difficult to lift than others. In particular, *inEnv* and *callcc* require special attention. We will first list the easy cases, followed by the rest. Although we present a number of liftings in this chapter, we will defer the formal explanation of why these listing are the desirable ones to Chapter 3.

2.4.1 The Easy Cases

RdEnv, *err* and *update* take a non-monadic type, and return a computation. They are handled by *lift*. For any monad transformer *t* applied to monad *m*, we have:

³An example of the features we cannot model is concurrent computation in multi-threaded programs. In addition, the state monad transformer is more general than what is needed to model output. The *output monad transformer* [Moggi, 1990] is also able to support the *output* operation:

```

type OutputT m a = m (String, a)
returnOutputT m x = returnm ("", x)
bindOutputT m m k = {(o1, a) ← m; (o2, b) ← k a; returnm (o1 ++ o2, b)}m
liftOutputT : m a → OutputT m a
liftOutputT c = {x ← c; returnm ("", x)}m
outputOutputT m : String → OutputT m
outputOutputT m s = returnm (s, ())

```

Investigating the properties of *OutputT* and its relationship with *StateT* is a topic for future research.

$$\begin{aligned}
rdEnv_{t\ m} &= lift_t\ rdEnv_m \\
err_{t\ m} &= lift_t \cdot err_m \\
update_{t\ m} &= lift_t \cdot update_m
\end{aligned}$$

Because *List* always is the base monad, we only have to consider cases when (possibly a sequence of) monad transformers are applied to *List*:

$$merge_{(t_1 \dots (t_n\ List) \dots)} = join_{(t_1 \dots (t_n\ List) \dots)} \cdot lift_{t_1} \cdot \dots \cdot lift_{t_n}$$

2.4.2 Lifting *Callcc*

The crucial issue in lifting *callcc* through a monad transformer, for example, *EnvT r*, is to specify how it interacts with the newly introduced environment *r*. The following lifting discards the *current* environment ρ' upon invoking the captured continuation *k*. The execution will continue in the environment ρ captured when *callcc* was first invoked. This is indeed how SML's *callcc* normally interacts with the environment.

$$\begin{aligned}
callcc_{EnvT\ r\ m} &: ((a \rightarrow r \rightarrow m\ b) \rightarrow r \rightarrow m\ a) \rightarrow r \rightarrow m\ a \\
callcc_{EnvT\ r\ m} &= \lambda\rho.callcc_m(\lambda k.f(\lambda a.\lambda\rho'.ka)\rho)
\end{aligned}$$

In lifting *callcc* through *StateT*, we have a choice of passing either the current state s_1 or the captured state s_0 . The former is the usual semantics for *callcc*, and the latter is useful in Tolmach and Appel's approach to debugging [Tolmach and Appel, 1990].

$$\begin{aligned}
callcc_{StateT\ s\ m} &: ((a \rightarrow s \rightarrow m(s, b)) \rightarrow s \rightarrow m(s, a)) \rightarrow s \rightarrow m(s, a) \\
callcc_{StateT\ s\ m\ f} &= \lambda s_0.callcc_m(\lambda k.f(\lambda a.\lambda s_1.k(s_1, a))\ s_0)
\end{aligned}$$

The above shows the usual *callcc* semantics, and can be changed to the "debugging" version by instead passing (s_0, a) to *k*:

$$callcc_{StateT\ s\ m\ f} = \lambda s_0.callcc_m(\lambda k.f(\lambda a.\lambda s_1.k(s_0, a))\ s_0)$$

Callcc can be lifted through *ErrT* as follows:

$$\begin{aligned} \text{callcc}_{\text{ErrT } m} & : ((a \rightarrow m(\text{Err } b)) \rightarrow m(\text{Err } a)) \rightarrow m(\text{Err } a) \\ \text{callcc}_{\text{ErrT } m} f & = \text{callcc}_m(\lambda k. f(\lambda a. k(\text{Ok } a))) \end{aligned}$$

2.4.3 Lifting *InEnv*

The liftings of *inEnv* through *EnvT* and *StateT* are similar:

$$\begin{aligned} \text{inEnv}_{\text{EnvT } r' m} & : r \rightarrow (r' \rightarrow m a) \rightarrow r' \rightarrow m a \\ \text{inEnv}_{\text{EnvT } r' m} \rho e & = \lambda \rho'. \text{inEnv}_m \rho (e \rho') \\ \text{inEnv}_{\text{StateT } s m} & : r \rightarrow (s \rightarrow m (s, a)) \rightarrow s \rightarrow m (s, a) \\ \text{inEnv}_{\text{StateT } s m} \rho e & = \lambda s. \text{inEnv}_m \rho (e s) \end{aligned}$$

A function of type:

$$m a \rightarrow m a$$

maps $m (\text{Err } a)$ to $m (\text{Err } a)$, thus *inEnv* stays the same after being lifted through *ErrT*.

We do not know of a desirable way to lift *inEnv* through *ContT*. This means that we always have to apply the continuation monad transformer before we apply environment monad transformers. In the following lifting, for example, the environment is not restored when *c* invokes *k*, and would thus reflect the history of dynamic execution.

$$\begin{aligned} \text{inEnv}_{\text{ContT } c m} \rho c & = \lambda k. \text{inEnv}_m \rho (c k) \\ \text{rdEnv}_{\text{ContT } c m} & = \text{lift rdEnv}_m \end{aligned}$$

2.5 Summary

Monad transformers and lifting are summarized in Figures 2.2 and 2.3. The most problematic case is the continuation monad transformer *ContT*. Not only are

<p>State:</p> <p>type $StateT\ s\ m\ a = s \rightarrow m\ (s, a)$</p> <p>$return_{StateT\ s\ m\ a} = \lambda s. return_m(s, a)$</p> <p>$bind_{StateT\ s\ m}\ e\ k =$ $\lambda s. \{(s', a) \leftarrow es; kas'\}_m$</p> <p>$update\ f = \lambda s. return_m(fs, s)$</p> <p>$lift_{StateT\ s}\ e = \lambda s. \{a \leftarrow e; return_m(s, a)\}_m$</p>	<p>Environment:</p> <p>type $EnvT\ r\ m\ a = r \rightarrow m\ a$</p> <p>$return_{EnvT\ r\ m\ a} = \lambda \rho. return_m\ a$</p> <p>$bind_{EnvT\ r\ m}\ e\ k =$ $\lambda \rho. \{a \leftarrow e\rho; ka\rho\}_m$</p> <p>$rdEnv = \lambda \rho. return_m\ \rho$</p> <p>$inEnv\ \rho\ c = \lambda \rho'. c\ \rho$</p> <p>$lift_{EnvT\ r}\ e = \lambda \rho. e$</p>
<p>Errors:</p> <p>type $Err\ a = Ok\ a \mid Err\ String$</p> <p>type $ErrT\ m\ a = m\ (Err\ a)$</p> <p>$return_{ErrT\ m} = return_m \cdot Ok$</p> <p>$bind_{ErrT\ m}\ e\ k =$ $\{a \leftarrow e;$ case a of $Ok\ x \rightarrow kx$ $Err\ s \rightarrow return_m(Err\ s)\}$</p> <p>$err = return_m \cdot Err$</p> <p>$lift_{ErrT} = map_m\ Ok$</p>	<p>Continuation:</p> <p>type $ContT\ c\ m\ a = (a \rightarrow m\ c) \rightarrow m\ c$</p> <p>$return_{ContT\ c\ m\ a} = \lambda k. ka$</p> <p>$bind_{ContT\ c\ m}\ e\ f = \lambda k. e(\lambda a. fak)$</p> <p>$callcc\ f = \lambda k. f(\lambda a. \lambda k'. ka)k$</p> <p>$lift_{ContT\ c} = bind_m$</p>

Figure 2.2: Monad transformers

operations relatively hard to lift though $ContT$, the $callcc$ operation also requires more work to lift through other monad transformers.

Equipped with the monad transformers, we can construct the underlying monad M to support all of the semantic building blocks in Section 2.1:

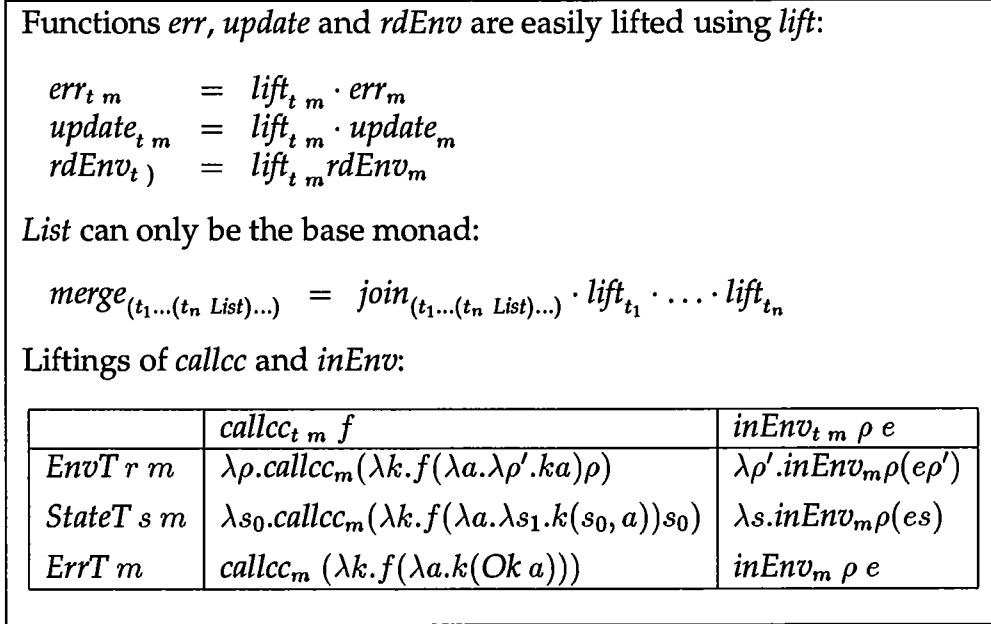


Figure 2.3: Liftings

```

type M a = EnvT Env           (environment)
           (ContT Answer      (continuation)
            (StateT Store     (store)
             (StateT IO       (input/output)
              (ErrT           (error reporting)
               List)))) a    (nondeterminism)

```

Env , $Answer$, $Store$, and IO are the types of environment, answer, store, and I/O channels, respectively. The order of some monad transformers can be changed. However, because of the limitations in lifting $inEnv$ through $ContT$, we cannot exchange the order of $EnvT$ and $ContT$.

By using a series of abstractions, modular monadic semantics turns the monolithic structure of traditional denotational semantics into reusable components. The modularity is manifested at two levels, high-level monadic building blocks and low-level monad transformers.

We have, however, only achieved part of our goal. Without a theory of monads and monad transformers, we would have to unfold the definitions of all kernel-level monadic operations (such as *bind* and *inEnv*) to reason about semantic building blocks and the source language. In the next chapter, we will present a theory that enables us to perform equational reasoning at a higher-level with a set of laws and axioms.

Chapter 3

A Theory of Monads and Monad Transformers

The purpose of developing a theory for monads and monad transformers is to reason about the monadic semantics without having to unfold the definitions of kernel-level monadic operations (e.g., *bind*, *inEnv*). Unfolding the monadic operations would defeat the purpose of the modular abstraction mechanism. Instead, we will make it possible to perform equational reasoning at a high level by providing a set of properties directly associated with various monadic operations. An example in Chapter 5 will further demonstrate that reasoning in the monadic framework offers modular proofs and more general results. In this chapter, we concentrate on the fundamental properties of monads and monad transformers.

This chapter begins with the formal definition of monads and monad transformers, based on Moggi's and Walder's earlier work. The main topics of this chapter are:

- how monadic axioms capture the properties of individual programming language features, and
- how natural liftings preserve existing features and capture the interactions

between the newly added feature and existing features.

This chapter ends with a discussion of the order of composing monad transformers.

3.1 Monad and Monad Transformers

In this section we give a formal definition of monads and monad transformers.

3.1.1 Monads

Definition 3.1.1 A monad M is a triple consisting of a type constructor and two functions: $(m, \text{return}_m, \text{bind}_m)$. Monads must satisfy the following laws [Moggi, 1990]:

$$\begin{aligned} \{b \leftarrow \text{return } a; k\ b\} &= k\ a && \text{(left unit)} \\ \{a \leftarrow e; \text{return } a\} &= e && \text{(right unit)} \\ \{v_1 \leftarrow e_1; \{v_2 \leftarrow e_2; e_3\}\} &= \{v_2 \leftarrow \{v_1 \leftarrow e_1; e_2\}; e_3\} && \text{(associativity)} \end{aligned}$$

Intuitively, the (left and right) unit laws say that trivial computations can be skipped in certain contexts; and the associativity law captures the very basic property of sequencing, one that we usually take for granted in imperative programming languages.

Note that in the associativity law, e_1 is in the scope of v_2 on the right hand side but not so on the left hand side. In applying this law, we must make sure that there is no unwanted name capture.

The type constructors *Id* and *List* introduced in Chapter 2 are well-known monads (presented in, for example, [Wadler, 1990]):

Proposition 3.1.1 *Id* and *List* are monads.

3.1.2 Monad Transformers

To capture monad transformers formally, we first introduce *monad morphisms* [Moggi, 1990]:

Definition 3.1.2 A **monad morphism** f between monads m and m' is a function of type:

$$f : m\ a \rightarrow m'\ a$$

satisfying:

$$\begin{aligned} f.\text{return}_m &= \text{return}_{m'} \\ f(\text{bind}_m\ m\ k) &= \text{bind}_{m'}(f\ m)(f \cdot k) \end{aligned}$$

Note that f is polymorphic in a . We can now define monad transformers as follows:

Definition 3.1.3 A **monad transformer** consists of a type constructor t and an associated function lift_t , where t maps any given monad $(m, \text{return}_m, \text{bind}_m)$ to a new monad $(t\ m, \text{return}_{t\ m}, \text{bind}_{t\ m})$. Furthermore, lift_t is a monad morphism between m and $t\ m$:

$$\text{lift}_t : m\ a \rightarrow t\ m\ a$$

Therefore lifting a trivial computation results in a trivial computation; lifting a sequence of computations is equivalent to first lifting them individually, and then combining them in the lifted monad.

The type constructors listed in Figure 2.2 satisfy the above definition.

Proposition 3.1.2 $\text{EnvT}\ r$, $\text{StateT}\ s$, ErrT , and $\text{ContT}\ c$ are monad transformers.

It is well known that these type constructors transform monads to monads. “ $\text{EnvT}\ r$ ” is the composable reader monad presented in [Jones and Duponcheel, 1993]. The remaining three were discovered by Moggi [Moggi, 1990]. Appendix A contains detailed proofs that the corresponding *lift* functions are indeed monad morphisms.

Monad transformers compose with each other (a property that follows immediately from the definition of monad morphisms):

Proposition 3.1.3 Given monad transformers t_1 and t_2 , $t_1 \cdot t_2$ is a monad transformer with:

$$\begin{aligned} \text{type } (t_1 \cdot t_2) m a &= t_1 (t_2 m) a \\ \text{lift}_{(t_1 \cdot t_2)} &= \text{lift}_{t_1} \cdot \text{lift}_{t_2} \end{aligned}$$

3.2 Environment Axioms

Environments have a profound impact on programming language semantics and compilation. For example, lexically scoped languages fit well into the environment model. The monadic framework provides us a way to capture the essential properties of environments as follows:

Proposition 3.2.1 The environment operations, $rdEnv$ and $inEnv$ satisfy the following axioms:

$$\begin{aligned} (inEnv \rho) \cdot \text{return} &= \text{return} && \text{(unit)} \\ inEnv \rho \{v \leftarrow e_1; e_2\} &= \{v \leftarrow inEnv \rho e_1; inEnv \rho e_2\} && \text{(distribution)} \\ inEnv \rho rdEnv &= \text{return } \rho && \text{(cancellation)} \\ inEnv \rho' (inEnv \rho e) &= inEnv \rho e && \text{(overriding)} \end{aligned}$$

Intuitively, a trivial computation cannot depend on the environment (the unit law); the environment stays the same across a sequence of computations (the distribution law); the environment does not change between a set and a read if there are no intervening computations (the cancellation law); and an inner environment supersedes an outer one (the overriding law). The distribution law, for example, is what distinguishes the environment from a store. A store does not distribute across a sequence of computations. It is updated as the computation progresses.

We can prove the environment axioms by first verifying that they hold after the environment monad transformer is applied, and then by making sure that they are preserved through the liftings of $rdEnv$ and $inEnv$. A detailed proof of these results is included in Appendix A.

In Chapter 5, we will present an example that uses the environment axioms to prove a property about compiling the source language.

The environment axioms provide an answer to the question: “what constitutes an environment?” We expect that useful monadic axioms can be derived for other features, following the earlier efforts on state [Hudak and Bloss, 1985] [Peyton Jones and Wadler, 1993] [Hudak, 1992], continuations [Felleisen *et al.*, 1986] [Felleisen and Hieb, 1992] and exceptions [Spivey, 1990].

3.3 Natural Liftings

In this section, we investigate what conditions a desirable lifting must satisfy. First we will formalize how types are transformed in the lifting process. We will then introduce the *natural lifting condition* and verify that the liftings we constructed in Section 2.4 are indeed natural.

3.3.1 Lifting Types

How does its type change when an operation is lifted? The set of operations we consider has the following types in monad m :

$\tau ::=$	A	(type constants)
	a	(type variables)
	$\tau \rightarrow \tau$	(functions)
	(τ, τ)	(products)
	$List \tau$	(lists)
	$m \tau$	(computations)

When an operation is lifted through the monad transformer t , its new type can be derived by substituting all occurrences of m in the type with $t m$. Formally, $[\cdot]_t$ is

the mapping of types across the monad transformer t :

$$\begin{aligned}
 [A]_t &= A \\
 [a]_t &= a \\
 [\tau_1 \rightarrow \tau_2]_t &= [\tau_1]_t \rightarrow [\tau_2]_t \\
 [(\tau_1, \tau_2)]_t &= ([\tau_1]_t, [\tau_2]_t) \\
 [List \tau]_t &= List [\tau]_t \\
 [m \tau]_t &= t m [\tau]_t
 \end{aligned}$$

3.3.2 Natural Lifting Condition

What properties should a particular lifting satisfy? Recall that in Section 2.4.3, we noted that the following was not a desirable lifting of $inEnv$ through $ContT$:

$$inEnv_{ContT \ c \ m \ r \ c} = \lambda k. inEnv_m \ r \ (c \ k)$$

The problem is that the environment is not restored when c invokes k , which is equivalent to, for example, not popping off the arguments after a function returns. This lifting is not desirable because the new feature (continuation) has disrupted the existing feature (environment).

Intuitively, any programs not using the added feature should behave in the same way after a monad transformer is applied. The monad morphism property of *lift* ensures that single computations are properly lifted. But some operations, such as *callcc*, have more complex types—they take computations as arguments. We extend Moggi’s original definition and define **natural liftings** as a family of relations \mathcal{L}_τ , indexed by type τ :

Definition 3.3.1 \mathcal{L}_τ is a **natural lifting** of operations of type τ along the monad transformer t if it satisfies:

$$\mathcal{L}_\tau \quad : \quad \tau \rightarrow [\tau]_t$$

$$\mathcal{L}_A \quad = \quad id \quad (1)$$

$$\mathcal{L}_a \quad = \quad id \quad (2)$$

$$\mathcal{L}_{\tau_1 \rightarrow \tau_2} \quad = \quad \lambda f. f' \text{ satisfying:}$$

$$\forall \mathcal{L}_{\tau_1}, \exists \mathcal{L}_{\tau_2}, \text{ such that: } f' \cdot \mathcal{L}_{\tau_1} = \mathcal{L}_{\tau_2} \cdot f \quad (3)$$

$$\mathcal{L}_{(\tau_1, \tau_2)} \quad = \quad \lambda(a, b). (\mathcal{L}_{\tau_1} a, \mathcal{L}_{\tau_2} b) \quad (4)$$

$$\mathcal{L}_{List \ \tau} \quad = \quad map_{List} \ \mathcal{L}_\tau \quad (5)$$

$$\mathcal{L}_m \ \tau \quad = \quad lift_t \cdot (map_m \ \mathcal{L}_\tau) \quad (6)$$

Despite the similarity between cases 5 and 6, case 5 is in fact more similar to case 4. Both cases 4 and 5 map τ across the some basic data type. In case 6, m is the monad on which the monad transformer t is applied.

Constant types (such as integer) and polymorphic types do not depend on any particular monad. (See cases 1 and 2.) On the other hand, we expect a lifted function, when applied to a value lifted from the domain of the original function, to return a lifting of the result of applying the original function to the unlifted value. This relationship is precisely captured by equation 3, which corresponds to the following commuting diagram:

$$\begin{array}{ccc}
 [\tau_1]_t & \xrightarrow{f'} & [\tau_2]_t \\
 \mathcal{L}_{\tau_1} \uparrow & & \uparrow \mathcal{L}_{\tau_2} \\
 \tau_1 & \xrightarrow{f} & \tau_2
 \end{array}$$

The liftings of tuples and lists are straightforward. Finally, the *lift* operator that comes with the monad transformer m lifts computations in m . Note that \mathcal{L}_τ is mapped to the result of the computation, which may involve other computations.

The above does *not* provide a constructive definition for a type-parametric lifting function \mathcal{L} . The “satisfying” clause in the third equation specifies a constraint, rather than a definition of f' . That is why we define \mathcal{L} as a relation rather than a function. In practice, we first find out by hand how to lift an operation through particular monad transformers, and then use the above equations to verify that such a lifting is indeed natural.

3.3.3 Examples

We now verify the natural lifting condition for the liftings in Section 2.4. The easy cases (*update*, *err* and *rdEnv*) are covered by the following theorem by Moggi [Moggi, 1990]:

Proposition 3.3.1 If function f 's domain does not involve any monadic type, then:

$$\text{lift}_t \cdot f$$

is a natural lifting of f through any monad transformer t .

Since the domain type (call it τ) does not involve the monad, the lifting of τ is τ itself.

The above theorem follows from the commutativity of the following diagram:

$$\begin{array}{ccc}
 \tau & \xrightarrow{\text{lift}_t \cdot f} & [\tau']_t \\
 \uparrow \text{id} & & \uparrow \text{lift}_t \\
 \tau & \xrightarrow{f} & \tau'
 \end{array}$$

We address the remaining cases (*merge*, *inEnv* and *callcc*) separately.

Proposition 3.3.2

$$\text{merge}_{(t_1 \dots (t_n \text{ List}) \dots)} = \text{join}_{(t_1 \dots (t_n \text{ List}) \dots)} \cdot \text{lift}_{t_1} \cdots \text{lift}_{t_n}$$

is a natural lifting of $\text{merge}_{\text{List}}$.

To prove that the lifting for *merge* is natural, we need the following property of *map* and *join*:

Lemma 3.3.1 If t is a monad transformer, m a monad, then:

$$\text{lift}_t \cdot \text{join}_m = \text{join}_{t \cdot m} \cdot \text{lift}_t \cdot (\text{map}_m \text{lift}_t)$$

Proof:

$$\begin{aligned} \text{lift}_t (\text{join}_m e) &= \text{lift}_t \{a \leftarrow e; a\}_m && (\text{join}) \\ &= \{a \leftarrow \text{lift}_t e; \text{lift}_t a\}_{t \cdot m} && (\text{monad morphism}) \\ &= \{a \leftarrow \text{lift}_t e; b \leftarrow \text{return}_{t \cdot m}(\text{lift}_t a); b\}_{t \cdot m} && (\text{left unit}) \\ &= \text{join}_{t \cdot m} \{a \leftarrow \text{lift}_t e; \text{return}_{t \cdot m}(\text{lift}_t a)\}_{t \cdot m} && (\text{join}) \\ &= \text{join}_{t \cdot m} \{a \leftarrow \text{lift}_t e; \text{lift}_t(\text{return}_m(\text{lift}_t a))\}_{t \cdot m} && (\text{monad morphism}) \\ &= \text{join}_{t \cdot m} (\text{lift}_t \{a \leftarrow e; \text{return}_m(\text{lift}_t a)\}_m) && (\text{monad morphism}) \\ &= \text{join}_{t \cdot m} (\text{lift}_t (\text{map}_m \text{lift}_t e)) && (\text{map}) \end{aligned}$$

We can now prove Proposition 3.3.2 by verifying that the following diagram commutes:

$$\begin{array}{ccc} \text{List}((t_1 \dots (t_n \text{List}) \dots) a) & \xrightarrow{\text{merge}_{(t_1 \dots (t_n \text{List}) \dots)}} & (t_1 \dots (t_n \text{List}) \dots) a \\ \uparrow \text{map}_{\text{List}}(\text{lift}_{t_1} \dots \text{lift}_{t_n}) & & \uparrow \text{lift}_{t_1} \dots \text{lift}_{t_n} \\ \text{List}(\text{List } a) & \xrightarrow{\text{merge}_{\text{List}}} & \text{List } a \end{array}$$

Indeed we have:

$$\begin{aligned} \text{merge}_{(t_1 \dots (t_n \text{List}) \dots)} \cdot \text{map}_{\text{List}}(\text{lift}_{t_1} \dots \text{lift}_{t_n}) & \\ &= \text{merge}_{(t_1 \dots t_n) \text{List}} \cdot \text{map}_{\text{List}} \text{lift}_{t_1 \dots t_n} \quad (3.1.3) \end{aligned}$$

$$= \text{lift}_{t_1 \dots t_n} \cdot \text{merge}_{\text{List}} \quad (3.3.1)$$

$$= \text{lift}_{t_1} \dots \text{lift}_{t_n} \cdot \text{merge}_{\text{List}} \quad (3.1.3)$$

Proposition 3.3.3

$$\begin{aligned}
inEnv_{EnvT\ r'\ m}\rho\ e &= \lambda\rho'.inEnv_m\ \rho\ (e\ \rho') \\
inEnv_{StateT\ s\ m}\rho\ e &= \lambda s.inEnv_m\ \rho\ (e\ s) \\
inEnv_{ErrT\ m}\rho\ e &= inEnv_m\ \rho\ e
\end{aligned}$$

are natural liftings of $inEnv_m$.

For $inEnv_t\ m$ to be a natural lifting, we need to prove that:

$$inEnv_t\ m\rho \cdot lift_t = lift_t \cdot inEnv_m\rho$$

Indeed we have:

$$\begin{aligned}
inEnv_{EnvT\ r'\ m}\rho\ (lift_{EnvT\ r'\ m}\ e) &= \lambda\rho'.inEnv_m\rho(lift_{EnvT\ r'\ m}\ e\rho') \quad (inEnv_{EnvT\ r'\ m}) \\
&= \lambda\rho'.inEnv_m\rho e \quad (lift_{EnvT\ r'\ m}) \\
&= lift_{EnvT\ r'\ m}(inEnv_m\rho e) \quad (lift_{EnvT\ r'\ m})
\end{aligned}$$

$$\begin{aligned}
inEnv_{StateT\ s\ m}\rho\ (lift_{StateT\ s\ m}\ e) &= \lambda s.inEnv_m\rho(lift_{StateT\ s\ m}\ es) \quad (inEnv_{StateT\ s\ m}) \\
&= \lambda s.inEnv_m\rho\{a \leftarrow e; return_m(s, a)\}_m \quad (lift_{StateT\ s\ m}) \\
&= \lambda s.\{a \leftarrow inEnv_m\rho e; return_m(s, a)\}_m \quad (\text{Prop. 3.2.1}) \\
&= lift_{StateT\ s\ m}(inEnv_m\rho e) \quad (lift_{StateT\ s\ m})
\end{aligned}$$

$$\begin{aligned}
inEnv_{ErrT\ m}\rho\ (lift_{ErrT\ m}\ e) &= inEnv_m\rho(lift_{ErrT\ m}\ e) \quad (inEnv_{ErrT\ m}) \\
&= inEnv_m\rho\{a \leftarrow e; return_m(Ok\ a)\}_m \quad (lift_{ErrT\ m},\ map_m) \\
&= \{a \leftarrow inEnv_m\rho e; return_m(Ok\ a)\}_m \quad (\text{Prop. 3.2.1}) \\
&= lift_{ErrT\ m}(inEnv_m\rho e) \quad (lift_{ErrT\ m})
\end{aligned}$$

Proposition 3.3.4

$$\begin{aligned}
callcc_{EnvT\ r\ m} &= \lambda\rho.callcc_m(\lambda k.f(\lambda a.\lambda\rho'.ka)\rho) \\
callcc_{StateT\ s\ m}\ f &= \lambda s_0.callcc_m(\lambda k.f(\lambda a.\lambda s_1.k(s_0, a))\ s_0) \\
callcc_{ErrT\ m}\ f &= callcc_m(\lambda k.f(\lambda a.k(Ok\ a)))
\end{aligned}$$

are natural liftings of $callcc_m$.

To prove Proposition 3.3.4, we apply Definition 3.3.1 to the type of $callcc$, and arrive at the following lemma:

Lemma 3.3.2

$callcc_{t\ m}$ is a natural lifting of $callcc_m$

iff:

$$\forall f, f'. (\forall k. f'(lift_t \cdot k) = lift_t(fk)) \Rightarrow callcc_{t\ m} f' = lift_t(callcc_m f)$$

Using Lemma 3.3.2, it is easy to show that $callcc_{EnvT\ r\ m}$ is a natural lifting of $callcc_m$:

$$\begin{aligned} callcc_{EnvT\ r\ m} f' &= \lambda\rho. callcc_m(\lambda k. f'(\lambda a. \lambda\rho'. ka)\rho) && (callcc_{EnvT\ r\ m}) \\ &= \lambda\rho. callcc_m(\lambda k. f'(\lambda a. lift_{EnvT\ r}(ka))\rho) && (lift_{EnvT\ r}) \\ &= \lambda\rho. callcc_m(\lambda k. lift_{EnvT\ r}(fk)\rho) && (\text{prerequisite of 3.3.2}) \\ &= \lambda\rho. callcc_m(\lambda k. fk) && (lift_{EnvT\ r}) \\ &= lift_{EnvT\ r}(callcc_m f) && (lift_{EnvT\ r}) \end{aligned}$$

Paterson [Paterson, 1995] showed a simple proof for the naturalness of $callcc_{ErrT\ m}$ using the free theorem [Wadler, 1989] for $callcc$:

$$\begin{aligned} &\forall g, h, f, f'. \\ &(\forall k, k'. k' \cdot g = map\ h \cdot k \Rightarrow f'k' = map\ g(fk)) \Rightarrow \\ &callcc\ f' = map\ g(callcc\ f) \end{aligned}$$

By specializing f' to $\lambda k. f''(k \cdot g)$, we can transform the free theorem to:

Lemma 3.3.3

$$\begin{aligned} &\forall g, h, f, f''. \\ &(\forall k, f''(map\ h \cdot k) = map\ g(fk)) \Rightarrow \\ &callcc(\lambda k. f''(k \cdot g)) = map\ g(callcc\ f) \end{aligned}$$

We use Lemma 3.3.3 to prove $callcc_{ErrT\ m}$ is a natural lifting. Let:

$$g = h = Ok$$

we have:

$$\begin{aligned}
callcc_{ErrT\ m} f'' &= callcc_m(\lambda k. f(k \cdot Ok)) && (callcc_{ErrT\ m}) \\
&= map_m Ok (callcc_m f) && (\text{free theorem and prerequisite in 3.3.2}) \\
&= lift_{ErrT} (callcc_m f) && (lift_{ErrT})
\end{aligned}$$

Thus $callcc_{ErrT\ m}$ is a natural lifting, following Lemma 3.3.2.

The free theorem, however, is not powerful enough to prove the naturalness of $callcc_{StateT\ s\ m}$. Instead, we introduce the following lemma, which is a slight variation of the free theorem.

Lemma 3.3.4

$$\begin{aligned}
&\forall g, h, f, f', s_0. \\
&(\forall k, f'(\lambda x. \lambda s. map(\lambda x. h(s, x))(kx))s_0 = map\ g\ (fk) \Rightarrow \\
&callcc(\lambda k. f'(\lambda x. \lambda s. k(gx))s_0) = map\ g\ (callcc\ f)
\end{aligned}$$

The proof of the lemma is in Appendix A. We will apply Lemma 3.3.4 with the following specialized definitions to prove $callcc_{StateT\ s\ m}$ is a natural lifting:

$$\begin{aligned}
g &= \lambda x. (s_0, x) \\
h &= \lambda x. x
\end{aligned}$$

The proof is carried out in two steps. First, we verify the prerequisite of Lemma 3.3.4, using the prerequisite of Lemma 3.3.2.

$$\begin{aligned}
&f'(\lambda x. \lambda s. map_m(\lambda x. h(s, x))(kx))s_0 \\
&= f'(\lambda x. \lambda s. map_m(\lambda x. (s, x))(kx))s_0 && (h) \\
&= f'(\lambda x. lift_{StateT\ s}(kx))s_0 && (lift_{StateT\ s}) \\
&= lift_{StateT\ s}(fk)s_0 && (\text{prerequisite of 3.3.2}) \\
&= map_m\ g\ (fk) && (lift_{StateT\ s})
\end{aligned}$$

Second, we use the result of Lemma 3.3.4 to establish the the sufficient and necessary condition in Lemma 3.3.2:

$$\begin{aligned}
callcc_{StateT\ s\ m}\ f'\ s_0 &= callcc_m(\lambda k. f'(\lambda a. \lambda s_1. k(s_0, a)))s_0 && (callcc_{StateT\ s\ m}) \\
&= map_m(\lambda a. (s_0, a)) (callcc_m\ f) && (3.3.4) \\
&= lift_{StateT\ s} (callcc_m\ f) && (lift_{StateT\ s})
\end{aligned}$$

Apply the above to Lemma 3.3.2, we have proved that $callcc_{StateT\ s\ m}$ is a natural lifting.

So far we have established that *all* the liftings in Figure 2.3 are natural. Note that the following lifting of $callcc_{StateT\ s\ m}$:

$$callcc_{StateT\ s\ m}\ f = \lambda s_0. callcc_m (\lambda k. f (\lambda a. \lambda s_1. k (s_1, a))) s_0$$

which passes the current state to the continuation, is *not* natural. Here is a counterexample discovered by Paterson [Paterson, 1995]. Let:

$$f'k = lift_{StateT\ s}(f(\lambda x. bind(kx s_1)(\lambda (s', x). return\ x)))$$

For any state s_1 , f' and f meet the condition:

$$\forall k. f'(lift_{StateT\ s} \cdot k) = lift_{StateT\ s}(fk)$$

However:

$$\begin{aligned}
callcc_{StateT\ s}\ (ContT\ c\ Id)f'\ s_0\ k &= f(\lambda x. \lambda k'. k(s_1, x))(\lambda x. k(s_0, x)) \\
lift_{StateT\ s}(callcc_{ContT\ c\ Id}f)\ s_0\ k &= f(\lambda x. \lambda k'. k(s_0, x))(\lambda x. k(s_0, x))
\end{aligned}$$

are different.

3.4 Ordering of Monad Transformers

The ordering of monad transformers has an impact on the resulting semantics. For example, we have seen that lifting $callcc$ through $StateT$ results in a “debugging” semantics. On the other hand, if we apply $ContT$ to a state monad, then we get the usual semantics for $callcc$. To demonstrate, we construct two monads:

type $M1\ a = ContT\ c\ (StateT\ Int\ Id)\ a$

type $M2\ a = StateT\ Int\ (ContT\ c\ Id)\ a$

The program segment:

$$callcc(\lambda k.\{- \leftarrow update(\lambda x.x + 1); k0\})$$

expands to:

$$\lambda k.\lambda s_0.k\ 0\ (s_0 + 1)$$

in $M1$, but to:

$$\lambda s_0.\lambda k.k\ (s_0, 0)$$

in $M2$.

The key difference is that one combination captures the state in the continuation, whereas the other combination does not.

In general we can swap the ordering of some monad transformers (such as between $StateT$ and $EnvT$), but doing so to others (such as $ContT$) may effect semantics. This is consistent with earlier experience in combining monads [King and Wadler, 1993], and, in practice, provides us with an opportunity to fine tune the resulting semantics.

Chapter 4

Modular Monadic Interpreters

We can transform a denotational semantics description into an executable interpreter by translating the mathematical notations into corresponding programming constructs. Modern functional languages such as Haskell [Hudak *et al.*, 1992] or SML [Milner *et al.*, 1990] are particularly suitable because these languages offer features such as algebraic data types and higher-order functions that match well with the mathematical notations used in denotational semantics.

While the static type system in Haskell or SML is capable of implementing traditional denotational semantics, implementing monadic modular semantics in a strongly typed language has proved to be a challenge. For example, Steele [Steele Jr., 1994] reported numerous difficulties when he built a modular monadic interpreter in Haskell. Although the Haskell type system can implement individual monads and monad transformers as type constructors, modular monadic semantics requires the type system to capture relationships among different monads and monad transformers.

We have successfully implemented a modular monadic interpreter in *Gofer* [Jones, 1991], whose *constructor classes* and *multi-parameter type classes* provide just the added power over Haskell's type classes¹ to allow precise and convenient

¹The newly defined Haskell 1.3 [Peterson and Hammond, 1996] supports constructor classes

```

type Term = OR TermA      -- arithmetic
             ( OR TermF   -- functions
             ( OR TermR   -- assignment
             ( OR TermL   -- lazy evaluation
             ( OR TermT   -- tracing
             ( OR TermC   -- callcc
               TermN     -- nondeterminism
             ))))

type M = EnvT Env        -- environment
        ( ContT Answer  -- continuations
        ( StateT Store  -- memory cells
        ( StateT String -- trace output
        ( ErrT          -- error reporting
          List          -- multiple results
        ))))

type Value = OR Int      -- integers
             ( OR Loc    -- memory locations
             ( OR Fun    -- functions
             ( )))

```

Figure 4.1: Gofer specification of a modular interpreter

expression of the typing relationships. Figure 4.1 gives the high-level definition of the interpreter for our source language. The rest of the chapter will explain how the type declarations expand into a full interpreter. For now just note that `OR` is equivalent to the domain sum operator, and that `Term`, `Value` and `M` denote the abstract syntax, runtime values, and the interpreter monad, respectively.

(but not multi-parameter type classes).

4.1 Extensible Union Types

We begin with a discussion of a key idea in our implementation: how values and terms may be expressed as *extensible union types*. This facility has nothing to do with monads.

The disjoint union of two types is implemented by the data type OR:

```
data OR a b = L a
           | R b
```

where L and R are used to perform the conventional injection of a summand type into the union; conventional pattern-matching is used for projection. However, such injections and projections only work if we know the exact structure of the union. When building modular interpreters, an extensible union may be arbitrarily nested or extended. We would like a *single* pair of injection and projection functions to work on all such constructions.

To achieve this, we define a multi-parameter type class to implement the summand/union type relationship, which we refer to as a “subtype” relationship:

```
class SubType sub sup where
  inj :: sub -> sup           -- injection
  prj :: sup -> Maybe sub    -- projection

data Maybe a = Just a
             | Nothing
```

The Maybe data type is used because the projection function may fail. We can now express the relationships between the summand and union types:

```
instance SubType a (OR a b) where
  inj      = L
  prj (L x) = Just x
  prj _    = Nothing

instance SubType a b => SubType a (OR c b) where
```

```
inj      = R . inj
prj (R a) = prj a
prj _    = Nothing
```

It would appear that we could have a more symmetric instance declaration in place of the second declaration above:

```
instance SubType a (OR b a) where
  inj      = R
  prj (R x) = Just x
  prj _    = Nothing
```

With this declaration, however, the Gofer type system complains that `(OR a a)` is an overlapping instance. The type system cannot determine which of the two injection/projection pairs are applicable if the programmer supplies, for example, `(OR Int Int)` as the union type.

Now we can see how the Value domain used in Figure 4.1, for example, is actually constructed:

```
type Value = OR Int (OR Loc (OR Fun ()))
type Fun   = M Value -> M Value
```

With these definitions the Gofer type system will infer that `Int`, `Loc`, and `Fun` are all “subtypes” of `Value`, and the coercion functions `inj` and `prj` will be generated automatically.²

4.2 Interpreter Building Blocks

As seen in Figure 4.1, the `Term` type is also constructed as an extensible union (of subterm types). We define additionally a class `InterpC` to characterize the term types that we wish to interpret:

²Most of the typing problems Steele [Steele Jr., 1994] encountered disappear with the use of our extensible union types; in particular, there is no need for Steele’s “towers” of data types.

```
class InterpC t where
  interp :: t -> M Value
```

The behavior of the evaluation function `interp` on unions of terms is given in the obvious way:

```
instance (InterpC t1, InterpC t2) =>
  InterpC (OR t1 t2) where
  interp (L t) = interp t
  interp (R t) = interp t
```

The interpreter is just the method associated with the top-level type `Term`:

```
interp :: Term -> M Value
```

The interpreter building blocks are straightforward translations of the semantic building blocks in Section 2.1 into instance declarations. For example, the arithmetic building block can be implemented as follows:

```
data TermA = Num Int
           | Add Term Term

instance InterpC TermA where
  interp (Num x)    = returnInj x
  interp (Add x y) = interp x `bindPrj` \i ->
                        interp y `bindPrj` \j ->
                        returnInj ((i + j) :: Int)

returnInj      = return . inj
m `bindPrj` k = m `bind` \a ->
  case (prj a) of
    Just x  -> k x
    Nothing -> error "type error"
```

Note the simple use of `inj` and `prj` to inject/project the integer result into/out of the `Value` domain, regardless of how `Value` is eventually defined (`returnInj` and `bindPrj` make this a tad easier). The error function is Gofer's runtime error function which is denotationally equivalent to \perp .

We omit Gofer implementation of other interpreter building blocks. They can be similarly translated from the corresponding monadic semantics.

Before discussing how to implement the monad transformers needed to construct the interpreter monad M , we introduce Gofer's constructor classes through a motivating example.

4.3 Constructor Classes

Constructor classes [Jones, 1993] support abstraction of common features among type constructors. Haskell, for example, provides the standard `map` function to apply a function to each element of a given list:

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

Meanwhile, we can define similar functions for a wide range of other data types. For example:

```
data Tree a = Leaf a
            | Node (Tree a) (Tree a)
```

```
mapTree :: (a -> b) -> Tree a -> Tree b
```

```
mapTree f (Leaf x)    = Leaf (f x)
mapTree f (Node l r) = Node (mapTree f l) (mapTree f r)
```

The `mapTree` function has similar type and functionality to those of `map`. With this in mind, it seems a shame that we have to use different names for each of these variants. Indeed, Gofer allows type variables to stand for *type constructors*, on which the Haskell type class system has been extended to support overloading. To solve the problem with `map`, we can introduce a new constructor class `Functor` (in a categorical sense):

```
class Functor f where
  map :: (a -> b) -> f a -> f b
```

Now the standard list (`List`) and the user-defined type constructor `Tree` can be defined as instances of `Functor`:

```
instance Functor List where
  map f [] = []
  map f (x:xs) = f x : map f xs

instance Functor Tree where
  map f (Leaf x) = Leaf (f x)
  map f (Node l r) = Node (map f l) (map f r)
```

Constructor classes are extremely useful for dealing with multiple instances of monads and monad transformers (which are all type constructors).

4.4 Monads

We follow a well known approach [Jones, 1993] to define monads using a constructor class:

```
class Monad m where
  return :: a -> m a
  bind :: m a -> (a -> m b) -> m b

  map :: (a -> b) -> (m a -> m b)
  join :: m (m a) -> m a

  map f m = m `bind` \a -> return (f a)
  join m = m `bind` id
```

`Map` and `join` are conveniently defined as default methods in terms of `bind` and `return`. A specific monad, such as `List`, is an instance of the `Monad` class:

```
instance Monad List where
  return x = [x]

  [] `bind` k = []
  (x:xs) `bind` k = k x ++ (xs `bind` k)
```


The interesting properties of a monad are the additional operations it supports. We can further define subclasses of `Monad`, each containing an additional set of operations. For example, `ListMonad` has one extra operation `merge`:

```
class Monad m => ListMonad m where
  merge :: List (m a) -> m a
```

The standard list monad `List` implements `merge` as follows:

```
instance ListMonad List where
  --merge :: List (List a) -> List a
  merge []      = []
  merge (x:xs) = x ++ (merge xs)
```

Other classes of monads, such as `StateMonad`, `EnvMonad`, `ContMonad` and `ErrMonad`, can be similarly defined. (See Appendix B for details.)

4.5 Monad Transformers

We implement monad transformers in the following constructor class definition:

```
class MonadT t where
  lift :: (Monad m, Monad (t m)) => m a -> t m a
```

To illustrate how individual instances are defined, we use the state monad transformer (`StateT`) as an example. The Gofer implementation of `EnvT`, `ContT`, and `ErrT` can be found in Appendix B.

From Section 2.3 we know that applying monad transformer `StateT s` to monad `m` results in a monad `StateT s m`. Because Gofer only allows us to partially apply a data type, not a type synonym, we introduce a dummy data constructor and define `StateT` as an algebraic data type³:

³Haskell 1.3[Peterson and Hammond, 1996] introduces a newtype construct that can be used to avoid the run-time penalty of dummy data constructors such as `StateM`.

```

data StateT s m a = StateM (s -> m (s,a))
unStateM (StateM x) = x

instance Monad m => Monad (StateT s m) where
  return x           = StateM (\s -> return (s,x))
  (StateM m) 'bind' k =
    StateM (\s0 -> m s0 'bind' \(s1, a) ->
            unStateM (k a) s1)

```

The definition follows exactly from Figure 2.2, except for dealing with the `StateM` data constructor. Note that `bind` and `return` are not recursive functions; the constructor class system automatically infers that the functions appearing on the right are for monad `m`.

Next, we define `StateT s` as a monad transformer:

```

instance MonadT (StateT s) where
  -- lift  :: m a -> StateT s m a
  lift m = StateM (\s -> m 'bind' \x -> return (s,x))

```

We introduce `StateMonad` as a subclass of `Monad` with an additional operation `update`:

```

class Monad m => StateMonad s m where
  update :: (s -> s) -> m s

```

Monad transformer `StateT s` adds the `update` function on `s` to any monad `m`:

```

instance Monad m => StateMonad s (StateT s m) where
  update f = StateM (\s -> return (f s, s))

```

Finally, we can lift `update` through any monad transformer by composing it with `lift` (see Proposition 3.3.1):

```

instance (StateMonad s m, MonadT t) =>
  StateMonad s (t m) where
  update = lift . update

```

As another example of lifting, we can apply any monad transformer to `List` and obtain a `ListMonad` (see Proposition 3.3.2):

```

instance (MonadT t, Monad m) => ListMonad (t m) where
  merge = join . lift

```

4.6 Summary

We have shown that modular interpreter building blocks and monad transformers can be implemented using two key features in Gofer type system: *constructor classes* and *multi-parameter type classes*. Our approach offers several benefits. First, it allows us to experiment with and debug our ideas. Second, the overloading mechanism greatly facilitates representing multiple instances of monads and monad transformers, eliminating the need for subscripts. Third, type checking guarantees that we have enough features in the underlying monad to support the set of building blocks needed for our source language. For example, if we had instead constructed the monad `M` in figure 4.1 without the `StateT String` monad transformer:

```
type M = EnvT Env          -- environment
      ( ContT Answer      -- continuations
      ( StateT Store      -- memory cells
      -- missing state component for IO
      ( ErrT              -- error reporting
      List                -- multiple results
      )))
```

then the Gofer type system would complain that `StateMonad String M` cannot be inferred from the definition of `M`.

Chapter 5

Compilation

In this chapter we investigate how to compile the source language from its monadic semantics specification. The target language we consider is fairly high-level, providing support for closures, tagged data structures, basic control-flow (such as conditionals) and garbage collection. How to implement a back-end that efficiently supports such target languages has been investigated by a number of compiler research efforts (e.g., the techniques developed for T [Kranz *et al.*, 1986], SML/NJ [Appel, 1992], and Haskell [Peyton Jones, 1992]).

Even though we do not tackle the problem of building compiler back-ends, our work provides insights into how we may build a *common back-end* capable of supporting a variety of source languages. Writing separate back-ends for different source languages leads to duplication of efforts. On the other hand, a common back-end has the following benefits:

- It simplifies the task of constructing compilers.
- It allows multiple source languages to interoperate by freely exchanging compatible runtime data.

Modular monadic semantics fits well with a common back-end, because it is suitable for specifying multiple source languages, and, as will be seen, it leads to an

efficient and provably correct compilation scheme. This is achieved in several steps. First, we require that our semantics be *compositional*: the arguments in recursive calls to E are substructures of the argument received on the left-hand side. From a theoretical point of view, it makes inductive proofs on programs possible. In practice, this guarantees that, given any abstract syntax tree, we can recursively unfold all calls to the interpreter, effectively removing runtime dispatch on the abstract syntax tree.

Our second step is to simplify the resulting monadic style code composed out of various monadic operations (such as *bind* and *inEnv*). As will be seen in Section 5.1, monad laws are useful in simplifying code; and environment axioms can be used to eliminate the costly interpretive overhead of environment lookups. In Section 5.2, we formally prove that all environment lookups can be removed.

The final step (Section 5.3) is to map monadic-style intermediate code to the target language. The main focus is on how to utilize the built-in target language features.

We will present a compilation method for our source language (defined in Section 1.5) following the above steps. To demonstrate that our techniques are applicable to realistic languages, in Section 5.4 we will use monadic semantics to target the STG Language (the intermediate language of the Glasgow Haskell Compiler) to the SML/NJ back-end.

5.1 Using Monad Laws to Transform Programs

Following the monadic semantics presented in Chapter 2, by unfolding all calls to the semantic function E , we can transform source-level programs into monadic-style code. For example, “ $((\lambda x.x + 1) 2)_v$ ” is transformed to:

$$\begin{aligned}
E\llbracket((\lambda v.v + 1) 2)_v\rrbracket = & \\
\{ f \leftarrow \{ \rho \leftarrow rdEnv; & \\
\quad return (\lambda x.inEnv \rho[x/\llbracket v \rrbracket]) \{ i \leftarrow \{ \rho \leftarrow rdEnv; & \\
\quad \quad \quad \rho\llbracket v \rrbracket \}; & \\
\quad \quad \quad j \leftarrow return 1; & \\
\quad \quad \quad return (i + j) \} \} \}; & \\
v \leftarrow return 2; & \\
f(return v) \} &
\end{aligned}$$

Even without any further simplifications, the above code is clear enough to describe the computation. By applying monad laws we can simplify it to:

$$\begin{aligned}
\{ \rho \leftarrow rdEnv; & \\
(\lambda x.inEnv \rho[x/\llbracket v \rrbracket]) \{ \rho \leftarrow rdEnv; & \\
\quad i \leftarrow \rho\llbracket v \rrbracket; & \\
\quad return (i + 1) \} \} (return 2) \} &
\end{aligned}$$

By applying the distribution, unit and cancellation environment axioms, followed by the unit monad law, we can further transform the example code to:

$$\begin{aligned}
\{ \rho \leftarrow rdEnv; & \\
(\lambda x.\{ i \leftarrow inEnv \rho[x/\llbracket v \rrbracket] x; return (i + 1) \}) (return 2) \} &
\end{aligned}$$

Note that explicit environment accesses have disappeared. Instead, the meta-language environment is directly used to support function calls. This is exactly what good partial evaluators achieve when they transform interpreters to compilers.

Note that the true computation in the original expression “ $((\lambda x.x + 1) 2)_v$ ” is left unreduced. With traditional denotational semantics, it is harder to distinguish the redexes introduced by the compilation process from computations in the source program. In the above example, we could safely further reduce the intermediate code:

$$\begin{aligned}
 & (\lambda x. \{i \leftarrow x; \text{return } (i + 1)\}) (\text{return } 2) \\
 \Rightarrow & \{i \leftarrow \text{return } 2; \text{return } (v + 1)\} && (\beta) \\
 \Rightarrow & \text{return } 3 && (\text{left unit})
 \end{aligned}$$

However, in general, unrestricted reductions for arbitrary source programs could result in unwanted compile-time exceptions, such as in “ $((\lambda x. 10/x) 0)_v$.”

5.2 A Natural Semantics

We successfully transformed away the explicit environment in the above example, but can we do the same for arbitrary source programs? If that is possible, we will have an effective compilation scheme that uses the target language environment for the source language, without any interpretive overhead.

It turns out that we can indeed prove such a general result by using monad laws and environment axioms. Following Wand [Wand, 1990], we define a “natural semantics” that translates source language variables to lexical variables in the meta-language, and we prove that it is equivalent to the standard semantics.

5.2.1 Definition of a Natural Semantics

We adopt Wand’s definition of a *natural semantics* (which differs from Kahn’s notion [Clément *et al.*, 1986]) to our functional sub-language. For any source language variable name v , we assume there is a corresponding variable name \bar{v} in the meta-language, and $\bar{\rho}$ is an environment that maps variable name v to \bar{v} .

Definition 5.2.1 The *natural semantics* for the source language is defined as follows:

$$\begin{aligned}
N[[v]] &= \bar{v} \\
N[[\lambda v.e]] &= \text{return}(\lambda \bar{v}.\text{inEnv } \bar{\rho} N[[e]]) \\
N[[e_1 e_2]_n] &= \{f \leftarrow N[[e_1]]; f(\text{inEnv } \bar{\rho} N[[e_2]])\} \\
N[[e_1 e_2]_v] &= \{f \leftarrow N[[e_1]]; v \leftarrow N[[e_2]]; f(\text{return } v)\} \\
N[[e_1 e_2]_l] &= \{ f \leftarrow N[[e_1]; \\
&\quad l \leftarrow \text{alloc}; \\
&\quad \text{let thunk} = \{ v \leftarrow \text{inEnv } \bar{\rho} N[[e_2]; \\
&\quad \quad _ \leftarrow \text{write}(l, \text{return } v); \\
&\quad \quad \text{return } v \} \\
&\quad \text{in } \{ _ \leftarrow \text{write}(l, \text{thunk}); \\
&\quad \quad f(\text{read } l) \} \}
\end{aligned}$$

Other source-level constructs, such as $+$, $:=$, and calcc , do not explicitly deal with the environment, and have the same natural semantics as the standard semantics.

The natural semantics uses the environment of the meta-language for variables in the source language.

5.2.2 Correspondence between Natural and Standard Semantics

The next theorem, a variation of Wand's [Wand, 1990], states that the standard semantics and natural semantics are equivalent, and thus guarantees that it is safe to implement function calls in the source language using the meta-language environment.

Theorem 5.2.1 For any source language program e , we have:

$$\text{inEnv } \bar{\rho} E[[e]] = \text{inEnv } \bar{\rho} N[[e]]$$

The detailed proof is in Appendix A. The basic technique is equational reasoning based on the rules of lambda calculus (e.g., β reduction), monad laws, and

environment axioms. We establish the theorem for each semantic building block, independent of:

- the existence of other building blocks, and
- the organization of the underlying monad.

Therefore the result holds for each building block as long as the underlying monad provides the necessary kernel-level support so that the monad laws and environment axioms hold. The proof can be reused, even after other features are added into the source language.

The proof is possible because both the source language and meta language are lexically scoped. If the source language supported dynamically scoped functions:

$$E[\lambda v.e] = \text{return}(\lambda c.\{\rho \leftarrow rdEnv; inEnv \rho[c/\llbracket v \rrbracket] E[e]\}),$$

where the caller-site environment is used within the function body, then the theorem would fail to hold.

5.2.3 Benefits of Reasoning in Monadic Style

In denotational semantics, adding a feature may change the structure of the entire semantics, forcing us to redo the induction for every case of abstract syntax. For example, Wand [Wand, 1990] pointed out that he could change to a continuation-based semantics, and prove the theorem, but only by modifying the proofs accordingly.

Modular monadic semantics, on the other hand, offers highly modularized proofs and more general results. This is particularly applicable to real programming languages, which usually carry a large set of features and undergo evolving designs.

5.3 Targeting Monadic Code

In general, it is more efficient to use target language built-in features instead of monadic combinators defined as higher-order functions. We have seen how the explicit environment can be “absorbed” into the meta-language. This section addresses the question of whether we can do the same for other features, such as stores and continuations.

5.3.1 The Target Language Monad

We can view a target language as having a built-in monad supporting a set of monadic operations. For example, the following table lists the correspondence between certain monadic operations and ML constructs:

Monadic operations	ML constructs
<i>return</i> x	x
$\{x \leftarrow c_1; c_2\}$	let val $x = c_1$ in c_2 end
<i>update</i>	ref , ! , := , print
<i>calcc</i>	calcc
<i>err</i>	raise Err

Note that the imperative features in ML (e.g., **:=** and **print**) supports a single-threaded store, whereas the monadic *update* operation more generally supports recoverable store. It is easy to verify that the monad laws are satisfied in the above context. For example, the ML **let** construct is associative (assuming no unwanted name capturings occur):

$$\boxed{\begin{array}{l} \mathbf{let\ val\ } v_2 = \mathbf{let\ val\ } v_1 = c_1 \\ \qquad \qquad \mathbf{in\ } c_2 \mathbf{\ end} \\ \mathbf{in\ } c_3 \mathbf{\ end} \end{array}} = \boxed{\begin{array}{l} \mathbf{let\ val\ } v_1 = c_1 \\ \mathbf{in\ let\ val\ } v_2 = c_2 \\ \qquad \mathbf{in\ } c_3 \mathbf{\ end\ end} \end{array}}$$

5.3.2 Utilizing Target Language Features

We now investigate how to utilize the features directly supported by the target language monad. Because of a technical limitation related to nondeterminism, we tentatively drop it from our source language. (We will discuss the support for nondeterminism later.) The underlying monad M becomes:

```
type  $M\ a = EnvT\ Env\ (ContT\ Answer\ (StateT\ Store\ (StateT\ IO\ (ErrT\ Id))))\ a$ 
```

Now we substitute the base monad Id with the built-in ML monad (call it M_{ML}):

```
type  $M_1\ a = EnvT\ Env\ (ContT\ Answer\ (StateT\ Store\ (StateT\ IO\ (ErrT\ M_{ML}))))\ a$ 
```

Note that M_1 supports two sets of kernel-level operations for continuation, store, I/O, and error reporting. The monadic code can choose to use the ML built-in operations instead of those implemented as higher-order functions. In addition, if we have used the natural semantics to transform away all environment accesses, then the $EnvT$ monad transformer is no longer useful. Because the natural lifting condition guarantees that adding or deleting an unused monad transformer does not effect the result of the computation, it suffices to run the target program on M_2 :

```
type  $M_2\ a = M_{ML}\ a$ 
```

which directly utilizes the more efficient ML built-in features.

Therefore, by using a monad with a set of primitive monadic combinators, we can expose the features embedded in the target language. It then becomes clear what is directly supported in the target language, and what needs to be compiled explicitly.

The above process would have been impossible had we been working with traditional denotational semantics. Various features clutter up and make it hard to determine whether it is safe to remove certain interpretation overhead, and how to achieve that.

We do not need to transform away all monad transformers. For example, the following monad is also capable of supporting the source language:

```
type  $M_3$   $a = \text{ContT Answer } M_{ML} a$ 
```

Because M_3 supports two *callcc* operations, the monadic code can either use the ML built-in *callcc* function, or use the *callcc* supported by the continuation monad transformer.

5.3.3 Limitations of This Approach

It is important to recognize the limitations of the transformation process:

1. Unlike other features, nondeterminism must be directly supported by the target language, since the nondeterminism monad (*List*) must be the base monad. This is why we put aside nondeterminism in the preceding discussion.
2. We have shown that the ordering of monad transformers (in particular, the cases involving *ContT*) has an impact on the resulting semantics. In practice, we need to make sure when we use one monad transformer instead of another, that the resulting change of ordering does not have unwanted effects on semantics. For example, if we had left one of the state monad transformers unreduced:

```
type  $M_4$   $a = \text{StateT Store } M_{ML} a$ 
```

we have effectively swapped the order of *StateT* and *ContT*. (The latter is now supported in M_{ML} .)

5.3.4 Implications for a Common Back-end

To overcome the above limitations, a common back-end must support a rich set of features needed by a wide range of source languages, thereby guaranteeing that we can always transform away the monad transformers.

The ordering of monad transformers only effects the semantics of *callcc*. To deal with situations where the order of monad transformers matters, the back-end can provide multiple variations of the a monadic operation, with each version implementing a variation of the semantics. For example, a back-end can support a special version of *callcc* that captures the current state for the purpose of debugging.

5.4 An Experiment: Retargeting a Haskell Compiler

In this section, we put some of our ideas to the test by describing how we have used them to retarget the Glasgow Haskell Compiler (GHC). Although monadic semantics is equally capable of specifying various *static* semantics in the compilation process—for example, the type checker and strictness analyzer in GHC are both written in the monadic style—our focus has been on the dynamic semantics of the source language. Therefore we begin with the STG language, the output of the GHC front-end. The front-end has already transformed away Haskell’s syntactic sugar, carried out type checking, and performed various program analysis (such as strictness analysis) and front-end optimization (such as inlining).

The goal of this experiment is to connect the STG language with the SML/NJ back-end [Appel, 1992]. We chose the SML/NJ back-end because it provided efficient support for closures, tagged data structures (for implementing algebraic data types), basic control flow (if and case constructs), and garbage collection. This is exactly the set of run-time support needed to implement monadic semantics.

5.4.1 The STG Language

The STG language [Peyton Jones, 1992], whose grammar is shown in Figure 5.1¹, is a small purely-functional language with an operational reading as well as the usual

¹We dropped the update flag and free variable list from the lambda form because they are not used in our framework.

Expression	$e ::=$ <ul style="list-style-type: none"> $\text{let } b \text{ in } e$ (local definition) $\text{ letrec } \{b_1; \dots; b_n\} \text{ in } e$ (local recursion) $\text{ case } e \text{ of}$ (case on data type) <ul style="list-style-type: none"> $c_1 \{v_{11}, \dots\} \rightarrow e_1;$ \dots $c_n \{v_{n1}, \dots\} \rightarrow e_n;$ $v \rightarrow e$ (default) $\text{ case } e \text{ of}$ (case on literals) <ul style="list-style-type: none"> $l_1 \rightarrow e_1;$ \dots $l_n \rightarrow e_n;$ $v \rightarrow e$ (default) $\ v \{a_1, \dots, a_n\}$ (function application) $\ c \{a_1, \dots, a_n\}$ (constructor application) $\ p \{a_1, \dots, a_n\}$ (primitive application) $\ l$ (literals) 	
Binding	$b ::= v = \lambda \{v_1, \dots, v_n\} . e$ (closure)	
Atom	$a ::= v$ (variables) $ \ l$ (literals)	
Variable	v	
Primitive	p	
Constructor	c	
Literal	l	

Figure 5.1: The STG language

denotational semantics. STG code has the following distinctive characteristics:

- Arguments to functions and constructors are either variables or constants (together called *atoms*).
- All constructors and primitives are saturated. Functions, however, are by default *curried*.
- Case expressions cause evaluation to happen and perform one-level pattern matching.
- Let and `letrec` bindings create closures, which could either be functions or delayed computations.

For example, given the Haskell program:

```
let f x y = x + y
in f 5
```

the corresponding STG code² is:

```
let f = λ{x,y}. + {x,y}
in f{5}
```

Because the STG language supports partial applications, we pass one argument (5) to f even though f takes two arguments.

5.4.2 Compiling the STG Language

The STG language is designed to run on the GHC back-end—the STG machine. The major differences between the STG machine and the SML/NJ back-end are:

²To ease presentation, we use a slightly simplified version of STG code than the output of the GHC front-end. In particular, we omit the primitive integer constructors (I#).

- The STG machine represents all boxed data types as self-updating closures. Self-updating closures work similarly as the lazy evaluation semantics in Section 2.1.4. The first evaluation of a closure causes the closure to update itself to a trivial computation that simply returns the cached result. The SML/NJ back-end, on the other hand, supports normal closures as well as separate operations to create and update reference cells.
- The STG machine directly supports curried function application. An efficient built-in mechanism checks if functions have received enough arguments, and if not, creates a closure that waits for more arguments. In comparison, the SML/NJ back-end creates a closure for each intermediate argument, even if the arguments are supplied at once. For example, suppose that g takes 3 arguments. We first supply g with 2 arguments, and supply the third argument later in the computation. The STG machine only creates one intermediate closure, whereas the SML/NJ back-end creates two closures, one for each of the first two arguments.

Because the SML/NJ back-end does not support updatable closures, we implement lazy evaluation using the *cell mode* [Bloss *et al.*, 1988]. Figure 5.2 gives an example of how an expression (factorial of 3) is delayed and cached in the cell mode. To create a delayed computation, we allocate a two-element cell. The first element contains a flag bit signaling whether the second element contains a delayed computation or a cached value. The lazy evaluation semantics is guaranteed because the first evaluation sets the flag bit and overwrites the delayed computation with the result.

Because a delayed computation incurs both space and time overhead, a well-known technique in implementing lazy functional languages is to use the *strictness* information so that we can directly pass values, instead of cells when it is safe to do so. The GHC front-end includes a *strictness analyzer*. For example, given the example program in Section 5.4.1:

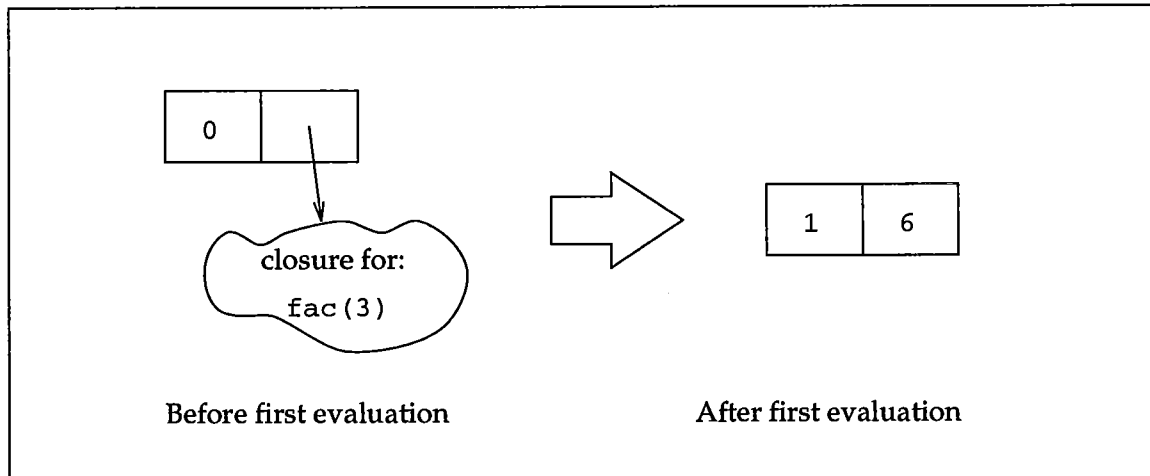


Figure 5.2: How an expression (`fac(3)`) is evaluated in cell mode

```
let f x y = x + y
in f 5
```

The GHC front-end is able to determine that `f` must evaluate both of its arguments, i.e., `f` is *strict* in `x` and `y`.

5.4.3 Monadic Semantics of the STG Language

We will present a portion of monadic semantics for the STG Language in this section. The complete semantics is listed in Appendix C.

As mentioned before, we can pass strict arguments directly without creating a delayed computation (i.e., updatable cells). The *Mode* type signals whether the data is passed as a value or a cell:

```
data Mode = C    -- (updatable cells)
          | V    -- (values)
```

Information about what mode each variable is in comes from the GHC strictness analyzer, which is part of the GHC front-end. We will use the *strictnessOf* function to obtain the mode of a given variable name. At run-time, a pair of functions convert between these two modes:

- *delay* e creates a cell containing e as the delayed computation.
- *force* e generates a demand on cell e and returns the value of the delayed computation.

We introduce a strictness environment to keep track of what mode each variables are in. The strictness environment maps variable names to strictness information:

```
data StrictnessInfo = VAR Mode           -- (variables)
                    | FUN [Mode1, ..., Moden] -- (functions)
```

VAR Mode implies that the variable is represented in the mode denoted by *Mode*. On the other hand, *FUN [Mode₁, ..., Mode_n]* implies that:

- The variable is in the *V* (value) mode.
- The variable is a known function.
- The function takes n arguments, which are in modes *Mode₁*, *Mode₂*, etc.

We first present the monadic semantics for `let` bindings. E_e denotes the semantics of STG expressions. Similarly, we will use E_r , E_a and E_v to denote the semantics for closures (appearing on the right-hand side of `let` bindings), atoms, and variables, respectively.

$$\begin{aligned}
E_e[\llbracket \text{let } v = \lambda\{v_1, \dots, v_n\}. e_1 \text{ in } e_2 \rrbracket] = & \\
\{ \rho \leftarrow rdEnv; & \\
\text{let} & \\
r = E_r[\llbracket \lambda\{v_1, \dots, v_n\}. e_1 \rrbracket] & \\
r' = stdEntry \bar{v} \llbracket \lambda\{v_1, \dots, v_n\}. e_1 \rrbracket & \\
body = inEnv \rho [FUN [strictnessOf v_1, \dots, strictnessOf v_n] / \llbracket v \rrbracket] E_e[\llbracket e_2 \rrbracket] & \\
\text{in} & \\
\{ (\bar{v}, \bar{v}) \leftarrow \mathbf{fix} (\lambda(\bar{v}, \bar{v}). \{ x \leftarrow r; & \\
& x' \leftarrow r'; \\
& \text{return } (x, x') \}); \\
body \} \} &
\end{aligned}$$

We adopt the “natural semantics” style, which maps source language variables to meta language variables. The fixed-point operator `fix` has type $(a \rightarrow m a) \rightarrow m a$. Note that meta language variables \bar{v} and \bar{v} inside the argument of `fix` scope over r and r' , whereas \bar{v} and \bar{v} on the left-hand side of \leftarrow scope over *body*. The environment component above carries the strictness information, not runtime values.

The main difference between the above semantics and the source language semantics in Section 5.2 is that here we map the STG variable v to two meta language variables \bar{v} and \bar{v} . Meta language variable \bar{v} denotes the *optimized function entry point*; whereas \bar{v} denotes the *standard entry point*.

- The optimized function entry point allows us to pass multiple arguments in an uncurried form and utilize the strictness information. The optimized function entry point is only useful when all the arguments needed by the function are available at the call site.
- The standard entry point is used when we need to invoke an unknown function, or when not all arguments are available. The standard entry point expects all arguments to be passed in the cell mode.

$$\begin{aligned}
E_e \llbracket v \{a_1, \dots, a_n\} \rrbracket = & \\
& \{ \rho \leftarrow rdEnv; \\
& \text{case } \rho \llbracket v \rrbracket \text{ of} \\
& \quad \text{VAR } _ \rightarrow stdApp (E_v \llbracket v \rrbracket) \\
& \qquad \qquad \qquad [delay E_a \llbracket a_1 \rrbracket, \dots, delay E_a \llbracket a_n \rrbracket] \} \\
& \text{FUN } [m_1, \dots, m_k] \rightarrow \\
& \quad \text{if } k > n \text{ then} \\
& \qquad \quad stdApp (E_v \llbracket v \rrbracket) \\
& \qquad \qquad \quad [delay E_a \llbracket a_1 \rrbracket, \dots, delay E_a \llbracket a_n \rrbracket] \\
& \quad \text{else if } k = n \text{ then} \\
& \qquad \quad optApp \bar{v} [delay? m_1 E_a \llbracket a_1 \rrbracket, \dots, delay? m_k E_a \llbracket a_n \rrbracket] \\
& \quad \text{else} \\
& \qquad \quad stdApp (optApp \bar{v} [delay? m_1 E_a \llbracket a_1 \rrbracket, \dots, delay? m_k E_a \llbracket a_k \rrbracket]) \\
& \qquad \qquad \quad [delay E_a \llbracket a_{k+1} \rrbracket, \dots, delay E_a \llbracket a_n \rrbracket]
\end{aligned}$$

If the function is an unknown variable (the VAR case), then we evaluate the variable and follow the standard calling convention. If the variable is a known function, then we use the optimized or standard calling convention depending on the number of arguments we have received. To apply an optimized entry point, we simply evaluate all the arguments and pass them to the known function. Note that lazy evaluation semantics is preserved because x_1, \dots, x_n below may denote either values or cells.

$$optApp f [e_1, \dots, e_n] = \{x_1 \leftarrow e_1; \dots; x_n \leftarrow e_n; f(x_1, \dots, x_n)\}$$

Standard entry point application $stdApp$ will be defined later.

The $delay?$ function is a variation of $delay$. It only delays the computation if the first argument is C .

$$delay? V e = e$$

$$delay? C e = delay e$$

Primitive applications in the STG language are always saturated. We use the strictness information to determine whether an argument should be passed as a value or a cell. \bar{p} denotes that meta language support for the primitive p .

$$\begin{aligned}
 E_e[[p \{a_1, \dots, a_n\}]] &= \text{let} \\
 &\quad [m_1, \dots, m_n] = \text{strictnessOf } p \\
 &\quad \text{in} \\
 &\quad \{x_1 \leftarrow \text{delay? } m_1 E_a[[a_1]]; \dots \\
 &\quad \quad x_n \leftarrow \text{delay? } m_n E_a[[a_n]]; \\
 &\quad \quad \bar{p} [x_1, \dots, x_n]\}
 \end{aligned}$$

Atoms are either variables or literal constants. In the STG language, functions, primitives and constructors always receive atoms as arguments.

$$\begin{aligned}
 E_a[[v]] &= E_v[[v]] \\
 E_a[[l]] &= \text{return } \bar{l}
 \end{aligned}$$

Variables can denote either a known function, a value, or a cell. E_v looks up variable in the strictness environment, and inserts *force* if necessary:

$$\begin{aligned}
 E_v[[v]] &= \{\rho \leftarrow rdEnv; \\
 &\quad \text{case } \rho[[v]] \text{ of} \\
 &\quad \quad \text{FUN } _ \rightarrow \text{return } \bar{v} \\
 &\quad \quad \text{VAR } m \rightarrow \text{force? } m \bar{v}\}
 \end{aligned}$$

Note that \bar{v} is the standard entry point of functions. Function *force?* is a variation of *force*:

$$\begin{aligned}
 \text{force? } V v &= \text{return } v \\
 \text{force? } C c &= \text{force } c
 \end{aligned}$$

Using the above semantics, the example code in Section 5.4.1 can be translated into the following monadic code:

```

{ ( $\bar{f}$ ,  $\bar{\bar{f}}$ )  $\leftarrow$  fix ( $\lambda(\bar{f}, \bar{\bar{f}}). \{ x \leftarrow \text{return}(\lambda(\bar{x}, \bar{y}). \text{return}(\bar{x} + \bar{y}));$ 
     $x' \leftarrow \text{stdEntry } \bar{f} \llbracket \lambda\{x, y\}. + \{x, y\} \rrbracket;$ 
     $\text{return } (x, x') \}$ );
   $\text{stdApp } (\text{return } \bar{\bar{f}}) [\text{delay } (\text{return } 5)] \}$ 

```

In the above code, we have transformed away the static strictness environment using monad laws, environment axioms, and meta-language properties.

The complete listing of monadic semantics for the STG language is given in Appendix C.

5.4.4 Implementing Standard Entries

The standard entry point must be able to handle any number of arguments. Since the standard entry is used for unknown functions, there is no strictness information available. All arguments are passed in the cell mode.

We will describe two ways to create a standard entry. The easiest approach is to take advantage of curried functions in the meta language.

$$\text{stdEntry } f \llbracket \lambda\{v_1, \dots, v_n\}. e \rrbracket =$$

$$\text{return } (\lambda x_1 \dots \text{return } (\lambda x_n. \text{optApp } f [\text{force? } (\text{strictnessOf } v_1) x_1, \dots,$$

$$\text{force? } (\text{strictnessOf } v_n) x_n]))$$

A new closure is created each time an argument becomes available. When a standard entry is applied to a number of arguments, intermediate closures will be created and immediately consumed by the next application. The result of evaluating e_1, \dots, e_n must be cells.

$$\text{stdApp } e [e_1, \dots, e_n] = \{f \leftarrow e; x_1 \leftarrow e_1; f_1 \leftarrow f x_1; \dots; x_n \leftarrow e_n; f_{n-1} x_n\}$$

To reduce the number of intermediate closures, standard entries can instead take multiple arguments at once (in a vector, for example). The standard entry determines whether enough arguments are available.

- When given exactly the right number of arguments, we invoke the optimized entry point.
- When given more than enough arguments, we invoke the optimized entry point with just the necessary arguments. The result of the application is an unknown function, and is in turn applied to the remaining arguments.
- When given not enough arguments, we return a closure that expects further arguments. Additional arguments, when they become available later, will be appended to the existing arguments. The standard entry will then be invoked again with all the arguments.

This is similar to the approach taken in the STG machine, where each self-updatable closure checks if enough arguments have been received.

The *stdEntry'* function builds a standard entry as described above. $\#[\dots]$ denotes an ordered sequence. In practice, arguments can be passed in vectors.

```

stdEntry' f  $\llbracket \lambda \{v_1, \dots, v_n\}. e \rrbracket = f_{std}$  where
   $f_{std} = \text{return } (\lambda \#[x_1, \dots, x_m].$ 
    if  $m = n$  then optApp f  $[x_1, \dots, x_n]$ 
    else if  $m > n$  then
       $\{f' \leftarrow \text{optApp } f [x_1, \dots, x_n]; f' \#[x_{n+1}, \dots, x_m]\}$ 
    else
       $\text{return } (\lambda \#[y_1, \dots, y_k]. f_{std} \#[x_1, \dots, x_m, y_1, \dots, y_k])$ 

```

Applying a function built with *stdEntry'* is straightforward. We simply pass all available arguments to the function. Expressions e_1, \dots, e_n must evaluate to cells.

```

stdApp' e  $[e_1, \dots, e_n] = \{f \leftarrow e; x_1 \leftarrow e_1; \dots; x_n \leftarrow e_n; f \#[x_1, \dots, x_n]\}$ 

```


5.4.5 Connecting to the SML/NJ Back-end

STG code generated from the GHC front-end is targeted to the SML/NJ back-end in two steps:

First, we derive monadic intermediate code by unfolding the monadic semantics for a given program. The strictness environment does not depend on any values or computations that are only available at run-time. Monad laws, environment axioms and meta-language properties allow us to propagate the strictness information across the program and transform away the strictness environment.

Second, the resulting monadic code is transformed to ML code with extensions that allow us to perform unsafe type casting. The ML type system cannot efficiently express the constructs needed to implement updatable cells. The first component of a cell is a boolean flag bit. The second component, however, may be a value or a closure depending on whether the cell has been evaluated. Using the unsafe operations, we can efficiently implement the following `CELL` signature. It captures an abstract type `cell` and two operations `delay` and `force`:

```
signature CELL =
  sig
    type 'a cell
    val delay : (unit -> 'a) -> 'a cell
    val force : 'a cell -> 'a
  end (* CELL *)
```

Note that we use unsafe castings only in the implementation of the above signature, and only in places where we know that they are safe. This guarantees that no type violations will actually occur at runtime.

The example code in Section 5.4.1 is translated into the following ML code: (The standard entry is implemented using currying.)

```
let val rec f : int * int -> int = (* optimized entry *)
      fn (x, y) => x + y
  and f_std = (* standard entry *)
      fn e1 => fn e2 =>
```

```
        let val x1 = force e1
            val x2 = force e2
        in f (x1, x2)
        end
in let val x = delay 5
    in f_std x
    end
end
```

The code is a little more verbose than it needs to be, but this does not affect performance because the SML front-end transforms away all the extra `let` bindings.

The reason we generated ML source code, instead of, for example, the continuation-passing style (CPS) intermediate language, is that we could use ML's type checker to verify the type safety of the generated ML code. We have verified that the CPS code generated by the SML/NJ front-end was indeed similar to what we would manually generate.

Using these ideas, we built a retargeted Haskell compiler based on the GHC version 0.26 front-end and on SML/NJ version 108.9. To verify the correctness of our semantics and implementation, we compiled and ran a set of Haskell programs with our system. The programs are listed in Table 5.1. Other than one (the list-based *set* utilities), all Haskell programs are from the `nofib` test suite [Partain, 1992]. They range from small toy examples to relatively large applications. We compared the result of running the same program using both the original GHC compiler and the one retargeted to SML/NJ back-end. The outputs were verified to be exactly the same.

Although it is not the main purpose of this exercise to obtain an efficient Haskell compiler, it is interesting to measure the performance of our retargeted compiler. The programs were run on a SparcStation 20 with 128 MB of memory. Table 5.2 lists the median time of 5 runs. The number of seconds include both system and user time.

The column marked as `GHC` is the number of seconds it took for GHC 0.26 to

Program	Lines	Description
anna	9611	Strictness analyzer
boyer	1016	Standard theorem-prover benchmark
calc	1137	Arbitrary precision calculator
compress	821	LZW compression program
infer	585	Hindley-Milner type checker
primetest	280	Primality testing for large numbers
prolog	637	Prolog Interpreter
queens	14	N queens
rsa	97	RSA encryption
set	102	List-based set operations

Table 5.1: Sample programs

Program	GHC	Currying	Vector
anna	6.40	13.45	13.10
boyer	2.87	4.52	5.02
calc	11.11	12.17	12.10
compress	44.69	73.97	75.66
infer	4.94	9.04	10.83
primetest	61.64	105.02	106.71
prolog	17.88	33.56	32.88
queens	62.8	61.64	59.98
rsa	13.17	25.2	25.46
set	17.15	18.67	14.86

Table 5.2: Timing results of sample programs (in seconds)

execute the programs. The **Currying** column contains the timing figures obtained by implementing standard entry points using ML curried functions. The **Vector** column contains the timing of implementing standard entry points using ML vectors to hold multiple arguments. The **Currying** approach corresponds to *stdEntry*, whereas the **Vector** approach corresponds to *stdEntry'*.

Overall, the GHC out-performs the retargeted compiler by about 60%. Our system performs well in a number of simple tests (*queens* and *set*), and on the arbitrary precision calculator (*calc*). Both *queens* and *set* are computation bound. They require little garbage collection. The *calc* program is unique in that it spends most of its time in a few C functions implementing arbitrary-precision integer arithmetic.

We have not yet performed a thorough analysis of where the overhead of our system is in other tests. Preliminary investigation suggests that our system spends significantly more time on garbage collection than GHC. The SML/NJ runtime does not appear to be well-tuned to handle the frequent updates required by lazy evaluation. In addition, cell-mode data structures consume more memory than self-updating closures in GHC.

The **Currying** and **Vector** versions of standard entries do not make a significant difference in most test programs. On one hand, the **Vector** version of standard entry requires fewer closures to be created. On the other hand, packaging multiple arguments in vectors requires additional memory and CPU time. The **Vector** approach paid off in one test program (*set*). In that program a predicate function expecting two arguments is first applied to one argument, and later to another. The **Vector** approach saves time by not creating the intermediate closure.

Although there has been little work in optimizing our system, it performs reasonably well comparing to the hand-crafted GHC compiler. The competitive performance should be attributed to the efficiency and versatility of the SML/NJ back-end and runtime system.

Chapter 6

Related Work, Future Work and Conclusion

In this thesis, we have demonstrated how monads and monad transformers can be used to provide more modular specification of programming language features than traditional denotational semantics. In addition, we have shown how the modularity offered in our framework can provide better support for equational reasoning, program transformation, interpreter construction, and semantics-directed compilation. More specifically, the contributions of the work presented in this thesis are as follows:

- We have constructed modular semantic building blocks that support a wide variety of source language features, including arithmetic expressions, call-by-value, call-by-name, lazy evaluation, references and assignment, tracing, first-class continuation, and nondeterminism. Although each of these features has been modeled using monads before, it is the first time all of them fit into a single modular framework.
- We have solved a number of open problems in how to lift operations through monad transformers. We have extended Moggi's [Moggi, 1990] natural lift-

ing condition to higher-order types, making it possible to reason about the relatively complex operations related to environment and continuation. In addition, we have shown how liftings capture the interactions between various programming language features.

- We have implemented modular semantic building blocks and monad transformers in Gofer [Jones, 1991]. This is the first implementation of a full-featured modular monadic interpreter using a strongly-typed language.
- We have investigated high-level monadic properties of programming language features (for example, the environment axioms). We have applied these properties to construct modular proofs and to perform semantics-directed compilation.
- We have constructed a retargeted Haskell compiler that demonstrates a practical use of modular monadic semantics.

6.1 Related Work

Our work is built on a number of previous attempts to better organize modular semantics, to more effectively reason about programming languages, and to more efficiently compile higher-order programs.

6.1.1 Modular Semantics

The lack of modularity of traditional denotational semantics [Stoy, 1977] has long been recognized [Mosses, 1984] [Lee, 1989].

Moggi first proposed to use *monads* and *monad transformers* to structure denotational semantics. Wadler popularized Moggi's ideas in the functional programming community by showing how monads could be used in a variety of settings, including incorporating imperative features [Peyton Jones and Wadler, 1993] and building

modular interpreters [Wadler, 1992]. Wadler [King and Wadler, 1993] discussed the issues in combining monads. *Pseudomonads* [Steele Jr., 1994] were proposed as a way to compose monads and thus build up an interpreter from smaller parts. However, implementing pseudomonads in the Haskell [Hudak *et al.*, 1992] type system turned out to be problematic.

Returning to Moggi's original ideas, Espinosa formulated a system called *Semantic Lego* [Espinosa, 1993] [Espinosa, 1995]. Espinosa's Scheme-based system was the first modular interpreter that incorporated monad transformers. Among his contributions, Espinosa pointed out that pseudomonads were really just a special kind of *monad transformer*, first suggested by Moggi as a way to leave a "hole" in a monad for further extension. Espinosa's work reminded the programming language community—who had become distracted by the use of monads—that Moggi himself, responsible in many ways for the interest in monadic programming, had actually focussed more on the importance of monad transformers.

Related approaches to enhance modularity include composing monads [Jones and Duponcheel, 1993] and stratified monads [Espinosa, 1994].

This thesis was motivated by the above line of work, which led to the solution [Liang *et al.*, 1995] of a number of open issues in how to lift operations through monad transformers, as well as how to implement modular interpreters in a strongly-typed language.

6.1.2 Reasoning with Monads

In his original note [Moggi, 1990], Moggi raised the issue of reasoning in the monadic framework. The monadic framework has been used to specify state monad laws [Wadler, 1990], and to reason about exceptions [Spivey, 1990]. A related, but more general, framework to reason about states is mutable abstract data types (MADTs) [Hudak, 1992].

This thesis extends previous work by presenting the environment axioms [Liang

and Hudak, 1996]. In addition, we demonstrate how these axioms, together with monad laws, can be used to reason about programs in a modular way.

6.1.3 Semantics-directed Compilation

Early efforts [Wand, 1984] [Paulson, 1982] were based on traditional denotational semantics. The resulting compilers were inefficient.

Action Semantics [Mosses, 1992] allows modular specification of programming language semantics. Action semantics and a related approach [Lee, 1989] have been successfully used to generate efficient compilers. While action semantics is easy to construct, extend, understand and implement, we note the following comments ([Mosses, 1992], page 5):

“Although the foundations of action semantics are firm enough, the *theory* for reasoning about actions (and hence about programs) is still rather weak, and needs further development. This situation is in marked contrast to that of denotational semantics, where the theory is strong, but severe pragmatic difficulties hinder its application to realistic programming languages.”

Our work essentially attempts to formulate actions in a denotational semantics framework. Monad transformers roughly correspond to *facets* in action semantics, although issues such as concurrency are beyond the power of our approach.

A related approach [Meijer] is to combine the standard initial algebra semantics approach with aspects of Action Semantics to derive compilers from denotational semantics.

One application of partial evaluation [Jones *et al.*, 1989] is to generate compilers from interpreters. A partial evaluator has been successfully applied to an action interpreter [Bondorf and Palsberg, 1993], and similar results can be achieved with monadic interpreters [Danvy *et al.*, 1991] as well.

Staging transformations [Jørring and Scherlis, 1986] are a class of general program transformation techniques for separating a given computation into stages.

Monad transformers make computational stages somewhat more explicit by separating compile-time features, such as the environment, from run-time features.

There have been several successful efforts (including [Kelsey and Hudak, 1989], [Appel and Jim, 1989], and others) to build efficient compilers for higher-order languages by transforming the source language into continuation-passing style (CPS). The suitability of a monadic form as an intermediate form has been observed by many researchers (including, for example, [Sabry and Felleisen, 1992] and [Hatcliff and Danvy, 1994]).

6.2 Future Work

6.2.1 Theory of Programming Language Features

We have used monads and monad transformers to study programming language features and their interactions. Plenty of work remains on extending the theory to handle other useful features we have not covered. As a result, we may be able to better understand and implement these features.

6.2.2 Monadic Program Transformation

We have demonstrated that monadic code is particularly suitable for program transformation. Because monadic semantics is no more than an abstraction of traditional denotational semantics, all equational reasoning techniques apply. Monadic semantics can thus be used to facilitate various program transformation techniques such as partial evaluation.

6.2.3 A Common Back-end for Modern Languages

The experience of building a retargeted Haskell compiler suggests the feasibility of a common back-end for modern languages. An efficient, well-thought-out system

such as SML/NJ is a strong candidate to serve as a common back-end for a variety of modern languages.

We can further develop our monadic semantics based compilation method into a compiler construction tool for a common back-end.

6.2.4 Concurrency

Concurrency is an important feature in many modern languages such as JavaTM [Gosling *et al.*, 1996]. The monadic framework covers the properties of *callcc*. Since *callcc* captures the activities occur during a thread context switch, we expect the results related to *callcc* will be useful in reasoning about multi-threaded concurrent systems.

6.3 Conclusion

We have demonstrated the power of modular monadic semantics in two ways. First, it is a powerful technique to specify and reason about programming language features. Second, it can be used in practice to construct modular interpreters and perform semantics-directed compilation.

The key benefit of our approach is modularity. The underlying mechanism is monad-based abstraction. Modular monadic semantics helps to bridge the gap between programming language theory and the complex practical languages.

Bibliography

- [Appel and Jim, 1989] Andrew W. Appel and Trevor Jim. Continuation-passing, closure-passing style. In *ACM Symposium on Principles of Programming Languages*, pages 193–302, January 1989.
- [Appel, 1992] Andrew W. Appel. *Compiling with Continuations*. Cambridge University Press, 1992.
- [Bloss *et al.*, 1988] Adrienne Bloss, Paul Hudak, and Jonathan Young. Code optimization for lazy evaluation. *Lisp and Symbolic Computation*, 1(1):147–164, 1988.
- [Bondorf and Palsberg, 1993] Anders Bondorf and Jens Palsberg. Compiling actions by partial evaluation. In *FPCA '93: Conference on Functional Programming Languages and Computer Architecture, Copenhagen, Denmark*, pages 308–317, New York, June 1993. ACM Press.
- [Clément *et al.*, 1986] D. Clément, Joëlle Despeyroux, Thierry Despeyroux, and Gilles Kahn. A simple applicative language: Mini-ML. In *Proceedings of the 1986 ACM Symposium on Lisp and Functional Programming*, pages 13–27, 1986.
- [Clinger and Rees, 1991] William Clinger and Jonathan Rees. Revised⁴ report on the algorithmic language Scheme. Available via World Wide Web as <http://www-swiss.ai.mit.edu/ftplib/scheme-reports/r4rs.ps>, November 1991.
- [Danvy *et al.*, 1991] Olivier Danvy, Jürgen Koslowski, and Karoline Malmkjær. Compiling monads. Technical Report CIS-92-3, Kansas State University, December 1991.
- [Espinosa, 1993] David Espinosa. Modular denotational semantics. Unpublished manuscript, December 1993.

- [Espinosa, 1994] David Espinosa. Building interpreters by transforming stratified monads. Unpublished manuscript, ftp from alt-dorf.ai.mit.edu:pub/dae, June 1994.
- [Espinosa, 1995] David Espinosa. *Semantic Lego*. PhD thesis, Columbia University, 1995.
- [Felleisen and Hieb, 1992] Matthias Felleisen and Robert Hieb. The revised report on the syntactic theories of sequential control and state. *Theoretical Computer Science*, 103:235–271, 1992.
- [Felleisen *et al.*, 1986] Matthias Felleisen, Daniel P. Friedman, Eugene Kohlbecker, and Bruce Duba. Reasoning with continuations. In *Proceedings of the Symposium on Logic in Computer Science*, pages 131–141, Cambridge, Massachusetts, June 1986. IEEE.
- [Gosling *et al.*, 1996] James Gosling, Bill Joy, and Guy Steele. *The Java™ Programming Language*. The Java Series. Addison-Wesley, 1996.
- [Hatcliff and Danvy, 1994] John Hatcliff and Olivier Danvy. A generic account of continuation-passing styles. In *21st ACM Symposium on Principles of Programming Languages (POPL '94)*, Portland, Oregon, pages 458–471, New York, January 1994. ACM Press.
- [Hudak and Bloss, 1985] Paul Hudak and Adriene Bloss. The aggregate update problem in functional programming languages. In *Conference Record of the Twelfth ACM Symposium on Principles of Programming Languages*, pages 300–314. Association for Computing Machinery, January 1985.
- [Hudak *et al.*, 1992] Paul Hudak, Simon Peyton Jones, and Philip Wadler. Report on the programming language Haskell: a non-strict, purely functional language, version 1.2. Technical Report YALEU/DCS/RR-777, Yale University Department of Computer Science, March 1992. Also in ACM SIGPLAN Notices, Vol. 27(5), May 1992.
- [Hudak, 1992] Paul Hudak. Mutable abstract datatypes. Research Report YALEU/DCS/RR-914, Yale University Department of Computer Science, December 1992.

- [Jones and Duponcheel, 1993] Mark P. Jones and Luc Duponcheel. Composing monads. Research Report YALEU/DCS/RR-1004, Yale University Department of Computer Science, New Haven, Connecticut, December 1993.
- [Jones *et al.*, 1989] Neil D. Jones, Peter Sestoft, and Harald Søndergaard. Mix: a self-applicable partial evaluator for experiments in compiler generation. *Lisp and Symbolic Computation*, 2:9–50, 1989.
- [Jones, 1991] Mark P. Jones. Introduction to gofer 2.20. Ftp from nebula.cs.yale.edu in the directory pub/haskell/gofer, September 1991.
- [Jones, 1993] Mark P. Jones. A system of constructor classes: Overloading and implicit higher-order polymorphism. In *FPCA '93: Conference on Functional Programming Languages and Computer Architecture, Copenhagen, Denmark*, pages 52–61, New York, June 1993. ACM Press.
- [Jørring and Scherlis, 1986] Ulrik Jørring and William Scherlis. Compilers and staging transformations. In *Proceedings Thirteenth ACM Symposium on Principles of Programming Languages, St. Petersburg, Florida*, pages 86–96, 1986.
- [Kelsey and Hudak, 1989] Richard Kelsey and Paul Hudak. Realistic compilation by program transformation. In *ACM Symposium on Principles of Programming Languages*, pages 181–192, January 1989.
- [King and Wadler, 1993] David J. King and Philip Wadler. Combining monads. In John Launchbury and Patrick Sansom, editors, *Functional Programming, Glasgow 1992*, pages 134–143, New York, 1993. Springer-Verlag.
- [Kishon *et al.*, 1991] Amir Kishon, Paul Hudak, and Charles Consel. Monitoring semantics: A formal framework for specifying, implementing and reasoning about execution monitors. In *Proceedings of the ACM SIGPLAN '91 Conference on Programming Language Design and Implementation*, pages 338–352, June 1991.
- [Kranz *et al.*, 1986] David Kranz, Richard Kelsey, Jonathan Rees, Paul Hudak, James Philbin, and Norman Adams. ORBIT: An optimizing compiler for Scheme. In *Proceedings SIGPLAN '86 Symposium*

- on Compiler Construction*, pages 219–233, July 1986. SIGPLAN Notices Volume 21, Number 7.
- [Lee, 1989] Peter Lee. *Realistic Compiler Generation*. Foundations of Computing. MIT Press, 1989.
- [Liang and Hudak, 1996] Sheng Liang and Paul Hudak. Modular denotational semantics for compiler construction. In Hanne Riis Nielson, editor, *ESOP '96: 6th European Symposium on Programming, Linköping, Sweden, Proceedings*, pages 219–234, New York, April 1996. Springer-Verlag. Lecture Notes in Computer Science 1058.
- [Liang *et al.*, 1995] Sheng Liang, Paul Hudak, and Mark Jones. Monad transformers and modular interpreters. In *22nd ACM Symposium on Principles of Programming Languages (POPL '95), San Francisco, California*, New York, January 1995. ACM Press.
- [Mac Lane, 1971] Saunders Mac Lane. *Categories for the Working Mathematician*. Graduate Texts in Mathematics. Springer-Verlag, 1971.
- [Meijer] Erik Meijer. More advice on proving a compiler correct: Improving a correct compiler. Submitted to *Journal of Functional Programming*.
- [Milner *et al.*, 1990] Robin Milner, Mads Tofte, and Robert Harper. *The Definition of Standard ML*. MIT Press, 1990.
- [Moggi, 1990] Eugenio Moggi. An abstract view of programming languages. Technical Report ECS-LFCS-90-113, Laboratory for Foundations of Computer Science, University of Edinburgh, Edinburgh, Scotland, 1990.
- [Mosses, 1984] Peter D. Mosses. A basic abstract semantic algebra. In Gilles Kahn, David B. MacQueen, and Gordon D. Plotkin, editors, *Semantics of Data Types: International Symposium, Sophia-Antipolis, France*, pages 87–107. Springer-Verlag, June 1984. Lecture Notes in Computer Science 173.
- [Mosses, 1992] Peter D. Mosses. *Action Semantics*, volume 26 of *Cambridge Tracts in Theoretical Computer Science*. Cambridge University Press, 1992.

- [Partain, 1992] Will Partain. The `nofib` benchmark suite of haskell programs. University of Glasgow, 1992.
- [Paterson, 1995] Ross A. Paterson, 1995. private communication.
- [Paulson, 1982] Laurence C. Paulson. A semantics-directed compiler generator. In *Proceedings of the Ninth ACM Symposium on Principles of Programming Languages, Albuquerque, New Mexico*, pages 224–233, 1982.
- [Peterson and Hammond, 1996] John Peterson and Kevin Hammond. Report on the programming language Haskell: a non-strict, purely functional language, version 1.3. Technical Report YALEU/DCS/RR-1106, Yale University Department of Computer Science, May 1996.
- [Peyton Jones and Wadler, 1993] Simon Peyton Jones and Philip Wadler. Imperative functional programming. In *Proceedings 20th Symposium on Principles of Programming Languages*, pages 71–84. ACM, January 1993.
- [Peyton Jones, 1992] S.L. Peyton Jones. Implementing lazy functional languages on stock hardware: The Spineless Tagless G-machine. *Journal of Functional Programming*, 2(2):127–202, April 1992.
- [Sabry and Felleisen, 1992] Amr Sabry and Matthias Felleisen. Reasoning about programs in continuation-passing style. In *Proceedings of the 1992 ACM Conference on LISP and Functional Programming*, pages 288–298. ACM Press, June 1992.
- [Spivey, 1990] Michael Spivey. A functional theory of exceptions. *Science of Computer Programming*, 14(1):25–42, June 1990.
- [Steele Jr., 1994] Guy L. Steele Jr. Building interpreters by composing monads. In *Conference Record of POPL '94: 21st ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, Portland, Oregon*, pages 472–492, New York, January 1994. ACM Press.
- [Stoy, 1977] Joseph Stoy. *Denotational Semantics: The Scott-Strachey Approach to Programming Language Theory*. MIT Press, 1977.

- [Tolmach and Appel, 1990] Andrew P. Tolmach and Andrew W. Appel. Debugging standard ML without reverse engineering. In *Proceedings of the 1990 ACM Conference on Lisp and Functional Programming*, Nice, France, June 1990.
- [Wadler, 1989] Philip L. Wadler. Theorems for free! In *Fourth Symposium on Functional Programming Languages and Computer Architecture*. ACM, September 1989. London.
- [Wadler, 1990] Philip L. Wadler. Comprehending monads. In *Proceedings of the 1990 ACM Conference on Lisp and Functional Programming*, 1990.
- [Wadler, 1992] Philip Wadler. The essence of functional programming. In *Conference Record of the Nineteenth Annual ACM Symposium on Principles of Programming Languages, Albuquerque, New Mexico*, pages 1–14, January 1992.
- [Wand, 1984] Mitchell Wand. A semantic prototyping system. *SIGPLAN Notices, ACM Symposium on Compiler Construction*, 19(6):213–221, 1984.
- [Wand, 1990] Mitchell Wand. A short proof of the lexical addressing algorithm. *Information Processing Letters*, 35:1–5, June 1990.
- [Yale Haskell Group, 1994] The Yale Haskell Group. The Yale Haskell users manual. Version Y2.2, September 1994.

Appendix A

Proofs

This appendix contains detailed proofs for many of the results given in the body of this thesis. For convenience, we repeat the statement of each result at the beginning of the corresponding proof.

Proposition 3.1.2 $EnvT\ r$, $StateT\ s$, $ErrT$, and $ContT\ c$ are monad transformers.

Proof: We need to show that the corresponding *lift* functions are monad morphisms.

Case $EnvT\ r$:

$$\begin{aligned} return_{EnvT\ r\ m}\ e &= \lambda\rho. return_m\ e && (return_{EnvT\ r\ m}) \\ &= lift_{EnvT\ r}\ (return_m\ e) && (lift_{EnvT\ r}) \end{aligned}$$

$$\begin{aligned} bind_{EnvT\ r\ m}\ (lift_{EnvT\ r}\ m)\ (\lambda a. (lift_{EnvT\ r}\ (ka))) & \\ &= bind_{EnvT\ r\ m}\ (\lambda\rho'. m)\ (\lambda a. (\lambda\rho'. ka))\ (lift_{EnvT\ r}) \\ &= \lambda\rho. \{a \leftarrow (\lambda\rho'. m)\ \rho; (\lambda\rho'. ka)\rho\}_m && (bind_{EnvT\ r\ m}) \\ &= \lambda\rho. bind_m\ e\ k && (\beta) \\ &= lift_{EnvT\ r}\ (bind_m\ e\ k) && (lift_{EnvT\ r}) \end{aligned}$$

Case $StateT\ s$:

$$\begin{aligned} return_{StateT\ s\ m}\ e & \\ &= \lambda s. return_m\ (s, e) && (return_{StateT\ s\ m}) \\ &= \lambda s. \{a \leftarrow return_m\ e; return_m\ (s, a)\}_m && (\text{left unit}) \\ &= lift_{StateT\ s}\ (return_m\ e) && (lift_{StateT\ s}) \end{aligned}$$

$$\begin{aligned}
& \text{bind}_{\text{StateT } s \ m} (\text{lift}_{\text{StateT } s} e) (\lambda a. \text{lift}_{\text{StateT } s} (ka)) \\
&= \lambda s. \{(s', a) \leftarrow \text{lift}_{\text{StateT } s} e \ s; \text{lift}_{\text{StateT } s} (ka) \ s'\}_m \quad (\text{bind}_{\text{StateT } s \ m}) \\
&= \lambda s. \{(s', a') \leftarrow \{a \leftarrow e; \text{return}_m(s, a)\}_m; \\
&\quad b \leftarrow ka'; \text{return}_m(s', b)\}_m \quad (\text{lift}_{\text{StateT } s}) \\
&= \lambda s. \{a \leftarrow e; (s', a') \leftarrow \text{return}_m(s, a); \\
&\quad b \leftarrow ka'; \text{return}_m(s', b)\}_m \quad (\text{associativity}) \\
&= \lambda s. \{a \leftarrow e; b \leftarrow ka; \text{return}_m(s, b)\}_m \quad (\text{left unit}) \\
&= \text{lift}_{\text{StateT } s} (\text{bind}_m e \ k) \quad (\text{lift}_{\text{StateT } s})
\end{aligned}$$

Case ErrT:

$$\begin{aligned}
& \text{return}_{\text{ErrT } m} e \\
&= \text{return}_m(\text{Ok } e) \quad (\text{return}_{\text{ErrT } m}) \\
&= \{a \leftarrow \text{return}_m e; \text{return}_m(\text{Ok } a)\}_m \quad (\text{left unit}) \\
&= \text{map}_m \text{Ok} (\text{return}_m e) \quad (\text{map}_m) \\
&= \text{lift}_{\text{ErrT}} (\text{return}_m e) \quad (\text{lift}_{\text{ErrT}}) \\
& \\
& \text{bind}_{\text{ErrT } m} (\text{lift}_{\text{ErrT}} e) (\lambda a. \text{lift}_{\text{ErrT}} (ka)) \\
&= \{a' \leftarrow \text{lift}_{\text{ErrT}} e; \\
&\quad \text{case } a' \text{ of } \text{Ok } a \rightarrow \text{lift}_{\text{ErrT}}(ka) \dots\}_m \quad (\text{bind}_{\text{ErrT } m}) \\
&= \{a' \leftarrow \{a \leftarrow e; \text{return}_m(\text{Ok } a)\}_m; \\
&\quad \text{case } a' \text{ of } \text{Ok } a \rightarrow \{b \leftarrow ka; \text{return}_m(\text{Ok } b)\}_m \dots\}_m \quad (\text{lift}_{\text{ErrT}}, \text{map}_m) \\
&= \{a \leftarrow e; a' \leftarrow \text{return}_m(\text{Ok } a); \\
&\quad \text{case } a' \text{ of } \text{Ok } a \rightarrow \{b \leftarrow ka; \text{return}_m(\text{Ok } b)\}_m \dots\}_m \quad (\text{associativity}) \\
&= \{a \leftarrow e; b \leftarrow ka; \text{return}_m(\text{Ok } b)\}_m \quad (\text{left unit}) \\
&= \text{lift}_{\text{ErrT}} (\text{bind}_m e \ k) \quad (\text{lift}_{\text{ErrT}}, \text{map}_m)
\end{aligned}$$

Case ContT:

$$\begin{aligned}
& \text{return}_{\text{ContT } c \ m} e \\
&= \lambda k. ke \quad (\text{return}_{\text{ContT } c \ m}) \\
&= \text{bind}_m (\text{return}_m e) \quad (\text{left unit}) \\
&= \text{lift}_{\text{ContT } c} (\text{return}_m e) \quad (\text{lift}_{\text{ContT } c})
\end{aligned}$$

$$\begin{aligned}
& \text{bind}_{\text{ContT } c \ m} (\text{lift}_{\text{ContT } c} e) (\lambda a. \text{lift}_{\text{ContT } c} (fa)) \\
&= \lambda k. (\text{lift}_{\text{ContT } c} e) (\lambda a. \text{lift}_{\text{ContT } c} (fa) k) \quad (\text{bind}_{\text{ContT } c \ m}) \\
&= \lambda k. \text{bind}_m e (\lambda a. \text{bind}_m (fa) k) \quad (\text{lift}_{\text{ContT } c}) \\
&= \lambda k. \text{bind}_m (\text{bind}_m e f) k \quad (\text{associativity}) \\
&= \text{lift}_{\text{ContT } c} (\text{bind}_m e f) \quad (\text{lift}_{\text{ContT } c})
\end{aligned}$$

Proposition 3.2.1 The environment operations, $rdEnv$ and $inEnv$ satisfy the following axioms:

$$\begin{aligned}
(\text{inEnv } \rho) \cdot \text{return} &= \text{return} && (\text{unit}) \\
\text{inEnv } \rho \{v \leftarrow e_1; e_2\} &= \{v \leftarrow \text{inEnv } \rho e_1; \text{inEnv } \rho e_2\} && (\text{distribution}) \\
\text{inEnv } \rho \text{rdEnv} &= \text{return } \rho && (\text{cancellation}) \\
\text{inEnv } \rho' (\text{inEnv } \rho e) &= \text{inEnv } \rho e && (\text{overriding})
\end{aligned}$$

Proof: We verify that: 1) $inEnv$ and $rdEnv$ satisfy the axioms after being introduced by $EnvT$, and that: 2) the axioms are preserved through $EnvT$, $StateT$, and $ErrT$. (There is no lifting of $inEnv$ through $ContT$.)

Base case:

$$\begin{aligned}
& \text{inEnv}_{\text{EnvT } r \ m} \rho (\text{return}_{\text{EnvT } r \ m} x) \\
&= \lambda \rho'. (\text{return}_{\text{EnvT } r \ m} x) \rho \quad (\text{inEnv}_{\text{EnvT } r \ m}) \\
&= \lambda \rho'. (\lambda \rho''. \text{return}_m x) \rho \quad (\text{return}_{\text{EnvT } r \ m}) \\
&= \lambda \rho'. \text{return}_m x \quad (\beta) \\
&= \text{return}_{\text{EnvT } r \ m} x \quad (\text{return}_{\text{EnvT } r \ m})
\end{aligned}$$

$$\begin{aligned}
& \text{inEnv}_{\text{EnvT } r \ m} \rho \{v \leftarrow e_1; e_2\}_{\text{EnvT } r \ m} \\
&= \lambda \rho'. \{v \leftarrow e_1; e_2\}_{\text{EnvT } r \ m} \rho \quad (\text{inEnv}_{\text{EnvT } r \ m}) \\
&= \lambda \rho'. \{v \leftarrow e_1 \rho; e_2 \rho\}_m \quad (\text{bind}_{\text{EnvT } r \ m}) \\
&= \lambda \rho'. \{v \leftarrow (\lambda \rho''. e_1 \rho) \rho'; (\lambda \rho''. e_2 \rho) \rho'\}_m \\
&= \lambda \rho'. \{v \leftarrow \text{inEnv}_{\text{EnvT } r \ m} \rho e_1 \rho'; \text{inEnv}_{\text{EnvT } r \ m} \rho e_2 \rho'\}_m \quad (\text{inEnv}_{\text{EnvT } r \ m}) \\
&= \{v \leftarrow \text{inEnv}_{\text{EnvT } r \ m} \rho e_1; \text{inEnv}_{\text{EnvT } r \ m} \rho e_2\}_{\text{EnvT } r \ m} \quad (\text{bind}_{\text{EnvT } r \ m})
\end{aligned}$$

$$\begin{aligned}
& inEnv_{EnvT \ r \ m} \rho \ rdEnv_{EnvT \ r \ m} \\
&= \lambda \rho'. rdEnv_{EnvT \ r \ m} \rho && (inEnv_{EnvT \ r \ m}) \\
&= \lambda \rho'. (\lambda \rho. return_m \rho) \rho && (rdEnv_{EnvT \ r \ m}) \\
&= \lambda \rho'. return_m \rho \\
&= return_{EnvT \ r \ m} \rho && (return_{EnvT \ r \ m})
\end{aligned}$$

$$\begin{aligned}
& inEnv_{EnvT \ r \ m} \rho' (inEnv_{EnvT \ r \ m} \rho e) \\
&= \lambda \rho''. inEnv_{EnvT \ r \ m} \rho e \rho' && (inEnv_{EnvT \ r \ m}) \\
&= \lambda \rho''. (\lambda \rho'. e \rho) \rho' && (inEnv_{EnvT \ r \ m}) \\
&= \lambda \rho''. e \rho && (\beta) \\
&= inEnv_{EnvT \ r \ m} \rho e && (inEnv_{EnvT \ r \ m})
\end{aligned}$$

Case $EnvT \ r'$:

$$\begin{aligned}
& inEnv_{EnvT \ r' \ m} \rho (return_{EnvT \ r' \ m} x) \\
&= \lambda \rho'. inEnv_m \rho (return_{EnvT \ r' \ m} x \rho') && (inEnv_{EnvT \ r' \ m}) \\
&= \lambda \rho'. inEnv_m \rho ((\lambda \rho''. return_m x) \rho') && (return_{EnvT \ r' \ m}) \\
&= \lambda \rho'. inEnv_m \rho (return_m x) && (\beta) \\
&= \lambda \rho'. return_m x && (ind. hypo.) \\
&= return_{EnvT \ r' \ m} x && (return_{EnvT \ r' \ m})
\end{aligned}$$

$$\begin{aligned}
& inEnv_{EnvT \ r' \ m} \rho \{v \leftarrow e_1; e_2\}_{EnvT \ r' \ m} \\
&= \lambda \rho'. inEnv_m \rho (\{v \leftarrow e_1; e_2\}_{EnvT \ r' \ m} \rho') && (inEnv_{EnvT \ r' \ m}) \\
&= \lambda \rho'. inEnv_m \rho \{v \leftarrow e_1 \rho'; e_2 \rho'\}_m && (bind_{EnvT \ r' \ m}) \\
&= \lambda \rho'. \{v \leftarrow inEnv_m \rho (e_1 \rho'); inEnv_m \rho (e_2 \rho')\}_m && (ind. hypo.) \\
&= \lambda \rho'. \{v \leftarrow inEnv_{EnvT \ r' \ m} \rho e_1 \rho'; inEnv_{EnvT \ r' \ m} \rho e_2 \rho'\}_m && (inEnv_{EnvT \ r' \ m}) \\
&= \{v \leftarrow inEnv_{EnvT \ r' \ m} \rho e_1; inEnv_{EnvT \ r' \ m} \rho e_2\}_{EnvT \ r' \ m} && (bind_{EnvT \ r' \ m})
\end{aligned}$$

$$\begin{aligned}
& inEnv_{EnvT \ r' \ m} \rho \ rdEnv_{EnvT \ r' \ m} \\
&= \lambda \rho'. inEnv_m \rho \ (rdEnv_{EnvT \ r' \ m} \rho') && (inEnv_{EnvT \ r' \ m}) \\
&= \lambda \rho'. inEnv_m \rho \ (\lambda \rho''. rdEnv_m \rho') && (rdEnv_{EnvT \ r' \ m}) \\
&= \lambda \rho'. inEnv_m \rho \ rdEnv_m && (\beta) \\
&= \lambda \rho'. return_m \rho && (\text{ind. hypo.}) \\
&= return_{EnvT \ r' \ m} \rho && (return_{EnvT \ r' \ m})
\end{aligned}$$

$$\begin{aligned}
& inEnv_{EnvT \ r' \ m} \rho' \ (inEnv_{EnvT \ r' \ m} \rho \ e) \\
&= \lambda \rho''. inEnv_m \rho' \ (inEnv_{EnvT \ r' \ m} \rho \ e \ \rho'') && (inEnv_{EnvT \ r' \ m}) \\
&= \lambda \rho''. inEnv_m \rho' \ ((\lambda \rho'. inEnv_m \rho \ (e \ \rho')) \ \rho'') && (inEnv_{EnvT \ r' \ m}) \\
&= \lambda \rho''. inEnv_m \rho' \ (inEnv_m \rho \ (e \ \rho'')) && (\beta) \\
&= \lambda \rho''. inEnv_m \rho \ (e \ \rho'') && (\text{ind. hypo.}) \\
&= inEnv_{EnvT \ r' \ m} \rho \ e && (inEnv_{EnvT \ r' \ m})
\end{aligned}$$

Case StateT s:

$$\begin{aligned}
& inEnv_{StateT \ s \ m} \rho \ (return_{StateT \ s \ m} \ x) \\
&= \lambda s. inEnv_m \rho \ (return_{StateT \ s \ m} \ x \ s) && (inEnv_{StateT \ s \ m}) \\
&= \lambda s. inEnv_m \rho \ ((\lambda s. return_m \ (s, x)) \ s) && (return_{StateT \ s \ m}) \\
&= \lambda s. inEnv_m \rho \ (return_m \ (s, x)) && (\beta) \\
&= \lambda s. return_m \ (s, x) && (\text{ind. hypo.}) \\
&= return_{StateT \ s \ m} \ x && (return_{StateT \ s \ m})
\end{aligned}$$

$$\begin{aligned}
& inEnv_{StateT \ s \ m} \rho \ \{v \leftarrow e_1; e_2\}_{StateT \ s \ m} \\
&= \lambda s. inEnv_m \rho \ (\{v \leftarrow e_1; e_2\}_{StateT \ s \ m} \ s) && (inEnv_{StateT \ s \ m}) \\
&= \lambda s. inEnv_m \rho \ \{(s', v) \leftarrow e_1 \ s; e_2 \ s'\}_m && (bind_{StateT \ s \ m}) \\
&= \lambda s. \{(s', v) \leftarrow inEnv_m \rho \ (e_1 \ s); inEnv_m \rho \ (e_2 \ s')\}_m && (\text{ind. hypo.}) \\
&= \lambda s. \{(s', v) \leftarrow inEnv_{StateT \ s \ m} \rho \ e_1 \ s; inEnv_{StateT \ s \ m} \rho \ e_2 \ s'\}_m && (inEnv_{StateT \ s \ m}) \\
&= \{v \leftarrow inEnv_{StateT \ s \ m} \rho \ e_1; inEnv_{StateT \ s \ m} \rho \ e_2\}_{StateT \ s \ m} && (bind_{StateT \ s \ m})
\end{aligned}$$

$$\begin{aligned}
& \mathit{inEnv}_{\text{StateT } s \ m} \ \rho \ \mathit{rdEnv}_{\text{StateT } s \ m} \\
&= \lambda s. \mathit{inEnv}_m \ \rho \ (\mathit{rdEnv}_{\text{StateT } s \ m} \ s) && (\mathit{inEnv}_{\text{StateT } s \ m}) \\
&= \lambda s. \mathit{inEnv}_m \ \rho \ \{\rho' \leftarrow \mathit{rdEnv}_m; \mathit{return}_m(s, \rho')\}_m && (\mathit{rdEnv}_{\text{StateT } s \ m}) \\
&= \lambda s. \{\rho' \leftarrow \mathit{inEnv}_m \ \rho \ \mathit{rdEnv}_m; \mathit{inEnv}_m \ \rho \ \mathit{return}_m(s, \rho')\}_m && (\text{ind. hypo.}) \\
&= \lambda s. \{\rho' \leftarrow \mathit{return}_m; \mathit{return}_m(s, \rho')\}_m && (\text{ind. hypo.}) \\
&= \lambda s. \mathit{return}_m(s, \rho) && (\text{left unit}) \\
&= \mathit{return}_{\text{StateT } s \ m} \ \rho && (\mathit{return}_{\text{StateT } s \ m})
\end{aligned}$$

$$\begin{aligned}
& \mathit{inEnv}_{\text{StateT } s \ m} \ \rho' \ (\mathit{inEnv}_{\text{StateT } s \ m} \ \rho \ e) \\
&= \lambda s. \mathit{inEnv}_m \ \rho' \ (\mathit{inEnv}_{\text{StateT } s \ m} \ \rho \ e \ s) && (\mathit{inEnv}_{\text{StateT } s \ m}) \\
&= \lambda s. \mathit{inEnv}_m \ \rho' \ ((\lambda s'. \mathit{inEnv}_m \ \rho \ (e \ s')) \ s) && (\mathit{inEnv}_{\text{StateT } s \ m}) \\
&= \lambda s. \mathit{inEnv}_m \ \rho' \ (\mathit{inEnv}_m \ \rho \ (e \ s)) && (\beta) \\
&= \lambda s. \mathit{inEnv}_m \ \rho \ (e \ s) && (\text{ind. hypo.}) \\
&= \mathit{inEnv}_{\text{StateT } s \ m} \ \rho \ e && (\mathit{inEnv}_{\text{StateT } s \ m})
\end{aligned}$$

Case ErrT:

$$\begin{aligned}
& \mathit{inEnv}_{\text{ErrT } m} \ \rho \ (\mathit{return}_{\text{ErrT } m} \ x) \\
&= \mathit{inEnv}_m \ \rho \ (\mathit{return}_{\text{ErrT } m} \ x) && (\mathit{inEnv}_{\text{ErrT } m}) \\
&= \mathit{inEnv}_m \ \rho \ (\mathit{return}_m(\text{Ok } x)) && (\mathit{return}_{\text{ErrT } m}) \\
&= \mathit{return}_m(\text{Ok } x) && (\text{ind. hypo.}) \\
&= \mathit{return}_{\text{ErrT } m} \ x && (\mathit{return}_{\text{ErrT } m})
\end{aligned}$$

$$\begin{aligned}
& inEnv_{ErrT\ m}\ \rho\ \{v \leftarrow e_1; e_2\}_{ErrT\ m} \\
&= inEnv_m\ \rho\ (\{v \leftarrow e_1; e_2\}_{ErrT\ m\ s}) && (inEnv_{ErrT\ m}) \\
&= inEnv_m\ \rho\ \{a \leftarrow e_1; \\
&\quad \text{case } a \text{ of} \\
&\quad\quad Ok\ v \rightarrow e_2 \\
&\quad\quad Err\ s \rightarrow return_m(Err\ s)\}_m \\
&= \{a \leftarrow inEnv_m\ \rho\ e_1; && (\text{ind. hypo.}) \\
&\quad \text{case } a \text{ of} \\
&\quad\quad Ok\ v \rightarrow inEnv_m\ \rho\ e_2 \\
&\quad\quad Err\ s \rightarrow inEnv_m\ \rho\ return_m(Err\ s)\}_m \\
&= \{a \leftarrow inEnv_m\ \rho\ e_1; && (\text{ind. hypo.}) \\
&\quad \text{case } a \text{ of} \\
&\quad\quad Ok\ v \rightarrow inEnv_m\ \rho\ e_2 \\
&\quad\quad Err\ s \rightarrow return_m(Err\ s)\}_m \\
&= \{v \leftarrow inEnv_{ErrT\ m}\ \rho\ e_1; inEnv_{ErrT\ m}\ \rho\ e_2\}_{ErrT\ m} && (bind_{ErrT\ m})
\end{aligned}$$

$$\begin{aligned}
& inEnv_{ErrT\ m}\ \rho\ rdEnv_{ErrT\ m} \\
&= inEnv_m\ \rho\ rdEnv_{ErrT\ m} && (inEnv_{ErrT\ m}) \\
&= inEnv_m\ \rho\ \{\rho' \leftarrow rdEnv_m; return_m(Ok\ \rho')\}_m && (rdEnv_{ErrT\ m}) \\
&= \{\rho' \leftarrow inEnv_m\ \rho\ rdEnv_m; inEnv_m\ \rho\ return_m(Ok\ \rho')\}_m && (\text{ind. hypo.}) \\
&= \{\rho' \leftarrow return_m\ \rho; return_m(Ok\ \rho')\}_m && (\text{ind. hypo.}) \\
&= return_m(Ok\ \rho) && (\text{left unit}) \\
&= return_{ErrT\ m}\ \rho && (return_{ErrT\ m})
\end{aligned}$$

$$\begin{aligned}
& inEnv_{ErrT\ m}\ \rho' (inEnv_{ErrT\ m}\ \rho\ e) \\
&= inEnv_m\ \rho' (inEnv_m\ \rho\ e) && (inEnv_{ErrT\ m}) \\
&= inEnv_m\ \rho\ e && (\text{ind. hypo.}) \\
&= inEnv_{ErrT\ m}\ \rho\ e && (inEnv_{ErrT\ m})
\end{aligned}$$

Lemma 3.3.4

$$\begin{aligned} & \forall g, h, f, f', s_0. \\ & (\forall k, f'(\lambda x. \lambda s. \text{map}(\lambda x. h(s, x))(kx))s_0 = \text{map } g(fk) \Rightarrow \\ & \quad \text{callcc}(\lambda k. f'(\lambda x. \lambda s. k(gx))s_0) = \text{map } g(\text{callcc } f)) \end{aligned}$$

Proof: We establish the lemma by covering the cases when *callcc* was first introduced by *ContT* and lifted through *EnvT*, *StateT*, and *ErrT*. (There is no lifting of *callcc* through *ContT*.)

Base case:

$$\begin{aligned} & \text{callcc}(\lambda k. f'(\lambda x. \lambda s. k(gx))s_0)k \\ &= (\lambda k. f'(\lambda x. \lambda s. k(gx))s_0)(\lambda a. \lambda k'. ka)k \\ &= f'(\lambda x. \lambda s. \lambda k'. k(gx))s_0k \\ &= f'(\lambda x. \lambda s. \lambda k'. (\lambda k''. k(gx))(\lambda x. k'(h(s, x))))s_0k \\ &= f'(\lambda x. \lambda s. \text{map}_{\text{ContT } c \ m}(\lambda x. h(s, x))(\lambda k''. k(gx)))s_0k \\ &= \text{map}_{\text{ContT } c \ m} g(f(\lambda x. \lambda k''. k(gx)))k \quad (\text{pre-condition}) \end{aligned}$$

$$\begin{aligned} \text{map}_{\text{ContT } c \ m} g(\text{callcc } f)k &= \text{callcc } f(\lambda x. k(gx)) \\ &= (\lambda k. f(\lambda a. \lambda k'. ka)k)(\lambda x. k(gx)) \\ &= f(\lambda a. \lambda k'. k(ga))(\lambda x. k(gx)) \\ &= \text{map}_{\text{ContT } c \ m} g(f(\lambda a. \lambda k'. k(ga)))k \end{aligned}$$

Case “ $t = \text{EnvT } r$ ”:

Let:

$$\begin{aligned} \underline{f'} k s_0 &= f'(\lambda x. \lambda s. \lambda \rho'. kxs)s_0\rho \\ \underline{f} k &= f(\lambda a. \lambda \rho'. ka)\rho \end{aligned}$$

We first verify that:

$$\begin{aligned}
& \underline{f}'(\lambda x. \lambda s. \text{map}_m(\lambda x. h(s, x))(kx))s_0 \\
&= f'(\lambda x. \lambda s. \lambda \rho'. \text{map}_m(\lambda x. h(s, x))(kx))s_0\rho \\
&= f'(\lambda x. \lambda s. \lambda \rho'. \text{bind}_m(kx)(\lambda x. \text{return}_m(h(s, x))))s_0\rho \\
&= f'(\lambda x. \lambda s. \lambda \rho'. \text{bind}_m((\lambda \rho''. kx)\rho') \\
&\quad (\lambda x. (\lambda \rho''. \text{return}_m(h(s, x)))\rho'))s_0\rho \\
&= f'(\lambda x. \lambda s. \text{bind}_{tm}(\lambda \rho''. kx)(\lambda x. \text{return}_{tm}(h(s, x))))s_0\rho \\
&= f'(\lambda x. \lambda s. \text{map}_{tm}(\lambda x. h(s, x))(\lambda \rho''. kx))s_0\rho \\
&\quad (\text{condition}) \\
&= \text{map}_{tm} g(f(\lambda x. \lambda \rho''. kx))\rho \\
&= \text{bind}_{tm}(f(\lambda x. \lambda \rho''. kx))(\lambda x. \text{return}_{tm}(gx))\rho \\
&= \text{bind}_m(f(\lambda x. \lambda \rho''. kx)\rho)(\lambda x. \text{return}_m(gx)) \\
&= \text{bind}_m(\underline{fk})(\lambda x. \text{return}_m(gx)) \\
&= \text{map}_m g(\underline{fk})
\end{aligned}$$

We now set out to prove:

$$\text{callcc}_{tm}(\lambda k. \underline{f}'(\lambda x. \lambda s. k(gx))s_0)\rho = \text{map}_{tm} g(\text{callcc}_{tm} f)\rho$$

$$\begin{aligned}
& \text{callcc}_{tm}(\lambda k. \underline{f}'(\lambda x. \lambda s. k(gx))s_0)\rho \\
&= \text{callcc}_m(\lambda k. (\lambda k. \underline{f}'(\lambda x. \lambda s. k(gx))s_0)(\lambda a. \lambda \rho'. ka)\rho) \\
&= \text{callcc}_m(\lambda k. \underline{f}'(\lambda x. \lambda s. \lambda \rho'. k(gx))s_0)\rho
\end{aligned}$$

$$\begin{aligned}
& \text{map}_{tm} g(\text{callcc}_{tm} f)\rho \\
&= \text{bind}_{tm}(\text{callcc}_{tm} f)(\lambda x. \text{return}_{tm}(gx))\rho \\
&= \text{bind}_m(\text{callcc}_{tm} f \rho)(\lambda x. \text{return}_m(gx)) \\
&= \text{map}_m g(\text{callcc}_m(\lambda k. f(\lambda a. \lambda \rho'. ka)\rho)) \\
&= \text{map}_m g(\text{callcc}_m \underline{f}) \\
&\quad (\text{induction hypo.}) \\
&= \text{callcc}_m(\lambda k. \underline{f}'(\lambda x. \lambda s. k(gx))s_0) \\
&= \text{callcc}_m(\lambda k. \underline{f}'(\lambda x. \lambda s. \lambda \rho'. k(gx))s_0)\rho
\end{aligned}$$

Case " $t = \text{StateT } s$:" (States of type s are underlined.)

Let:

$$\underline{f}' k s_0 = f'(\lambda x. \lambda s. \lambda \underline{s_1}. k(\underline{s_0}, x) s) s_0 s_0$$

$$\underline{f} k = f(\lambda a. \lambda \underline{s_1}. k(\underline{s_0}, a)) s_0$$

$$\underline{g} = \lambda(\underline{s_1}, x). (\underline{s_1}, gx)$$

$$\underline{h} = \lambda(s, (\underline{s_2}, x)). h(\underline{s_2}, (s, x))$$

We first verify that:

$$\begin{aligned} & \underline{f}'(\lambda x. \lambda s. \text{map}_m(\lambda x. \underline{h}(s, x))(kx)) s_0 \\ &= f'(\lambda x. \lambda s. \lambda \underline{s_1}. \text{map}_m(\lambda x. \underline{h}(s, x))(k(\underline{s_0}, x))) s_0 s_0 \\ &= f'(\lambda x. \lambda s. \lambda \underline{s_1}. \text{bind}_m(k(\underline{s_0}, x))(\lambda x. \text{return}_m(\underline{h}(s, x)))) s_0 s_0 \\ &= f'(\lambda x. \lambda s. \lambda \underline{s_1}. \text{bind}_m((\lambda \underline{s_2}. k(\underline{s_0}, x)) \underline{s_1}) \\ & \quad (\lambda(\underline{s_2}, x). (\lambda \underline{s_3}. \text{return}_m(\underline{h}(s, (\underline{s_2}, x)))) \underline{s_2})) s_0 s_0 \\ &= f'(\lambda x. \lambda s. \text{bind}_{tm}(\lambda \underline{s_2}. k(\underline{s_0}, x))(\lambda x. \text{return}_{tm}(h(s, x)))) s_0 s_0 \\ &= f'(\lambda x. \lambda s. \text{map}_{tm}(\lambda x. h(s, x))(\lambda \underline{s_2}. k(\underline{s_0}, x))) s_0 s_0 \\ & \quad (\text{condition}) \\ &= \text{map}_{tm} g (f(\lambda x. \lambda \underline{s_2}. k(\underline{s_0}, x))) s_0 \\ &= \text{bind}_{tm} (f(\lambda x. \lambda \underline{s_2}. k(\underline{s_2}, x))) (\lambda x. \text{return}_{tm}(gx)) s_0 \\ &= \text{bind}_m (f(\lambda x. \lambda \underline{s_2}. k(\underline{s_0}, x)) s_0) (\lambda(\underline{s_2}, x). \text{return}_m(\underline{s_2}, gx)) \\ &= \text{bind}_m (\underline{f}k) (\lambda(\underline{s_2}, x). \text{return}_m(\underline{s_2}, gx)) \\ &= \text{map}_m \underline{g} (\underline{f}k) \end{aligned}$$

We now set out to prove:

$$\text{callcc}_{tm}(\lambda k. f'(\lambda x. \lambda s. k(gx)) s_0) s_0 = \text{map}_{tm} g (\text{callcc}_{tm} f) s_0$$

$$\begin{aligned} & \text{callcc}_{tm}(\lambda k. f'(\lambda x. \lambda s. k(gx)) s_0) s_0 \\ &= \text{callcc}_m(\lambda k. (\lambda k. f'(\lambda x. \lambda s. k(gx)) s_0) (\lambda a. \lambda \underline{s_1}. k(\underline{s_0}, a)) s_0) \\ &= \text{callcc}_m(\lambda k. f'(\lambda x. \lambda s. \lambda \underline{s_1}. k(\underline{s_0}, gx)) s_0) s_0 \end{aligned}$$

$$\begin{aligned}
& \text{map}_{tm} \, g \, (\text{callcc}_{tm} \, f) \, \underline{s_0} \\
&= \text{bind}_{tm} (\text{callcc}_{tm} \, f) (\lambda x. \text{return}_{tm} (gx)) \, \underline{s_0} \\
&= \text{bind}_m (\text{callcc}_{tm} \, f \, \underline{s_0}) (\lambda (\underline{s_1}, x). \text{return}_m (\underline{s_1}, gx)) \\
&= \text{map}_m (\lambda (\underline{s_1}, x). (\underline{s_1}, gx)) (\text{callcc}_m (\lambda k. f (\lambda a. \lambda \underline{s_1}. k (\underline{s_0}, a)) \, \underline{s_0})) \\
&= \text{map}_m \, \underline{g} \, (\text{callcc}_m \, f) \\
&\quad (\text{induction hypo.}) \\
&= \text{callcc}_m (\lambda k. \underline{f}' (\lambda x. \lambda s. k (gx)) \, s_0) \\
&= \text{callcc}_m (\lambda k. \underline{f}' (\lambda x. \lambda s. \lambda \underline{s_1}. k (\underline{g} (\underline{s_0}, x))) \, s_0 \, \underline{s_0}) \\
&= \text{callcc}_m (\lambda k. \underline{f}' (\lambda x. \lambda s. \lambda \underline{s_1}. k (\underline{s_0}, gx)) \, s_0 \, \underline{s_0})
\end{aligned}$$

Case " $t = \text{ErrT}$:"

Let:

$$\begin{aligned}
\underline{f}' \, k \, s_0 &= \underline{f}' (\lambda x. \lambda s. k (\text{Ok } x) \, s) \, s_0 \\
\underline{f} \, k &= \underline{f} (\lambda a. k (\text{Ok } a)) \\
\underline{g} &= \lambda a. \text{case } a \text{ of } \text{Ok } x \rightarrow \text{Ok } (gx) \\
&\quad \text{Err } s \rightarrow \text{Err } s \\
\underline{h} &= \lambda (s, a). \text{case } a \text{ of } \text{Ok } x \rightarrow \text{Ok } (h(s, x)) \\
&\quad \text{Err } s \rightarrow \text{Err } s
\end{aligned}$$

We first verify that:

$$\begin{aligned}
& \underline{f}'(\lambda x. \lambda s. \text{map}_m(\lambda x. \underline{h}(s, x))(kx))s_0 \\
&= f'(\lambda x. \lambda s. \text{map}_m(\lambda x. \underline{h}(s, x))(k(Ok\ x)))s_0 \\
&= f'(\lambda x. \lambda s. \text{bind}_m(k(Ok\ x))(\lambda x. \text{return}_m(\underline{h}(s, x))))s_0 \\
&= f'(\lambda x. \lambda s. \text{bind}_m(k(Ok\ x)) \\
&\quad (\lambda x. \text{return}_m(\text{case } x \text{ of } Ok\ y \rightarrow Ok\ (h(s, y)) \\
&\quad\quad\quad Err\ s \rightarrow Err\ s)))s_0 \\
&= f'(\lambda x. \lambda s. \text{bind}_m(k(Ok\ x)) \\
&\quad (\lambda x. \text{case } a \text{ of } Ok\ y \rightarrow \text{return}_m(Ok\ (h(s, y))) \\
&\quad\quad\quad Err\ s \rightarrow \text{return}_m(Err\ s)))s_0 \\
&= f'(\lambda x. \lambda s. \text{bind}_{tm}(k(Ok\ x))(\lambda x. \text{return}_{tm}(h(s, x))))s_0 \\
&= f'(\lambda x. \lambda s. \text{map}_{tm}(\lambda x. h(s, x))(k(Ok\ x)))s_0 \\
&\quad (\text{condition}) \\
&= \text{map}_{tm}\ g(f(\lambda x. k(Ok\ x))) \\
&= \text{bind}_{tm}(f(\lambda x. k(Ok\ x)))(\lambda x. \text{return}_{tm}(gx)) \\
&= \text{bind}_m(\underline{fk})(\lambda a. \text{case } a \text{ of } Ok\ x \rightarrow \text{return}_m(Ok\ (gx)) \\
&\quad\quad\quad Err\ s \rightarrow \text{return}_m(Err\ s)) \\
&= \text{map}_m\ \underline{g}\ (\underline{fk})
\end{aligned}$$

We now set out to prove:

$$\text{callcc}_{tm}(\lambda k. f'(\lambda x. \lambda s. k(gx))s_0) = \text{map}_{tm}\ g(\text{callcc}_{tm}\ f)$$

$$\begin{aligned}
& \text{callcc}_{tm}(\lambda k. f'(\lambda x. \lambda s. k(gx))s_0) \\
&= \text{callcc}_m(\lambda k. (\lambda k. f'(\lambda x. \lambda s. k(gx))s_0)(\lambda a. k(Ok\ a))) \\
&= \text{callcc}_m(\lambda k. f'(\lambda x. \lambda s. k(Ok\ (gx))))s_0
\end{aligned}$$

$$\begin{aligned}
& \text{map}_{tm} g (\text{callcc}_{tm} f) \\
&= \text{bind}_{tm} (\text{callcc}_{tm} f) (\lambda x. \text{return}_{tm} (gx)) \\
&= \text{bind}_m (\text{callcc}_{tm} f) (\lambda a. \text{case } a \text{ of } Ok\ x \rightarrow \text{return}_m (Ok\ (gx)) \\
&\quad\quad\quad Err\ s \rightarrow \text{return}_m (Err\ s)) \\
&= \text{map}_m \underline{g} (\text{callcc}_m \underline{f}) \\
&\quad\quad\quad (\text{induction hypo.}) \\
&= \text{callcc}_m (\lambda k. \underline{f}' (\lambda x. \lambda s. k (\underline{g} x)) s_0) \\
&= \text{callcc}_m (\lambda k. \underline{f}' (\lambda x. \lambda s. k (\underline{g} (Ok\ x))) s_0) \\
&= \text{callcc}_m (\lambda k. \underline{f}' (\lambda x. \lambda s. k (Ok\ (gx))) s_0)
\end{aligned}$$

Theorem 5.2.1 For any source language program e , we have:

$$\text{inEnv } \bar{\rho} \ E[[e]] = \text{inEnv } \bar{\rho} \ N[[e]]$$

Proof:

We prove the theorem by induction over the structure of expressions.

Arithmetic expressions:

$$\begin{aligned}
\text{inEnv } \bar{\rho} \ E[[n]] &= \text{inEnv } \bar{\rho} \ (\text{return } n) \quad (E) \\
&= \text{inEnv } \bar{\rho} \ N[[n]] \quad (N)
\end{aligned}$$

$$\begin{aligned}
& \text{inEnv } \bar{\rho} \ E[[e_1 + e_2]] \\
&= \text{inEnv } \bar{\rho} \ \{v_1 \leftarrow E[[e_1]]; v_2 \leftarrow E[[e_2]]; \text{return } (v_1 + v_2)\} \quad (E) \\
&= \{v_1 \leftarrow \text{inEnv } \bar{\rho} \ E[[e_1]]; v_2 \leftarrow \text{inEnv } \bar{\rho} \ E[[e_2]]; \\
&\quad\quad\quad \text{inEnv } \bar{\rho} \ \text{return } (v_1 + v_2)\} \quad (\text{distribution}) \\
&= \{v_1 \leftarrow \text{inEnv } \bar{\rho} \ N[[e_1]]; v_2 \leftarrow \text{inEnv } \bar{\rho} \ N[[e_2]]; \\
&\quad\quad\quad \text{inEnv } \bar{\rho} \ \text{return } (v_1 + v_2)\} \quad (\text{ind. hypo.}) \\
&= \text{inEnv } \bar{\rho} \ \{v_1 \leftarrow N[[e_1]]; v_2 \leftarrow N[[e_2]]; \text{return } (v_1 + v_2)\} \quad (\text{distribution}) \\
&= \text{inEnv } \bar{\rho} \ N[[e_1 + e_2]] \quad (N)
\end{aligned}$$

Functions:

$$\begin{aligned}
inEnv \bar{\rho} E[v] &= inEnv \bar{\rho} \{ \rho \leftarrow rdEnv; \rho[v] \} && (E) \\
&= \{ \rho \leftarrow inEnv \bar{\rho} rdEnv; inEnv \bar{\rho} (\rho[v]) \} && (\text{distribution}) \\
&= \{ \rho \leftarrow return \bar{\rho}; inEnv \bar{\rho} (\rho[v]) \} && (\text{cancellation}) \\
&= inEnv \bar{\rho} (\bar{\rho}[v]) && (\text{left unit}) \\
&= inEnv \bar{\rho} \bar{v} && (\bar{\rho}) \\
&= inEnv \bar{\rho} N[v] && (N)
\end{aligned}$$

$$\begin{aligned}
inEnv \bar{\rho} E[\lambda v. e] &= inEnv \bar{\rho} \{ \rho \leftarrow rdEnv; return(\lambda c. inEnv \rho [c/[v]] E[e]) \} && (E) \\
&= \{ \rho \leftarrow inEnv \bar{\rho} rdEnv; \\
&\quad inEnv \bar{\rho} (return(\lambda c. inEnv \rho [c/[v]] E[e])) \} && (\text{distribution}) \\
&= \{ \rho \leftarrow return \bar{\rho}; return(\lambda c. inEnv \rho [c/[v]] E[e]) \} && (\text{cancel., unit}) \\
&= return(\lambda c. inEnv \bar{\rho} [c/[v]] E[e]) && (\text{left unit}) \\
&= return(\lambda \bar{v}. inEnv \bar{\rho} [\bar{v}/[v]] E[e]) && (\alpha \text{ renaming}) \\
&= return(\lambda \bar{v}. inEnv \bar{\rho} E[e]) && (\bar{\rho}) \\
&= return(\lambda \bar{v}. inEnv \bar{\rho} N[e]) && (\text{ind. hypo.}) \\
&= inEnv \bar{\rho} N[\lambda v. e] && (N)
\end{aligned}$$

$$\begin{aligned}
inEnv \bar{\rho} E[(e_1 e_2)_n] &= inEnv \bar{\rho} \{ f \leftarrow E[e_1]; \rho \leftarrow rdEnv; f(inEnv \rho E[e_2]) \} && (E) \\
&= \{ f \leftarrow inEnv \bar{\rho} E[e_1]; \rho \leftarrow inEnv \bar{\rho} rdEnv; \\
&\quad inEnv \bar{\rho} (f(inEnv \rho E[e_2])) \} && (\text{distribution}) \\
&= \{ f \leftarrow inEnv \bar{\rho} E[e_1]; \rho \leftarrow return \bar{\rho}; \\
&\quad inEnv \bar{\rho} (f(inEnv \rho E[e_2])) \} && (\text{cancellation}) \\
&= \{ f \leftarrow inEnv \bar{\rho} E[e_1]; inEnv \bar{\rho} (f(inEnv \bar{\rho} E[e_2])) \} && (\text{left unit}) \\
&= \{ f \leftarrow inEnv \bar{\rho} N[e_1]; inEnv \bar{\rho} (f(inEnv \bar{\rho} N[e_2])) \} && (\text{ind. hypo.}) \\
&= inEnv \bar{\rho} \{ f \leftarrow N[e_1]; f(inEnv \bar{\rho} N[e_2]) \} && (\text{distribution}) \\
&= inEnv \bar{\rho} N[(e_1 e_2)_n] && (N)
\end{aligned}$$

$$\begin{aligned}
& \text{inEnv } \bar{\rho} E[(e_1 e_2)_v] \\
&= \text{inEnv } \bar{\rho} \{f \leftarrow E[e_1]; v \leftarrow E[e_2]; f(\text{return } v)\} \quad (E) \\
&= \{f \leftarrow \text{inEnv } \bar{\rho} E[e_1]; v \leftarrow \text{inEnv } \bar{\rho} E[e_2]; \\
&\quad \text{inEnv } \bar{\rho} (f(\text{return } v))\} \quad (\text{distribution}) \\
&= \{f \leftarrow \text{inEnv } \bar{\rho} N[e_1]; v \leftarrow \text{inEnv } \bar{\rho} N[e_2]; \\
&\quad \text{inEnv } \bar{\rho} (f(\text{return } v))\} \quad (\text{ind. hypo.}) \\
&= \text{inEnv } \bar{\rho} \{f \leftarrow N[e_1]; v \leftarrow N[e_2]; f(\text{return } v)\} \quad (\text{distribution}) \\
&= \text{inEnv } \bar{\rho} N[(e_1 e_2)_v] \quad (N)
\end{aligned}$$

References and assignment:

We can prove:

$$\begin{aligned}
\text{inEnv } \bar{\rho} E[\mathbf{ref } e] &= \text{inEnv } \bar{\rho} N[\mathbf{ref } e] \\
\text{inEnv } \bar{\rho} E[\mathbf{deref } e] &= \text{inEnv } \bar{\rho} N[\mathbf{deref } e] \\
\text{inEnv } \bar{\rho} E[e_1 := e_2] &= \text{inEnv } \bar{\rho} N[e_1 := e_2]
\end{aligned}$$

the same way we established the case for $[(e_1 + e_2)]$.

Lazy evaluation:

$$\begin{aligned}
& \text{inEnv } \bar{\rho} E[(e_1 e_2)_l] \\
&= \text{inEnv } \bar{\rho} \{f \leftarrow E[e_1]; l \leftarrow \text{alloc}; \rho \leftarrow \text{rdEnv}; \\
&\quad \mathbf{let } \mathbf{thunk} = \{v \leftarrow \text{inEnv } \rho E[e_2]; \dots\} \mathbf{in } \dots\} \quad (E) \\
&= \{f \leftarrow \text{inEnv } \bar{\rho} E[e_1]; l \leftarrow \text{inEnv } \bar{\rho} \text{alloc}; \\
&\quad \rho \leftarrow \text{inEnv } \bar{\rho} \text{rdEnv}; \\
&\quad \text{inEnv } \bar{\rho} (\mathbf{let } \mathbf{thunk} = \{v \leftarrow \text{inEnv } \rho E[e_2]; \dots\} \mathbf{in } \dots)\} \quad (\text{distribution}) \\
&= \{f \leftarrow \text{inEnv } \bar{\rho} E[e_1]; l \leftarrow \text{inEnv } \bar{\rho} \text{alloc}; \\
&\quad \text{inEnv } \bar{\rho} (\mathbf{let } \mathbf{thunk} = \{v \leftarrow \text{inEnv } \bar{\rho} E[e_2]; \dots\} \mathbf{in } \dots)\} \quad (\text{can., l. unit}) \\
&= \{f \leftarrow \text{inEnv } \bar{\rho} N[e_1]; l \leftarrow \text{inEnv } \bar{\rho} \text{alloc}; \\
&\quad \text{inEnv } \bar{\rho} (\mathbf{let } \mathbf{thunk} = \{v \leftarrow \text{inEnv } \bar{\rho} N[e_2]; \dots\} \mathbf{in } \dots)\} \quad (\text{ind. hypo.}) \\
&= \text{inEnv } \bar{\rho} \{f \leftarrow N[e_1]; l \leftarrow \text{alloc}; \\
&\quad \mathbf{let } \mathbf{thunk} = \{v \leftarrow \text{inEnv } \bar{\rho} N[e_2]; \dots\} \mathbf{in } \dots\} \quad (\text{distribution}) \\
&= \text{inEnv } \bar{\rho} N[(e_1 e_2)_l] \quad (N)
\end{aligned}$$

Tracing:

Again, we can prove:

$$\text{inEnv } \bar{\rho} E[l @ e] = \text{inEnv } \bar{\rho} N[l @ e]$$

the same way we established the case for $\llbracket e_1 + e_2 \rrbracket$.

First-class continuations:

We can prove:

$$\text{inEnv } \bar{\rho} E[\text{callcc}] = \text{inEnv } \bar{\rho} N[\text{callcc}]$$

the same way we established the case for $\llbracket n \rrbracket$.

Nondeterminism:

First we establish a lemma:

$$\begin{aligned}
& \text{inEnv } \rho (\text{merge } (\text{map } (\lambda x. \text{inEnv } \rho x) e)) \\
&= \text{inEnv } \rho (\text{join } (\text{lift } (\text{map } (\lambda x. \text{inEnv } \rho x) e))) && (\text{merge}) \\
&= \text{inEnv } \rho (\text{join } (\text{lift } \{x \leftarrow e; \text{return } (\text{inEnv } \rho x)\})) && (\text{map}) \\
&= \text{inEnv } \rho (\text{join } \{x \leftarrow \text{lift } e; \text{lift } (\text{return } (\text{inEnv } \rho x))\}) && (\text{monad morphism}) \\
&= \text{inEnv } \rho (\text{join } \{x \leftarrow \text{lift } e; \text{return } (\text{inEnv } \rho x)\}) && (\text{monad morphism}) \\
&= \text{inEnv } \rho \{x \leftarrow \text{lift } e; a \leftarrow \text{return } (\text{inEnv } \rho x); a\} && (\text{join}) \\
&= \text{inEnv } \rho \{x \leftarrow \text{lift } e; \text{inEnv } \rho x\} && (\text{left unit}) \\
&= \{x \leftarrow \text{inEnv } \rho (\text{lift } e); \text{inEnv } \rho (\text{inEnv } \rho x)\} && (\text{distribution}) \\
&= \{x \leftarrow \text{inEnv } \rho (\text{lift } e); \text{inEnv } \rho x\} && (\text{overriding}) \\
&= \text{inEnv } \rho \{x \leftarrow \text{lift } e; x\} && (\text{distribution}) \\
&= \text{inEnv } \rho (\text{join } (\text{lift } e)) && (\text{join}) \\
&= \text{inEnv } \rho (\text{merge } e) && (\text{merge})
\end{aligned}$$

Now we can prove:

$$\begin{aligned}
& \text{inEnv } \bar{\rho} E[\{e_0, e_1, \dots\}] \\
&= \text{inEnv } \bar{\rho} \text{ merge } [E[e_0], E[e_1], \dots] && (E) \\
&= \text{inEnv } \bar{\rho} \text{ merge } [\text{inEnv } \bar{\rho} E[e_0], \text{inEnv } \bar{\rho} E[e_1], \dots] && (\text{lemma}) \\
&= \text{inEnv } \bar{\rho} \text{ merge } [\text{inEnv } \bar{\rho} N[e_0], \text{inEnv } \bar{\rho} N[e_1], \dots] && (\text{ind. hypo.}) \\
&= \text{inEnv } \bar{\rho} \text{ merge } [N[e_0], N[e_1], \dots] && (\text{lemma}) \\
&= \text{inEnv } \bar{\rho} N[\{e_0, e_1, \dots\}] && (N)
\end{aligned}$$

Appendix B

Gofer Code for Monad Transformers

This section lists the Gofer implementation for three monad transformers (environment, continuation and error reporting) and their associated liftings.

```
-- Environment monad transformer -----  
  
data EnvT r m a = EnvM (r -> m a)  
unEnvM (EnvM x) = x  
  
instance Monad m => Monad (EnvT r m) where  
  return a          = EnvM (\r -> return a)  
  (EnvM m) 'bind' k = EnvM (\r -> m r 'bind' \a ->  
                           unEnvM (k a) r)  
  
instance MonadT (EnvT r) where  
  -- lift :: m a -> EnvT r m a  
  lift m = EnvM (\r -> m)  
  
class Monad m => EnvMonad r m where  
  inEnv  :: r -> m a -> m a  
  rdEnv  :: m r  
  
instance Monad m => EnvMonad r (EnvT r m) where  
  inEnv r (EnvM m) = EnvM (\_ -> m r)  
  rdEnv           = EnvM (\r -> return r)
```

```

-- lift EnvMonad through EnvT
instance (MonadT (EnvT r'), EnvMonad r m) =>
    EnvMonad r (EnvT r' m) where
    inEnv r (EnvM m) = EnvM (\r' -> inEnv r (m r'))
    rdEnv           = lift rdEnv

-- lift EnvMonad through StateT
instance (MonadT (StateT s), EnvMonad r m) =>
    EnvMonad r (StateT s m) where
    inEnv r (StateM m) = StateM (\s -> inEnv r (m s))
    rdEnv              = lift rdEnv

-- lift EnvMonad through ErrT
instance (MonadT ErrT, EnvMonad r m) =>
    EnvMonad r (ErrT m) where
    inEnv r (ErrM m) = ErrM (inEnv r m)
    rdEnv           = lift rdEnv

-- Error monad transformer -----
data Err a = Ok a | Err String
data ErrT m a = ErrM (m (Err a))
unErrM (ErrM x) = x

instance Monad m => Monad (ErrT m) where
    return          = ErrM . return . Ok
    (ErrM m) `bind` k = ErrM (m `bind` \a ->
        case a of
            Ok x      -> unErrM (k x)
            Err msg   -> return (Err msg))

instance MonadT ErrT where
    -- lift :: m a -> ErrT m a
    lift c = ErrM (map Ok c)

class Monad m => ErrMonad m where
    err :: String -> m a

instance Monad m => ErrMonad (ErrT m) where
    err = ErrM . return . Err

```

```

instance (ErrMonad m, MonadT t) => ErrMonad (t m) where
    err = lift . err

-- Continuation monad transformer -----

data ContT ans m a = ContM ((a -> m ans) -> m ans)
unContM (ContM x) = x

instance Monad m => Monad (ContT ans m) where
    return x          = ContM (\k -> k x)
    (ContM m) `bind` f =
        ContM (\k -> m (\a -> unContM (f a) k))

instance MonadT (ContT ans) where
    -- lift :: m a -> ContT ans m a
    lift m = ContM (\f -> m `bind` f)

class Monad m => ContMonad m where
    callcc :: ((a -> m b) -> m a) -> m a

instance Monad m => ContMonad (ContT ans m) where
    callcc f =
        ContM (\k -> unContM (f (\a -> ContM (\_ -> k a))) k)

-- lift callcc through EnvT
instance (MonadT (EnvT r), ContMonad m) =>
    ContMonad (EnvT r m) where
    callcc f = EnvM (\r -> callcc (\k ->
        unEnvM (f (\a -> EnvM (\r -> k a))) r))

-- lift callcc through StateT
instance (MonadT (StateT s), ContMonad m) =>
    ContMonad (StateT s m) where
    callcc f = StateM (\s -> callcc (\k -> unStateM
        (f (\a -> StateM (\s1 -> k (s, a)))) s))

-- lift callcc through ErrT
instance (MonadT ErrT, ContMonad m) =>
    ContMonad (ErrT m) where

```

```
callcc f = ErrM (callcc (\k ->
    unErrM (f (\a -> ErrM (k (Ok a)))))
```

Appendix C

Monadic Semantics of the STG Language

We split the overall semantics for the STG language into a number of semantic functions. There are separate functions for literals, variables, atoms, function applications, constructors, primitives, case expressions, let bindings, and right-hand sides.

We use the following notations introduced in Chapter 5.

- The strictness environment ρ .
- Meta-language variables \bar{v} and $\bar{\bar{v}}$ corresponding to the optimized and standard entries of functions.
- Conversion functions *delay* and *force*.
- Helper functions *optApp*, *stdApp*, and *stdEntry*.

C.1 Literals, Variables and Atoms

Literals map to their meta-language counterparts:

$$E_e[[l]] = \text{return } \bar{l}$$

The strictness environment keeps track of whether a given variable is a known function. Although variables can be either a value or a cell, known functions are always a value.

$$E_v[[v]] = \{\rho \leftarrow rdEnv;$$

$$\quad \text{case } \rho[[v]] \text{ of}$$

$$\quad \text{FUN } _ \rightarrow \text{return } \bar{v}$$

$$\quad \text{VAR } m \rightarrow \text{force? } m \bar{v}\}$$

Note that \bar{v} is the standard entry point of functions. Function *force?* is a variation of *force*:

$$\text{force? } V v = \text{return } v$$

$$\text{force? } C c = \text{force } c$$

Atoms are either literals or variables:

$$E_a[[v]] = E_v[[v]]$$

$$E_a[[l]] = E_l[[l]]$$

C.2 Function Applications

In the STG language, a function application with no arguments denotes an expression made up of a single variable. Normal function applications take at least one argument.

If the function is known, and there are enough arguments available, the optimized entry can be used; otherwise, the standard entry is used.

$$\begin{aligned}
E_e[v \{\}] &= E_v[v] \\
E_e[v \{a_1, \dots, a_n\}] &= \\
&\quad \{\rho \leftarrow rdEnv; \\
&\quad \text{case } \rho[v] \text{ of} \\
&\quad \quad \text{FUN } [m_1, \dots, m_k] \rightarrow \\
&\quad \quad \quad \text{if } k > n \text{ then} \\
&\quad \quad \quad \quad \text{stdApp } (E_v[v]) \\
&\quad \quad \quad \quad \quad [delay E_a[a_1], \dots, delay E_a[a_n]] \\
&\quad \quad \quad \quad \text{else if } k = n \text{ then} \\
&\quad \quad \quad \quad \quad \text{optApp } \bar{v} [delay? m_1 E_a[a_1], \dots, delay? m_k E_a[a_n]] \\
&\quad \quad \quad \quad \text{else} \\
&\quad \quad \quad \quad \quad \text{stdApp } (\text{optApp } \bar{v} [delay? m_1 E_a[a_1], \dots, delay? m_k E_a[a_k]]) \\
&\quad \quad \quad \quad \quad \quad [delay E_a[a_{k+1}], \dots, delay E_a[a_n]] \\
&\quad \quad \text{VAR } - \rightarrow \text{stdApp } (E_v[v]) \\
&\quad \quad \quad [delay E_a[a_1], \dots, delay E_a[a_n]] \}
\end{aligned}$$

C.3 Constructors and Primitives

Constructor and primitive applications are always saturated. Therefore we can take advantage of the strictness information.

$$\begin{aligned}
E_e[c \{a_1, \dots, a_n\}] &= \text{let} \\
&\quad [m_1, \dots, m_n] = \text{strictnessOf } c \\
&\quad \text{in} \\
&\quad \{x_1 \leftarrow \text{delay? } m_1 E_a[a_1]; \dots \\
&\quad \quad x_n \leftarrow \text{delay? } m_n E_a[a_n]; \\
&\quad \quad \bar{c} [x_1, \dots, x_n]\}
\end{aligned}$$

$$\begin{aligned}
E_e \llbracket p \{a_1, \dots, a_n\} \rrbracket &= \mathbf{let} \\
&\quad [m_1, \dots, m_n] = \mathit{strictnessOf} \, p \\
&\mathbf{in} \\
&\quad \{x_1 \leftarrow \mathit{delay?} \, m_1 \, E_a \llbracket a_1 \rrbracket; \dots \\
&\quad \quad x_n \leftarrow \mathit{delay?} \, m_n \, E_a \llbracket a_n \rrbracket; \\
&\quad \quad \bar{p} \, [x_1, \dots, x_n]\}
\end{aligned}$$

C.4 Case Expressions

The meta-language case construct is used to interpret STG case expressions. In the default case, the strictness environment is augmented with the pattern variable.

$$\begin{aligned}
E_e \llbracket \mathit{case} \, e \, \mathit{of} \\
\quad l_1 \rightarrow e_1; \dots \\
\quad l_n \rightarrow e_n; \\
\quad v \rightarrow e_d \rrbracket &= \{ \rho \leftarrow \mathit{rdEnv}; \\
&\quad x \leftarrow E_e \llbracket e \rrbracket; \\
&\quad \mathbf{case} \, x \, \mathbf{of} \\
&\quad \quad \bar{l}_1 \rightarrow E_e \llbracket e_1 \rrbracket; \dots \\
&\quad \quad \bar{l}_n \rightarrow E_e \llbracket e_n \rrbracket; \\
&\quad \quad \bar{v} \rightarrow \mathit{inEnv} \, \rho \, [\mathit{VAR} \, V / \llbracket v \rrbracket] \, E_e \llbracket e_d \rrbracket \}
\end{aligned}$$

The semantics for case expression on data types is very similar. The strictness environment is augmented with all pattern variables. The variable matching the default case is always a value, because expression e has already been evaluated. On the other hand, variables matching the components of a data type can be either a value or a cell, depending on the strictness property of the data constructor.

$$\begin{aligned}
& E_e \llbracket \text{case } e \text{ of} \\
& \quad c_1 \{v_{11}, \dots\} \rightarrow e_1; \dots \\
& \quad c_n \{v_{n1}, \dots\} \rightarrow e_n; \\
& \quad v \qquad \qquad \rightarrow e_d \rrbracket = \\
& \{ \rho \leftarrow rdEnv; \\
& \quad x \leftarrow E_e \llbracket e \rrbracket; \\
& \quad \text{case } x \text{ of} \\
& \quad \quad \bar{c}_1 \{ \bar{v}_{11}, \dots, \} \rightarrow inEnv \rho [VAR (strictnessOf v_{11}) / \llbracket v_{11} \rrbracket, \dots] E_e \llbracket e_1 \rrbracket \dots \\
& \quad \quad \bar{c}_n \{ \bar{v}_{n1}, \dots, \} \rightarrow inEnv \rho [VAR (strictnessOf v_{n1}) / \llbracket v_{n1} \rrbracket, \dots] E_e \llbracket e_n \rrbracket \\
& \quad \quad \bar{v} \qquad \qquad \rightarrow inEnv \rho [VAR V / \llbracket v \rrbracket] E_e \llbracket e_d \rrbracket \}
\end{aligned}$$

C.5 Let Bindings

To simplify the presentation, we introduce a helper function that maps right-hand-sides to strictness information:

$$\begin{aligned}
info \llbracket \lambda \{ \} . e \rrbracket &= VAR C \\
info \llbracket \lambda \{ v_1, \dots, v_n \} . e \rrbracket &= FUN [strictnessOf v_1, \dots, strictnessOf v_n]
\end{aligned}$$

The semantics of let bindings is as follows. Note that we need the **fix** operator even for the non-recursive binding. The reason is that the standard entry point is defined in terms of the optimized entry point, but may also be used inside the body of function definition.

