

Part XIII

Modalities

Chapter 37

Computational Effects

In this chapter we study the use of types to segregate expressions that use computational effects from those that do not. While it is difficult to be precise about what constitutes a computational effect, a rough-and-ready rule is *any behavior that constrains the order of execution beyond that the requirements imposed by the flow of data*. For example, since the order in which input or output is performed clearly matters to the meaning of a program, these operations may be classified as effects. Similarly, mutation of data structures (as described in Chapter 35) is clearly sensitive to the order in which they are executed, and so mutation should also be classified as an effect.

The trouble with computational effects is precisely that they constrain the order of evaluation. This inhibits the use of parallelism (Chapter 44) or laziness (Chapter 42), and generally makes it harder to reason about the behavior of a piece of code. But it should ideally be possible to take advantage of these concepts when effects are *not* used, rather than always planning for the possibility that they might be used. One way to achieve this is to introduce a distinction, called a *modality*, between two modes of expression:

1. The *pure* expressions that are executed solely for their value, and that may engender no effects.
2. The *impure* expressions, or *commands*, that are executed for their value and their effect.

The mode distinction gives rise to a new form of type, called the *lax modality*, or *monad*, whose elements are unevaluated commands. These commands can be passed as pure data, or activated for use by a special form of command.

37.1 A Modality for Effects

The syntax of $\mathcal{L}\{\text{cmd}\}$ is given by the following grammar:

<i>Category</i>	<i>Item</i>	<i>Abstract</i>	<i>Concrete</i>
Type	τ	$::= \text{cmd}(\tau)$	$\tau \text{ cmd}$
Expr	e	$::= x$	x
		$ \text{cmd}(m)$	$\text{cmd}(m)$
Comm	m	$::= \text{return}(e)$	$\text{return}(e)$
		$ \text{letcmd}(e; x.m)$	$\text{let cmd}(x) \text{ be } e \text{ in } m$

The language $\mathcal{L}\{\text{cmd}\}$ distinguishes two modes of expression, the pure (effect-free) *expressions*, and the impure (effect-capable) *commands*. The modal type $\text{cmd}(\tau)$ consists of suspended commands that, when evaluated, yield a value of type τ . The expression $\text{cmd}(m)$ introduces an unevaluated command as a value of modal type. The command $\text{return}(e)$ returns the value of e as its value, without engendering any effects. The command $\text{letcmd}(e; x.m)$ activates the suspended command obtained by evaluating the expression e , then continue by evaluating the command m . This form sequences evaluation of commands so that there is no ambiguity about the order in which effects occur during evaluation.

The static semantics of $\mathcal{L}\{\text{cmd}\}$ consists of two forms of typing judgement, $e : \tau$, stating that the expression e has type τ , and $m \sim \tau$, stating that the command m only yields values of type τ . Both of these judgement forms are considered with respect to hypotheses of the form $x : \tau$, which states that a variable x has type τ . The rules defining the static semantics of $\mathcal{L}\{\text{cmd}\}$ are as follows:

$$\frac{\Gamma \vdash m \sim \tau}{\Gamma \vdash \text{cmd}(m) : \text{cmd}(\tau)} \quad (37.1a)$$

$$\frac{\Gamma \vdash e : \tau}{\Gamma \vdash \text{return}(e) \sim \tau} \quad (37.1b)$$

$$\frac{\Gamma \vdash e : \text{cmd}(\tau) \quad \Gamma, x : \tau \vdash m \sim \tau'}{\Gamma \vdash \text{letcmd}(e; x.m) \sim \tau'} \quad (37.1c)$$

The dynamic semantics of an instance of $\mathcal{L}\{\text{cmd}\}$ is specified by two transition judgements:

1. *Evaluation* of expressions, $e \mapsto e'$.
2. *Execution* of commands, $m \mapsto m'$.

The rules of expression evaluation are carried over from the effect-free setting without change. There is, however, an additional form of value, the encapsulated command:

$$\overline{\text{cmd}(m) \text{ val}} \quad (37.2)$$

Observe that $\text{cmd}(m)$ is a value regardless of the form of m . This is because the command is not executed, but only encapsulated as a form of value.

The rule of execution enforce the sequential execution of commands.

$$\frac{e \mapsto e'}{\text{return}(e) \mapsto \text{return}(e')} \quad (37.3a)$$

$$\frac{e \text{ val}}{\text{return}(e) \text{ final}} \quad (37.3b)$$

$$\frac{e \mapsto e'}{\text{letcmd}(e; x.m) \mapsto \text{letcmd}(e'; x.m)} \quad (37.3c)$$

$$\frac{m_1 \mapsto m'_1}{\text{letcmd}(\text{cmd}(m_1); x.m_2) \mapsto \text{letcmd}(\text{cmd}(m'_1); x.m_2)} \quad (37.3d)$$

$$\frac{\text{return}(e) \text{ final}}{\text{letcmd}(\text{cmd}(\text{return}(e)); x.m) \mapsto [e/x]m} \quad (37.3e)$$

Rules (37.3a) and (37.3c) specify that the expression part of a `return` or `let` command is to be evaluated before execution can proceed. Rule (37.3b) specifies that a `return` command whose argument is a value is a final state of command execution. Rule (37.3d) specifies that a `letcomp` activates an encapsulated command, and Rule (37.3e) specifies that a completed command passes its return value to the body of the `let`.

37.2 Imperative Programming

The `bind` construct imposes a sequential evaluation order on commands, according to which the encapsulated command is executed prior to execution of the body of the `bind`. This gives rise to a familiar programming idiom, called *sequential composition*, which we now derive from the lax modality.

Since there are only two constructs for forming commands, the `bind` and the `return` command, it is easy to see that a command of type τ always has the form

$$\text{let cmd}(x_1) \text{ be } e_1 \text{ in } \dots \text{let cmd}(x_n) \text{ be } e_n \text{ in return}(e),$$

where $e_1 : \tau_1 \text{ cmd}, \dots, e_n : \tau_n \text{ cmd}$, and $x_1 : \tau_1, \dots, x_n : \tau_n \vdash e : \tau$. The dynamic semantics of $\mathcal{L}\{\text{cmd}\}$ specifies that this is evaluated by evaluating the expression, e_1 , to an encapsulated command, m_1 , then executing m_1 for its value and effects, then passing this value to e_2 , and so forth, until finally the value determined by the expression e is returned.

To execute m_1 and m_2 in sequence, where m_2 may refer to the value of m_1 via a variable x_1 , we may write

$$\text{let cmd}(x_1) \text{ be cmd}(m_1) \text{ in } m_2.$$

This encapsulates, and then immediately activates, the command m_1 , binding its value to x_1 , and continuing by executing m_2 . More generally, to execute a sequence of commands in order, passing the value of each to the next, we may write

$$\text{let cmd}(x_1) \text{ be cmd}(m_1) \text{ in } \dots \text{let cmd}(x_{k-1}) \text{ be cmd}(m_{k-1}) \text{ in } m_k.$$

Notationally, this quickly gets out of hand. We therefore introduce the *do syntax*, which is reminiscent of the notation used in many imperative programming languages. The binary do construct, $\text{do}\{x \leftarrow m_1 ; m_2\}$, stands for the command

$$\text{let cmd}(x) \text{ be cmd}(m_1) \text{ in } m_2,$$

which executes the commands m_1 and m_2 in sequence, passing the value of m_1 to m_2 via the variable x . The general do construct,

$$\text{do}\{x_1 \leftarrow m_1 ; \dots ; x_k \leftarrow m_k ; \text{return}(e)\},$$

is defined by iteration of the binary do as follows:

$$\text{do}\{x_1 \leftarrow m_1 ; \dots \text{do}\{x_k \leftarrow m_k ; \text{return}(e)\} \dots\}.$$

This notation is reminiscent of that used in many well-known programming languages. The point here is that sequential composition of commands arises from the presence of the lax modality in the language. In other words conventional imperative programming languages are implicitly structured by this type, even if the connection is not made explicit.

37.3 Integrating Effects

The modal separation of expressions from commands ensures that the semantics of expression evaluation is not compromised by the possibility of

effects. One consequence of this restriction is that it is impossible to define an expression $x : \tau \text{ cmd} \vdash \text{run } x : \tau$ whose behavior is to unbundle the command bound to x , execute it, and return its value as the value of the entire expression. For if such an expression were to exist, expression evaluation would engender effects, ruining the very distinction we are trying to preserve!

The only way for a command to occur inside of an expression is for it to be encapsulated as a value of modal type. To execute such a command it is necessary to bind it to a variable using the bind construct, which is itself a form of command. This is the essential means by which effects are confined to commands, and by which expressions are ensured to remain pure. Put another way, it is impossible to define an *expression* $\text{run } e$ of type τ , where $e : \tau \text{ cmd}$, whose value is the result of running the command encapsulated in the value of e . There is, however, a *command* $\text{run } e$ defined by

$$\text{let cmd}(x) \text{ be } e \text{ in return}(x),$$

which executes the encapsulated command and returns its value.

Now consider the extension of $\mathcal{L}\{\text{cmd}\}$ with function types. Recall from Chapter 13 that a function has the form $\lambda(x:\tau. e)$, where e is a (pure) expression. In the context of $\mathcal{L}\{\text{cmd}\}$ this implies that no function may engender an effect when applied! For example, it is not possible to write a function of the form $\lambda(x:\text{unit}. \text{print "hello"})$ that, when applied, outputs the string `hello` to the screen.

This may seem like a serious limitation, but this apparent “bug” is actually an important “feature.” To see why, observe that the type of the foregoing function would, in the absence of the lax modality, be something like $\text{unit} \rightarrow \text{unit}$. Intuitively, a function of this type is either the identity function, the constant function returning the null tuple (this is, in fact, the identity function), or a function that diverges or incurs an error when applied (in the presence of such possibilities). But, above all, it *cannot* be the function that prints `hello`.

However, let us consider the closely related type $\text{unit} \rightarrow (\text{unit cmd})$. This is the type of functions that, when applied, yield an *encapsulated command*, of type unit . One such function is

$$\lambda(x:\text{unit}. \text{cmd}(\text{print "hello"})).$$

This function *does not* output to the screen when applied, since no pure function can have an effect, but it *does* yield a command that, when executed, performs this output. Thus, if e is the above function, then the

command

$$\text{let cmd}(_) \text{ be } e(\langle \rangle) \text{ in return}(\langle \rangle) \quad (37.4)$$

executes the encapsulated command yielded by e when applied, engendering the intended effect, and returning the trivial element of unit type.

The importance of this example lies in the distinction between the type $\text{unit} \rightarrow \text{unit}$, which can only contain uninteresting functions such as the identity, and the type $\text{unit} \rightarrow (\text{unit cmd})$, which reveals in its type that the result of applying it is an encapsulated command that may, when executed, engender an arbitrary effect. In short, the type reveals the reliance on effects. The function type retains its meaning, and, in combination with the lax modality, provides a type of *procedures* that yield a command when applied. A *procedure call* is implemented by combining function application with the modal bind operation in the manner illustrated by expression (37.4).

37.4 Exercises

Chapter 38

Monadic Exceptions

As we saw in Chapter 37, if an expression can raise an exception, then the order of evaluation of sub-expressions of an expression is significant. For example, the expression $e_1 + e_2$ is not in general equivalent to $e_2 + e_1$, even though addition is commutative. This is so because in the presence of exceptions an expression of type `nat` need not evaluate to a number—it can, instead, raise an exception. If e_1 is `raise(L)` and e_2 is `raise(R)`, then we may use a handler to distinguish the two addition expressions from each other, yielding, say, zero in the one case and one in the other.

The semantics of expressions may be preserved even in the presence of exceptions if we confine them to the monad by making the primitives for raising and handling exceptions commands, rather than expressions. In this chapter we study a variation on $\mathcal{L}\{\text{cmd}\}$ in which exceptions are treated as an impurity to be confined to commands. It should be noted, however, that this approach is unsatisfactory for two related reasons. First, because the monad imposes a strict sequential execution order on commands, the programmer must specify an evaluation order whenever an exception might be raised. Second, if any exception can appear *somewhere* in a program, then it must be structured as though an exception could appear *anywhere*. This is because there is no means of “escaping the monad”—an impurity somewhere infects all parts of the program that depend on its result.

38.1 Monadic Exceptions

The most natural formalization of exceptions in the monadic framework is to regard an exception as an alternative outcome of evaluation of a com-

mand. That is, a command, when executed, may engender effects, and then either return a value (as in $\mathcal{L}\{\text{cmd}\}$) or raise an exception. The language $\mathcal{L}\{\text{comm exc}\}$ is an extension of $\mathcal{L}\{\text{cmd}\}$ to account for this additional outcome of execution. The following grammar specifies the characteristic features of $\mathcal{L}\{\text{comm exc}\}$:

<i>Category</i>	<i>Item</i>	<i>Abstract</i>	<i>Concrete</i>
Comm	m	$::= \text{raise}[\tau](e)$ $\text{letcomp}(e; x.m_1; y.m_2)$	$\text{raise}(e)$ $\text{let cmd}(x) \text{ be } e \text{ in } m_1 \text{ ow}(y) \text{ in } m_2$

This grammar extends that of $\mathcal{L}\{\text{cmd}\}$ with a new primitive command, $\text{raise}(e)$, that raises an exception with value e . It also modifies the grammar of $\mathcal{L}\{\text{cmd}\}$ to generalize the monadic bind construct to include an exception handler. The command

$$\text{let cmd}(x) \text{ be } e \text{ in } m_1 \text{ ow}(y) \text{ in } m_2$$

executes the encapsulated command determined by evaluation of e . If it returns normally, then the return value is bound to x and the command m_1 is executed. If, instead, it raises an exception, the exception value is bound to y and the command m_2 is executed instead. The monadic bind construct of $\mathcal{L}\{\text{cmd}\}$ is to be regarded as short-hand for the command

$$\text{let cmd}(x) \text{ be } e \text{ in } m \text{ ow}(y) \text{ in } \text{raise}(y),$$

which behaves as before in the case of a normal return, and propagates any exception in that case of an exceptional return.

The static semantics of these constructs is given by the following rules:

$$\frac{\Gamma \vdash e : \tau_{\text{exn}}}{\Gamma \vdash \text{raise}[\tau](e) \sim \tau} \quad (38.1a)$$

$$\frac{\Gamma \vdash e : \text{cmd}(\tau) \quad \Gamma, x : \tau \vdash m_1 \sim \tau' \quad \Gamma, y : \tau_{\text{exn}} \vdash m_2 \sim \tau'}{\Gamma \vdash \text{letcomp}(e; x.m_1; y.m_2) \sim \tau'} \quad (38.1b)$$

The dynamic semantics of these commands consists of a transition system of the form $m \mapsto m'$ defined by the following rules:

$$\frac{e \mapsto e'}{\text{return}(e) \mapsto \text{return}(e')} \quad (38.2a)$$

$$\frac{e \mapsto e'}{\text{raise}[\tau](e) \mapsto \text{raise}[\tau](e')} \quad (38.2b)$$

$$\frac{e \mapsto e'}{\text{letcomp}(e; x.m_1; y.m_2) \mapsto \text{letcomp}(e'; x.m_1; y.m_2)} \quad (38.2c)$$

$$\frac{m \mapsto m'}{\text{letcomp}(\text{cmd}(m); x.m_1; y.m_2) \mapsto \text{letcomp}(\text{cmd}(m'); x.m_1; y.m_2)} \quad (38.2d)$$

$$\frac{e \text{ val}}{\text{letcomp}(\text{cmd}(\text{return}(e)); x.m_1; y.m_2) \mapsto [e/x]m_1} \quad (38.2e)$$

$$\frac{e \text{ val}}{\text{letcomp}(\text{cmd}(\text{raise}[\tau](e)); x.m_1; y.m_2) \mapsto [e/y]m_2} \quad (38.2f)$$

38.2 Programming With Monadic Exceptions

The chief virtue of monadic exceptions is also its chief vice. A value of type $\text{nat} \rightarrow \text{nat}$ remains a function that, when applied to a natural number, returns a natural number (or, in the case of partial functions, may diverge). If a function can raise an exception when called, then it must be given the weaker type $\text{nat} \rightarrow \text{nat cmd}$, which specifies that, when applied, it yields an encapsulated computation that, when executed, may raise an exception. Two such functions cannot be directly composed, since their types are no longer compatible. Instead we must explicitly sequence their execution. For example, to compose f and g of this type, we may write

$$\lambda(x:\text{nat}. \text{do } \{y \leftarrow \text{run } g(x); z \leftarrow \text{run } f(y); \text{return}(z)\}).$$

Here we have used the `do` syntax introduced in Chapter 37, which according to our conventions above, implicitly propagates exceptions arising from the application of f and g to their surrounding context.

This distinction may be regarded as either a boon or a bane, depending on how important it is to indicate in the type whether a function might raise an exception when called. For programmer-defined exceptions one may wish to draw the distinction, but the situation is less clear for other forms of run-time errors. For example, if division by zero is to be regarded as a form of exception, then the type of division must be

$$\text{nat} \rightarrow \text{nat} \rightarrow \text{nat cmd}$$

to reflect this possibility. But then one cannot then use division in an ordinary arithmetic expression, because its result is not a number, but an encapsulated command. One response to this might be to consider division by zero, and other related faults, not as handle-able exceptions, but rather as fatal errors that abort computation. In that case there is no difference between such an error and divergence: the computation never terminates, and this condition cannot be detected during execution. Consequently, operations such as division may be regarded as partial functions, and may therefore be used freely in expressions without taking special pains to manage any errors that may arise.

38.3 Exercises

Chapter 39

Monadic State

In Chapter 35 we introduced the type of cells of a given type so as to distinguish mutable from immutable data structures. In that chapter we left open the question of how to integrate mutation into a full-scale language. There are two main methods of doing so, one that permits great flexibility in the use of mutable storage at the expense of weakening the meaning of the typing judgement considerably, and one that retains the meaning of the typing judgement, but impairs the use of storage effects considerably. As this description suggest, each approach has its benefits and drawbacks, with neither being clearly preferable to the other in all circumstances.

The simplest, and most obvious, approach, which we will call the *integral* style, is to enrich a purely functional language, such as $\mathcal{L}\{\text{nat} \rightarrow\}$ or its extensions, with the mechanisms of $\mathcal{L}\{\text{ref}\}$. This results in an integration of imperative and functional programming in which the programmer may, at will, use or eschew mutation at any point within a program. For example, if we start with a purely functional data structure such as a tree structure represented as a recursive type, and then we wish to instrument this structure with, say, a reference count for profiling purposes, we may simply revise the definition of the type to, say, attach a mutable cell to each node that maintains the profiling information. It is usually straightforward to extend the implementation of the tree operations to account for the additional information at the nodes.

The chief drawback of the integral approach is that the meaning of the typing judgement changes drastically. The judgement $e : \tau$ no longer means “if e evaluates to a value v , then v is a value of type τ .” Instead, the judgement now means that, in addition, that evaluation of e can engender arbitrary *side effects* on any data structure to which e has (direct or indirect)

access. (Indeed, side effects are so-called precisely because they act “on the side,” without their influence being reflected in the type of the expression.) Consequently, the type $\text{nat} \rightarrow \text{nat}$ can no longer be understood as the type of partial functions on the natural numbers, but must also admit the possibility of side effects during its execution. As a case in point, in a language without mutation the type $\text{unit} \rightarrow \text{unit}$ is quite trivial, containing only the identity and the divergent function, whereas in a language with integral mutation, this type contains arbitrarily complex functions that mutate storage, with the type revealing nothing about this behavior.

The integral approach works best with a strict language, in which the order of evaluation of sub-expressions is fully determined by its form, and is not sensitive to the evolution of the computation. Laziness complicates the integral approach because it makes it much harder to predict when expressions are evaluated. In the absence of storage effects this is of no concern (at least for correctness, if not efficiency), but in the presence of storage effects, the indeterminacy of evaluation order is disastrous. It is difficult to tell exactly when, or how often, a cell will be allocated or assigned, greatly complicating reasoning about program correctness.

The *monadic* approach to storage effects is to confine operations that may affect storage to the command level of $\mathcal{L}\{\text{cmd}\}$. This ensures that the expression level remains pure, so that it is compatible with both an eager and a lazy interpretation. The chief benefit of the monadic style is that it makes explicit in the types any reliance on storage effects. The chief drawback of the monadic style is that it makes explicit in the types any reliance on storage effects. While it can be useful to document the use of storage effects, it can also be a hindrance to program development. For example, if we wish to instrument a piece of pure code with code for profiling, then we must restructure it to permit modifications to the store for profiling purposes, even though its functionality has not changed.

39.1 Storage Effects

The language $\mathcal{L}\{\text{cmd ref}\}$ is an extension of $\mathcal{L}\{\text{cmd}\}$ (described in Chapter 37) with mutable references into $\mathcal{L}\{\text{cmd}\}$. The syntax of $\mathcal{L}\{\text{cmd ref}\}$

extends that of $\mathcal{L}\{\text{cmd}\}$ with the following constructs:

<i>Category</i>	<i>Item</i>	<i>Abstract</i>	<i>Concrete</i>
Type	τ	$::= \text{ref}(\tau)$	$\tau \text{ ref}$
Expr	e	$::= l$	l
Comm	m	$::= \text{ref}(e)$	$\text{ref}(e)$
		$\text{get}(e)$	$! e$
		$\text{set}(e_1; e_2)$	$e_1 \leftarrow e_2$

Locations are pure expressions, whereas the primitives for reference cells are forms of command.

The static semantics of $\mathcal{L}\{\text{cmd ref}\}$ extends that of $\mathcal{L}\{\text{cmd}\}$ with the following rules:

$$\frac{}{\Lambda, l : \tau \Gamma \vdash l : \text{ref}(\tau)} \quad (39.1a)$$

$$\frac{\Lambda \Gamma \vdash e : \tau}{\Lambda \Gamma \vdash \text{ref}(e) \sim \text{ref}(\tau)} \quad (39.1b)$$

$$\frac{\Lambda \Gamma \vdash e : \text{ref}(\tau)}{\Lambda \Gamma \vdash \text{get}(e) \sim \tau} \quad (39.1c)$$

$$\frac{\Lambda \Gamma \vdash e_1 : \text{ref}(\tau) \quad \Lambda \Gamma \vdash e_2 : \tau}{\Lambda \Gamma \vdash \text{set}(e_1; e_2) \sim \tau} \quad (39.1d)$$

Here we make explicit the location typing assumptions, Λ , as well as the variable typing assumptions, Γ .

The dynamic semantics of $\mathcal{L}\{\text{cmd ref}\}$ is structured into two parts:

1. A transition relation $e \mapsto e'$ for expressions.
2. A transition relation $m @ \mu \mapsto m' @ \mu'$ for commands.

Expressions are evaluated without regard to context, since they engender no effects, whereas commands are evaluated relative to a memory, on which they may have an effect.

The rules defining the dynamic semantics of the monad constructs are as follows.

$$\overline{\text{cmd}(m) \text{ val}} \quad (39.2a)$$

$$\frac{e \mapsto e'}{\text{return}(e) @ \mu \mapsto \text{return}(e') @ \mu} \quad (39.2b)$$

$$\frac{e \mapsto e'}{\text{letcmd}(e; x.m) @ \mu \mapsto \text{letcmd}(e'; x.m) @ \mu} \quad (39.2c)$$

$$\frac{m_1 @ \mu \mapsto m'_1 @ \mu'}{\text{letcmd}(\text{cmd}(m_1); x.m_2) @ \mu \mapsto \text{letcmd}(\text{cmd}(m'_1); x.m_2) @ \mu'} \quad (39.2d)$$

$$\frac{e \text{ val}}{\text{letcmd}(\text{cmd}(\text{return}(e)); x.m) @ \mu \mapsto [e/x]m @ \mu} \quad (39.2e)$$

The transition rules for the monadic elimination form is somewhat unusual. First, the expression e is evaluated to obtain a suspended command. Once such a command has been obtained, execution continues by evaluating it in the current memory, updating that memory as appropriate during its execution. This process ends once the suspended command is a return statement, in which case this value is passed to the body of the `letcomp`.

The transition rules for evaluation of storage commands are as follows:

$$\overline{l \text{ val}} \quad (39.3a)$$

$$\frac{e \mapsto e'}{\text{ref}(e) @ \mu \mapsto \text{ref}(e') @ \mu} \quad (39.3b)$$

$$\frac{e \text{ val}}{\text{ref}(e) @ \mu \mapsto \text{return}(l) @ \mu \otimes \langle l : e \rangle} \quad (39.3c)$$

$$\frac{e \mapsto e'}{\text{get}(e) @ \mu \mapsto \text{get}(e') @ \mu} \quad (39.3d)$$

$$\frac{e \text{ val}}{\text{get}(l) @ \mu \otimes \langle l : e \rangle \mapsto \text{return}(e) @ \mu \otimes \langle l : e \rangle} \quad (39.3e)$$

$$\frac{e_1 \mapsto e'_1}{\text{set}(e_1; e_2) @ \mu \mapsto \text{set}(e'_1; e_2) @ \mu} \quad (39.3f)$$

$$\frac{e_1 \text{ val} \quad e_2 \mapsto e'_2}{\text{set}(e_1; e_2) @ \mu \mapsto \text{set}(e_1; e'_2) @ \mu} \quad (39.3g)$$

$$\frac{e \text{ val}}{\text{set}(l; e) @ \mu \otimes \langle l : e' \rangle \mapsto \text{return}(e) @ \mu \otimes \langle l : e \rangle} \quad (39.3h)$$

39.2 Integral versus Monadic Effects

The chief motivation for introducing monads is to make explicit in the types any reliance on computational effects. In the case of storage effects this is not always an advantage. The problem is that any use of storage forces the computation to be within the monad, and there is no way to get back out—once in the monad, always in the monad. This rules out the use of so-called *benign effects*, which may be used internally in some computation that is, for all outward purposes, entirely pure. One example of this is provided by splay trees, which may be used to implement a functional dictionary abstraction, but which rely heavily on mutation for their implementation in order to ensure efficiency. A simpler example, which we consider in detail, is provided by the use of backpatching to implement recursion as described in Chapter 35.

When formulated using monads to expose the use of storage effects, the backpatching implementation, *fact*, of the factorial function is as follows:

```
do {
  r ← return (ref (λ n:nat. comp(return (n))))
  ; f ← return (λ n:nat. ...)
  ; _ ← set (r, f)
  ; return f
}
```

where the elided λ -abstraction is given as follows:

```
λ(n:nat.
  ifz(n,
    comp(return(1)),
    n'.comp(
      do {
        f' ← get(r)
        ; return (n*f'(n'))
      })))
```

Observe that each branch of the conditional test returns a suspended command. In the case that the argument is zero, the command simply returns the value 1. Otherwise, it fetches the contents of the associated reference cell, applies this to the predecessor, and returns the result of the appropriate calculation.

We may check that that $fact \sim \text{nat} \rightarrow (\text{nat cmd})$, which exposes two aspects of this code:

1. The command that builds the recursive factorial function is impure, because it allocates and assigns to the reference cell used to implement backpatching.
2. The body of the factorial function is itself impure, because it accesses the reference cell to effect the recursive call.

The consequence is that the factorial function may no longer be used as a (pure) function! In particular, we cannot apply *fact* to an argument in an expression; it must be executed as a command. We must write

```
do {
  f ← return(fact)
  ; x ← let comp (x:nat) be f(n) in return x
  ; return x
}
```

to bind the function computed by the expression *fact* to the variable *f*; apply this to *n*, yielding the result; and return this to the caller.

The difficulty is that the use of a reference cell to implement recursion is a benign effect, one that does not affect the purity of the function expression itself, nor of its applications. But the type system for effects studied here is incapable of recognizing this fact, and for good reason. It is extremely difficult, in general, to determine whether or not the use of effects in some region of a program is benign. As a stop-gap measure, one way around this is to introduce an operation of type $\tau \text{ cmd} \rightarrow \tau$, which may be used to exit the monad. But this ruins the very distinction we are trying to enforce, to segregate pure expressions from impure commands.

39.3 Exercises

Chapter 40

Comonads

Monads arise naturally for managing effects that both *influence* and are *influenced by* the context in which they arise. This is particularly clear for storage effects, whose context is a memory mapping locations to values. The semantics of the storage primitives makes reference to the memory (to retrieve the contents of a location) and makes changes to the memory (to change the contents of a location or allocate a new location). These operations must be sequentialized in order to be meaningful (that is, the precise order of execution matters), and we cannot expect to escape the context since locations are values that give rise to dependencies on the context. As we shall see in Chapter 46 other forms of effect, such as input/output or interprocess communication, are naturally expressed in the context of a monad.

By contrast the use of monads for exceptions as in Chapter 38 is rather less natural. Raising an exception does not influence the context, but rather imposes the requirement on it that a handler be present to ensure that the command is meaningful even when an exception is raised. One might argue that installing a handler influences the context, but it does so in a nested, or stack-like, manner. A new handler is installed for the duration of execution of a command, and then discarded. The handler does not persist across commands in the same sense that locations persist across commands in the case of the state monad. Moreover, installing a handler may be seen as restoring purity in that it catches any exceptions that may be raised and, assuming that the handler does not itself raise an exception, yields a pure value. A similar situation arises with fluid binding (as described in Chapter 34). A reference to a symbol imposes the demand on the context to provide a binding for it. The binding of a symbol may be changed, but only

for the duration of execution of a command, and not persistently. Moreover, the reliance on symbol bindings within a specified scope confines the impurity to that scope.

The concept of a *comonad* captures the concept of an effect that *imposes a requirement* on its context of execution, but that does not persistently alter that context beyond its execution. Computations that rely on the context to provide some capability may be thought of as impure, but the impurity is confined to the extent of the reliance—outside of this context the computation may be once again regarded as pure. One may say that monads are appropriate for *global*, or *persistent*, effects, whereas comonads are appropriate for *local*, or *ephemeral*, effects.

40.1 A Comonadic Framework

The central concept of the comonadic framework for effects is the *constrained typing judgement*, $e : \tau [\chi]$, which states that an expression e has type τ (as usual) provided that the context of its evaluation satisfies the constraint χ . The nature of constraints varies from one situation to another, but will include at least the trivially true constraint, \top , and the conjunction of constraints, $\chi_1 \wedge \chi_2$. We sometimes write $e : \tau$ to mean $e : \tau [\top]$, which states that expression e has type τ under no constraints.

The syntax of the comonadic framework, $\mathcal{L}\{\text{comon}\}$, is given by the following grammar:

Category	Item		Abstract	Concrete
Type	τ	::=	$\text{box}[\chi](\tau)$	$\square_{\chi} \tau$
Const	χ	::=	tt	\top
			$\text{and}(\chi_1; \chi_2)$	$\chi_1 \wedge \chi_2$
Expr	e	::=	$\text{box}(e)$	$\text{box}(e)$
			$\text{unbox}(e)$	$\text{unbox}(e)$

A type of the form $\square_{\chi} \tau$ is called a *comonad*; it represents the type of unevaluated expressions that impose constraint χ on their context of execution. The constraint \top is the trivially true constraint, and the constraint $\chi_1 \wedge \chi_2$ is the conjunction of two constraints. The expression $\text{box}(e)$ is the introduction form for the comonad, and the expression $\text{unbox}(e)$ is the corresponding elimination form.

The judgement χ true expresses that the constraint χ is satisfied. This judgement is partially defined by the following rules, which specify the

meanings of the trivially true constraint and the conjunction of constraints.

$$\overline{\text{tt true}} \quad (40.1a)$$

$$\frac{\chi_1 \text{ true} \quad \chi_2 \text{ true}}{\text{and}(\chi_1; \chi_2) \text{ true}} \quad (40.1b)$$

$$\frac{\text{and}(\chi_1; \chi_2) \text{ true}}{\chi_1 \text{ true}} \quad (40.1c)$$

$$\frac{\text{and}(\chi_1; \chi_2) \text{ true}}{\chi_2 \text{ true}} \quad (40.1d)$$

We will make use of hypothetical judgements of the form $\chi_1 \text{ true}, \dots, \chi_n \text{ true} \vdash \chi \text{ true}$, where $n \geq 0$, expressing that χ is derivable from χ_1, \dots, χ_n , as usual.

The static semantics is specified by generic hypothetical judgements of the form

$$x_1 : \tau_1 [\chi_1], \dots, x_n : \tau_n [\chi_n] \vdash e : \tau [\chi].$$

As usual we write Γ for a finite set of hypotheses of the above form.

The static semantics of the core constructs of $\mathcal{L}\{\text{comon}\}$ is defined by the following rules:

$$\frac{\chi' \vdash \chi}{\Gamma, x : \tau [\chi] \vdash x : \tau [\chi']} \quad (40.2a)$$

$$\frac{\Gamma \vdash e : \tau [\chi]}{\Gamma \vdash \text{box}(e) : \square_\chi \tau [\chi']} \quad (40.2b)$$

$$\frac{\Gamma \vdash e : \square_\chi \tau [\chi'] \quad \chi' \vdash \chi}{\Gamma \vdash \text{unbox}(e) : \tau [\chi']} \quad (40.2c)$$

Rule (40.2b) states that a boxed computation has comonadic type under an arbitrary constraint. This is valid because a boxed computation is a value, and hence imposes no constraint on its context of evaluation. Rule (40.2c) states that a boxed computation may be activated provided that the ambient constraint, χ' , is at least as strong as the constraint χ of the boxed computation. That is, any requirement imposed by the boxed computation must be met at the point at which it is unboxed.

Rules (40.2) are formulated to ensure that the constraint on a typing judgement may be strengthened arbitrarily.

Lemma 40.1 (Constraint Strengthening). *If $\Gamma \vdash e : \tau [\chi]$ and $\chi' \vdash \chi$, then $\Gamma \vdash e : \tau [\chi']$.*

Proof. By rule induction on Rules (40.2). \square

Intuitively, if a typing holds under a weaker constraint, then it also holds under any stronger constraint as well.

At this level of abstraction the dynamic semantics of $\mathcal{L}\{\text{comon}\}$ is trivial.

$$\overline{\text{box}(e) \text{ val}} \quad (40.3a)$$

$$\frac{e \mapsto e'}{\text{unbox}(e) \mapsto \text{unbox}(e')} \quad (40.3b)$$

$$\overline{\text{unbox}(\text{box}(e)) \mapsto e} \quad (40.3c)$$

In specific applications of $\mathcal{L}\{\text{comon}\}$ the dynamic semantics will also specify the context of evaluation with respect to which constraints are to interpreted.

The role of the comonadic type in $\mathcal{L}\{\text{comon}\}$ is explained by considering how one might extend the language with, say, function types. The crucial idea is that the comonad isolates the dependence of a computation on its context of evaluation so that such constraints do not affect the other type constructors. For example, here are the rules for function types expressed in the context of $\mathcal{L}\{\text{comon}\}$:

$$\frac{\Gamma, x : \tau_1 [\text{tt}] \vdash e_2 : \tau_2 [\text{tt}]}{\Gamma \vdash \text{lam}[\tau_1](x.e_2) : \text{arr}(\tau_1; \tau_2) [\chi]} \quad (40.4a)$$

$$\frac{\Gamma \vdash e_1 : \tau_2 \rightarrow \tau [\chi] \quad \Gamma \vdash e_2 : \tau_2 [\chi]}{\Gamma \vdash \text{ap}(e_1; e_2) : \tau [\chi]} \quad (40.4b)$$

These rules are formulated so as to ensure that constraint strengthening remains admissible. Rule (40.4a) states that a λ -abstraction has type $\tau_1 \rightarrow \tau_2$ under any constraint χ provided that its body has type τ_2 under the trivially true constraint, assuming that its argument has type τ_1 under the trivially true constraint. By demanding that the body be well-formed under no constraints we are, in effect, insisting that its body be boxed if it is to impose a constraint on the context at the point of application. Under a call-by-value evaluation order, the argument x will always be a value, and hence imposes no constraints on its context.

Let the expression $\text{unbox_app}(e_1; e_2)$ be an abbreviation for $\text{unbox}(\text{ap}(e_1; e_2))$, which applies e_1 to e_2 , then activates the result. The derived static semantics

for this construct is given by the following rule:

$$\frac{\Gamma \vdash e_1 : \tau_2 \rightarrow \Box_\chi \tau [\chi'] \quad \Gamma \vdash e_2 : \tau_2 [\chi'] \quad \chi' \vdash \chi}{\Gamma \vdash \text{unbox_app}(e_1; e_2) : \tau [\chi']} \quad (40.5)$$

In words, to apply a function with impure body to an argument, the ambient constraint must be strong enough to type the function and its argument, and must be at least as strong as the requirements imposed by the body of the function. We may view a type of the form $\tau_1 \rightarrow \Box_\chi \tau_2$ as the type of functions that, when applied to a value of type τ_1 , yield a value of type τ_2 engendering local effects with requirements specified by χ .

Similar principles govern the extension of $\mathcal{L}\{\text{comon}\}$ with other types such as products or sums.

40.2 Comonadic Effects

In this section we discuss two applications of $\mathcal{L}\{\text{comon}\}$ to managing local effects. The first application is to exceptions, using constraints to specify whether or not an exception handler must be installed to evaluate an expression so as to avoid an uncaught exception error. The second is to fluid binding, using constraints to specify which symbols must be bound during execution so as to avoid accessing an unbound symbol. The first may be considered to be an instance of the second, in which we think of the exception handler as a distinguished symbol whose binding is the current exception continuation.

40.2.1 Exceptions

To model exceptions we extend $\mathcal{L}\{\text{comon}\}$ as follows:

<i>Category</i>	<i>Item</i>		<i>Abstract</i>		<i>Concrete</i>
Const	χ	::=	\uparrow		\uparrow
Expr	e	::=	$\text{raise}[\tau](e)$		$\text{raise}(e)$
			$\text{handle}(e_1; x.e_2)$		$\text{try } e_1 \text{ ow } x \Rightarrow e_2$

The constraint \uparrow specifies that an expression may raise an exception, and hence that its context is required to provide a handler for it.

The static semantics of $\mathcal{L}\{\text{comon}\}$ is extended with the following rules:

$$\frac{\Gamma \vdash e : \tau_{\text{exn}} [\chi] \quad \chi \vdash \uparrow}{\Gamma \vdash \text{raise}[\tau](e) : \tau [\chi]} \quad (40.6a)$$

$$\frac{\Gamma \vdash e_1 : \tau [\chi \wedge \uparrow] \quad \Gamma, x : \tau_{\text{exn}} \vdash e_2 : \tau [\chi]}{\Gamma \vdash \text{handle}(e_1; x.e_2) : \tau [\chi]} \quad (40.6b)$$

Rule (40.6a) imposes the requirement for a handler on the context of a raise expression, in addition to any other conditions that may be imposed by its argument. (The rule is formulated so as to ensure that constraint strengthening remains admissible.) Rule (40.6b) transforms an expression that requires a handler into one that may or may not require one, according to the demands of the handling expression. If e_2 does not demand a handler, then χ may be taken to be the trivial constraint, in which case the overall expression is pure, even though e_1 is impure (may raise an exception).

The dynamic semantics of exceptions is as given in Chapter 28. The interesting question is to explore the additional assurances given by the comonadic type system given by Rules (40.6). Intuitively, we may think of a stack as a constraint transformer that turns a constraint χ into a constraint χ' by composing frames, including handler frames. Then if e is an expression of type τ imposing constraint χ and k is a τ -accepting stack transforming constraint χ into constraint \top , then evaluation of e on k cannot yield an uncaught exception. In this sense the constraints reflect the reality of the execution behavior of expressions.

To make this precise, we define the judgement $k : \tau [\chi]$ to mean that k is a stack that is suitable as an execution context for an expression $e : \tau [\chi]$. The typing rules for stacks are as follows:

$$\overline{\epsilon : \tau [\top]} \quad (40.7a)$$

$$\frac{k : \tau' [\chi'] \quad f : \tau [\chi] \Rightarrow \tau' [\chi']}{k; f : \tau [\chi]} \quad (40.7b)$$

Rule (40.7a) states that the empty stack must not impose any constraints on its context, which is to say that there must be no uncaught exceptions at the end of execution. Rule (40.7b) simply specifies that a stack is a composition of frames. The typing rules for frames are easily derived from the static semantics of $\mathcal{L}\{\text{comon}\}$. For example,

$$\frac{x : \tau_{\text{exn}} \vdash e : \tau [\chi]}{\text{handle}(-; x.e) : \tau [\chi \wedge \uparrow] \Rightarrow \tau [\chi]} \quad (40.8)$$

This rule states that a handler frame transforms an expression of type τ demanding a handler into an expression of type τ that may, or may not, demand a handler, according to the form of the handling expression.

The formation of states is defined essentially as in Chapter 27.

$$\frac{k : \tau[\chi] \quad e : \tau[\chi]}{k \triangleright e \text{ ok}} \quad (40.9a)$$

$$\frac{k : \tau[\chi] \quad e : \tau[\chi] \quad e \text{ val}}{k \triangleleft e \text{ ok}} \quad (40.9b)$$

Observe that a state of the form $\epsilon \triangleright \text{raise}(e)$, where $e \text{ val}$, is ill-formed, because the empty stack is well-formed only under no constraints on the context.

Safety ensures that no uncaught exceptions can arise. This is expressed by defining final states to be only those returning a value to the empty stack.

$$\frac{e \text{ val}}{\epsilon \triangleleft e \text{ final}} \quad (40.10)$$

In contrast to Chapter 28, we *do not* consider an uncaught exception state to be final!

Theorem 40.2 (Safety). 1. *If $s \text{ ok}$ and $s \mapsto s'$, then $s' \text{ ok}$.*

2. *If $s \text{ ok}$ then either $s \text{ final}$ or there exists s' such that $s \mapsto s'$.*

Proof. These are proved by rule induction on the dynamic semantics and on the static semantics, respectively, proceeding along standard lines. \square

40.2.2 Fluid Binding

Using comonads we may devise a type system for fluid binding that ensures that no unbound symbols are accessed during execution. This is achieved by regarding the mapping of symbols to their values to be the context of execution, and introducing a form of constraint stating that a specified symbol must be bound in the context.

Let us consider a comonadic static semantics for $\mathcal{L}\{\text{fluid}\}$ defined in Chapter 34. For this purpose we consider atomic constraints of the form $\text{bd}(a)$, stating that the symbol a has a binding.

The static semantics of fluid binding consists of judgements of the form $\Sigma \Gamma \vdash e : \tau[\chi]$, where Σ consists of hypotheses of the form $a : \tau$ assigning a type to a symbol.

$$\frac{\Sigma \vdash a : \tau \quad \chi \vdash \text{bd}(a)}{\Sigma \Gamma \vdash \text{get}[a] : \tau[\chi]} \quad (40.11a)$$

$$\frac{\Sigma \vdash a : \tau \quad \Sigma \Gamma \vdash e_1 : \tau [\chi] \quad \Sigma \Gamma \vdash e_2 : \tau [\chi \wedge \text{bd}(a)]}{\Sigma \Gamma \vdash \text{put}[a](e_1; e_2) : \tau [\chi]} \quad (40.11b)$$

Rule (40.11a) records the demand for a binding for the symbol a incurred by retrieving its value. Rule (40.11b) propagates the fact that the symbol a is bound to the body of the fluid binding.

The dynamic semantics is as specified in Chapter 34. The safety theorem for the comonadic type system for fluid binding states that no unbound symbol error may ever arise during execution. We define the judgement $\theta \models \chi$ to mean that $a \in \text{dom}(\theta)$ whenever $\chi \vdash \text{bd}(a)$.

Theorem 40.3 (Safety). 1. If $e : \tau [\chi]$ and $e \mapsto_{\theta} e'$, then $e' : \theta [\chi]$.

2. If $e : \tau [\chi]$ and $\theta \models \chi$, then either e val or there exists e' such that $e \mapsto_{\theta} e'$.

The comonadic static semantics for $\mathcal{L}\{\text{fluid}\}$ may be extended to $\mathcal{L}\{\text{fluidnew}\}$, which also permits dynamic symbol generation. The main difficulty is to manage the interaction between the scopes of symbols and their occurrences in types. First, it is straightforward to define the judgement $\Sigma \vdash \chi$ constr to mean that χ is a constraint involving only those symbols a such that $\Sigma \vdash a : \tau$ for some τ . Using this we may also define the judgement $\Sigma \vdash \tau$ type analogously. This judgement is used to impose a restriction on symbol generation to ensure that symbols do not escape their scope:

$$\frac{\Sigma, a : \sigma \Gamma \vdash e : \tau \quad \Sigma \vdash \tau \text{ type}}{\Sigma \Gamma \vdash \text{new}[\sigma](a.e) : \langle \sigma \rangle \tau} \quad (40.12)$$

This imposes the requirement that the result type of a computation involving a dynamically generated symbol must not mention that symbol. Otherwise the type $\langle \sigma \rangle \tau$ would involve a symbol that makes no sense with respect to the ambient symbol context, Σ . In practical terms this means that the expression, e , must ensure that its type imposes no residual requirements involving the symbol a introduced by the binder.

For example, an expression such as

$$\text{gen}(v(a:\text{nat}. \text{put } a \text{ is } z \text{ in } \lambda(x:\text{nat}. \text{box}(\dots \text{get } a \dots))))$$

is necessarily ill-typed. The type of the λ -abstraction must be of the form $\text{nat} \rightarrow \square_{\chi} \tau$, where $\chi \vdash \text{bd}(a)$, reflecting the dependence of the body of the function on the binding of a . This type is propagated through the fluid binding for a , since it holds only for the duration of evaluation of the λ -abstraction itself, which is immediately returned as its value. Since

the type of the λ -abstraction involves the symbol a , the second premise of Rule (40.12) is not met, and the expression is ill-typed. This is as it should be, for we cannot guarantee that the dynamically generated symbol replacing a during evaluation will, in fact, be bound when the body of the function is executed.

However, if we move the binding for a into the scope of the λ -abstraction,

$$\text{gen}(v(a:\text{nat}.\lambda(x:\text{nat}.\text{box}(\text{put } a \text{ is z in } \dots \text{get } a \dots))))),$$

then the type of the λ -abstraction may have the form $\text{nat} \rightarrow \square_{\chi} \tau$, where χ need not constrain a to be bound. The reason is that the fluid binding for a discharges the obligation to bind a within the body of the function. Consequently, the condition on Rule (40.12) is met, and the expression is well-typed. Indeed, each evaluation of the body of the λ -abstraction initializes the fresh copy of a generated during evaluation, so no unbound symbol error can arise during execution.

40.3 Exercises

