

The Science of Functional Programming

Applied Functional Type Theory, with Examples in Scala

by Sergei Winitzki, Ph.D.

draft version 0.2, April 29, 2019

Published by **lulu.com** 2019

Copyright © 2018-2019 by Sergei Winitzki

Published and printed by lulu.com

Permission is granted to copy, distribute and/or modify this document under the terms of the GNU Free Documentation License, Version 1.2 or any later version published by the Free Software Foundation; with no Invariant Sections, no Front-Cover Texts, and no Back-Cover Texts. A copy of the license is included in the section entitled “GNU Free Documentation License” (Appendix [G](#)).

ISBN XXXXXX

This book presents the theoretical knowledge that helps to write code in the functional programming paradigm. Detailed explanations and derivations are logically developed and accompanied by worked examples tested in the Scala interpreter. Exercises with solutions are provided. Readers need to have a working knowledge of basic Scala (e.g. be able to write code that reads a small text file and prints all word counts, sorted in descending order by count). Readers should also have a basic command of school-level mathematics; for example, be able to simplify the expressions such as $\frac{1}{x-2} - \frac{1}{x+2}$ and $\frac{d}{dx} ((x+1)e^{-x})$.

Contents

1	Values, types, expressions, functions	1
1.1	Translating mathematics into code	1
1.1.1	First examples	1
1.1.2	Nameless functions	3
1.1.3	Nameless functions and bound variables	6
1.1.4	Aggregating data from lists	9
1.1.5	Filtering	11
1.2	Worked examples: transformation and aggregation	13
1.3	Summary	16
1.4	Exercises	18
1.5	Discussion	19
1.5.1	Functional programming as a paradigm	19
1.5.2	Functional programming languages	20
1.5.3	The mathematical meaning of variables	20
1.5.4	Iteration without loops	22
1.5.5	Nameless functions in mathematical notation	23
1.5.6	Named and nameless expressions and their uses	25
1.5.7	Nameless functions: historical perspective	27
2	Mathematical induction	29
2.1	Tuple types	29
2.1.1	First examples	29
2.1.2	Pattern matching on tuples	31
2.1.3	Using tuples with collections	33
2.1.4	Using dictionaries (Scala's <code>Maps</code>) as collections	34
2.1.5	Worked examples: Tuples and collections	37
2.1.6	Exercises: Tuples and collections	41
2.2	Discussion	43
2.2.1	Total and partial functions	43
2.2.2	Scope of pattern matching variables	44

3	The formal logic of types	45
3.1	Types of higher-order functions	45
3.1.1	Curried functions	45
3.1.2	Calculations with nameless functions	47
3.1.3	Short syntax for function applications	49
3.1.4	Higher-order functions	49
3.1.5	Worked examples: higher-order functions	51
3.2	Discussion	52
3.2.1	Scope of bound variables	52
3.3	Exercises	54
3.4	Disjunctive types	54
3.4.1	Discussion	54
3.5	The Curry-Howard correspondence	54
3.5.1	Discussion	54
4	Functors	55
4.1	Discussion	55
4.2	Practical use	55
4.3	Laws and structure	55
5	Type-level functions and type classes	57
5.1	Combining typeclasses	57
5.2	Inheritance	57
5.3	Functional dependencies	57
5.4	Discussion	57
6	Filterable functors	59
6.1	Practical use	59
6.1.1	Discussion	59
6.2	Laws and structure	59
6.2.1	Discussion	59
7	Semimonads and monads	61
7.1	Practical use	61
7.1.1	Discussion	61
7.2	Laws and structure	61
7.2.1	Discussion	61

8	Applicative functors, contrafunctors, and profunctors	63
8.1	Practical use	63
8.1.1	Discussion	63
8.2	Laws and structure	63
9	Traversable functors and profunctors	65
9.1	Discussion	65
10	Free constructions	67
10.1	Discussion	67
11	Monad transformers	69
11.1	Practical use	69
11.2	Laws and structure	69
11.2.1	Properties of monadic morphisms	70
11.3	Functor composition using transformed monads	70
11.4	Stacking composition of monad transformers	71
11.4.1	Stacking two monads	71
11.4.2	Stacking multiple monads	74
11.5	Monad transformers that use functor composition	75
11.5.1	Motivation for the <code>swap</code> function	76
11.5.2	Laws for <code>swap</code>	79
11.5.3	Intuition behind the laws of <code>swap</code>	84
11.5.4	Deriving <code>swap</code> from <code>flatten</code>	85
11.5.5	Laws of monad transformer liftings	89
11.5.6	Laws of monad transformer runners	90
11.6	Composed-inside transformers: Linear monads	92
11.7	Composed-outside transformers: Rigid monads	92
11.7.1	Rigid monad construction 1: contrafunctor choice	93
11.7.2	Rigid monad construction 2: composition	102
11.7.3	Rigid monad construction 3: product	103
11.7.4	Rigid monad construction 4: selector monad	103
11.8	Irregular monad transformers	104
11.8.1	The state monad transformer <code>StateT</code>	104
11.8.2	The continuation monad transformer <code>ContT</code>	104
11.8.3	The codensity monad transformer <code>CodT</code>	104
11.9	Discussion	104
11.10	Rigid functors	104

12 Recursive types	107
12.1 Fixpoints and type recursion schemes	107
12.2 Row polymorphism and OO programming	107
12.3 Column polymorphism	107
12.4 Discussion	107
13 Co-inductive typeclasses. Comonads	109
13.1 Practical use	109
13.2 Laws and structure	109
13.3 Co-free constructions	109
13.4 Co-free comonads	109
13.5 Comonad transformers	109
13.6 Discussion	109
14 Irregular type classes*	111
14.1 Distributive functors	111
14.2 Lenses and prisms	111
14.3 Discussion	111
15 Value-dependent types*	113
16 Conclusions and discussion	115
17 Essay: Software engineers and software artisans	117
17.1 Engineering disciplines	117
17.2 Artisanishp: Trades and crafts	118
17.3 Programmers today are artisans, not engineers	118
17.3.1 No requirement of formal study	118
17.3.2 No mathematical formalism to guide software <i>development</i>	120
17.3.3 Programmers avoid academic terminology	121
17.4 Towards software engineering	122
17.5 Do we need software <i>engineering</i> , or are artisans good enough?	124
18 Essay: Towards functional data engineering with Scala	127
18.1 Data is math	127
18.2 Functional programming is math	128
18.3 The power of abstraction	129
18.4 Scala is Java on math	130
18.5 Conclusion	131

A	Notations	133
A.1	Summary table	133
A.2	Explanations	134
B	Glossary of terms	139
B.1	On the current misuse of the term “algebra”	141
C	Scala syntax and features	145
C.0.1	Function syntax	145
C.0.2	Scala collections	145
D	Intuitionistic (constructive) propositional logic	147
E	Category theory	149
F	What is “applied functional type theory”	151
F.1	The import and scope of AFTT	151
F.2	Approach to presentation	155
F.3	Intended audience	156
G	GNU Free Documentation License	157
G.0.0	Applicability and definitions	157
G.0.1	Verbatim copying	158
G.0.2	Copying in quantity	158
G.0.3	Modifications	158
H	A humorous disclaimer	161
	Index	163

1 Values, types, expressions, functions

1.1 Translating mathematics into code

1.1.1 First examples

We begin by writing Scala code for some computational tasks.

Factorial of 10 Find the product of integers from 1 to 10 (the **factorial** of 10).
First, we write a mathematical formula for the result:

$$\prod_{k=1}^{10} k \quad .$$

We can then write Scala code in a way that resembles this formula:

```
scala> (1 to 10).product  
res0: Int = 3628800
```

The Scala interpreter indicates that the result is the value 3628800 of type `Int`. If we need to define a name for this value (e.g. to use it later), we use the “`val`” syntax and write

```
scala> val fac10 = (1 to 10).product  
fac10: Int = 3628800  
  
scala> fac10 == 3628800  
res1: Boolean = true
```

The code `(1 to 10).product` is an **expression**, which means two things: (1) it can be computed (e.g. with the Scala interpreter) and yields a value, and (2) the

1 Values, types, expressions, functions

code can be inserted anywhere as part of a larger expression; for example, we could write

```
scala> 100 + (1 to 10).product + 100
res0: Int = 3629000
```

Factorial as a function Define a function that takes an integer n and computes the factorial of n .

A mathematical formula for this function can be written as

$$f(n) = \prod_{k=1}^n k \quad .$$

The corresponding Scala code is

```
def f(n:Int) = (1 to n).product
```

In Scala's `def` syntax, we need to specify the type of a function's argument; in this case, we write `n:Int`. In the usual mathematical notation, types of arguments are either not written at all, or written separately from the formula:

$$f(n) = \prod_{k=1}^n k, \quad \forall n \in \mathbb{N} \quad .$$

This indicates that n must be from the set of non-negative integers (denoted by \mathbb{N} in mathematics). This is quite similar to specifying the type `Int` in the Scala code. So, the argument's type in the code specifies the *domain* of a function.

Having defined the function `f`, we can now apply it to an integer argument:

```
scala> f(10)
res6: Int = 3628800
```

It is an error to apply `f` to a non-integer value, e.g. to a string:

```
scala> f("abc")
<console>:13: error: type mismatch;
 found   : String("abc")
 required: Int
    f("abc")
      ^
```

1.1.2 Nameless functions

The formula and the code, as written above, both involve *naming* the function as “*f*”. Sometimes we do not really need to name that function, – for instance, if the function is used only once. We could denote “nameless” mathematical functions like this:

$$x \Rightarrow (\text{some formula}) \quad .$$

Then the mathematical notation for the nameless factorial function is

$$n \Rightarrow \prod_{k=1}^n k \quad .$$

This reads as “a function that maps *n* to the product of all *k* where *k* goes from 1 to *n*”. The Scala expression implementing this mathematical formula is

```
(n: Int) => (1 to n).product
```

This expression shows Scala’s syntax for a **nameless** function. Here,

```
n: Int
```

is the function’s **argument**, while

```
(1 to n).product
```

is the function’s **body**. The double arrow symbol `=>` separates the argument from the body.¹

Functions in Scala (whether named or nameless) are treated as values, which means that we can also define a Scala value `fac` as

```
scala> val fac = (n: Int) => (1 to n).product
fac: Int => Int = <function1>
```

We see that the value `fac` has the type `Int => Int`, which means that the function takes an integer argument and returns an integer result value. What is the value of the function `fac` *itself*? As we have just seen, the Scala interpreter

¹In Scala, the two ASCII characters `=>` and the single Unicode character \Rightarrow have the same meaning. I use the symbol \Rightarrow (pronounced “maps to”) in this book. However, when I write calculations *by hand*, I tend to write \rightarrow instead of \Rightarrow since it is faster. Several programming languages, such as OCaml and Haskell, use the symbols `->` or the Unicode equivalent, \rightarrow , for the function arrow.

1 Values, types, expressions, functions

prints `<function1>` as the “value” of `fac`. An alternative Scala interpreter called `ammonite`² prints something like this,

```
scala@ val fac = (n: Int) => (1 to n).product
fac: Int => Int =
  ammonite.$sess.cmd0$$$Lambda$1675/2107543287@1e44b638
```

This seems to indicate some identifying number, or perhaps a memory location.

I usually imagine that a function value represents a block of compiled machine code, — code that will actually run and evaluate the function body when the function is applied.

Once defined, a function value can be applied to an argument:

```
scala> fac(10)
res1: Int = 3628800
```

However, functions can be used without naming them. We can directly apply a nameless factorial function to an integer argument 10 instead of writing `fac(10)`:

```
scala> ((n: Int) => (1 to n).product)(10)
res2: Int = 3628800
```

One would not often write code like this because there is no advantage in creating a nameless function and then applying it right away to an argument. This is so because we can evaluate the expression

```
((n: Int) => (1 to n).product)(10)
```

by substituting 10 instead of `n` in the function body, which gives us

```
(1 to 10).product
```

If a nameless function uses the argument several times, for example

```
((n: Int) => n*n*n + n*n)(12345)
```

it is still better to substitute the argument and to eliminate the nameless function. We could have written

```
12345*12345*12345 + 12345*12345
```

²See <https://ammonite.io/>

but, of course, we want to avoid repeating the value 12345. To achieve that, we may define `n` as a value in an **expression block** like this:

```
scala> { val n = 12345; n*n*n + n*n }
res3: Int = 322687002
```

Defined in this way, the value `n` is visible only within the expression block. Outside the block, another value named `n` could be defined independently of this `n`. For this reason, the definition of `n` is called a **local value** definition.

Nameless functions are most useful when they are themselves arguments of other functions, as we will see next.

Checking integers for being prime Define a function that takes an integer n and determines whether n is a prime number.

A simple mathematical formula for this function can be written as

$$\text{is_prime}(n) = \forall k \in [2, n-1] : n \neq 0 \pmod k \quad . \quad (1.1)$$

This formula has two clearly separated parts: first, a range of integers from 2 to $n-1$, and second, a requirement that all these integers should satisfy a given condition, $n \neq 0 \pmod k$.

This mathematical expression is translated into the Scala code as

```
def is_prime(n: Int) =
  (2 to n-1).forall(k => n % k != 0)
```

In this code, the two parts of the mathematical formula are implemented in a way that is closely similar to the mathematical notation, except for the arrow after k .

We can now apply the function `is_prime` to some integer values:

```
scala> is_prime(12)
res3: Boolean = false

scala> is_prime(13)
res4: Boolean = true
```

As we can see from the output above, the function `is_prime` returns a value of type `Boolean`. Therefore, the function `is_prime` has type `Int => Boolean`.

A function that returns a `Boolean` value is called a **predicate**.

In Scala, it is optional—but strongly recommended—to specify the return type of named functions. The required syntax looks like this,

1 Values, types, expressions, functions

```
def is_prime(n: Int): Boolean =  
  (2 to n-1).forall(k => n % k != 0)
```

However, we do not need to specify the type `Int` for the argument `k` of the nameless function `k => n % k != 0`. This is because the Scala compiler knows that `k` is going to iterate over the *integer* elements of the range `(2 to n-1)`, which effectively forces `k` to be of type `Int`.

1.1.3 Nameless functions and bound variables

The Scala code for `is_prime` differs from the mathematical formula (1.1) in two ways.

One difference is that the interval $[2, n - 1]$ comes first in the Scala expression. To understand this, look at the ways Scala allows programmers to define syntax.

The Scala syntax such as `(2 to n-1).forall(k => ...)` means to apply a function called `forall` to *two* arguments: the first argument is the range `(2 to n-1)`, and the second argument is the nameless function `(k => ...)`. In Scala, the **infix** syntax `x.f(z)`, or equivalently `x f z`, means that a function `f` is applied to its *two* arguments, `x` and `z`. In the ordinary mathematical notation, this would be $f(x, z)$. Infix notation is often easier to read and is also widely used in mathematics, for instance when we write $x + y$ rather than something like *plus*(x, y).

A single-argument function could be also defined with infix notation, and then the syntax is `x.f`, as in the expression `(1 to n).product` we have seen before.

The infix methods `.product` and `.forall` are already provided in the Scala standard library, so it is natural to use them. If we want to avoid the infix syntax, we could define a function `for_all` with two arguments and write code like this,

```
for_all(2 to n-1, k => n % k != 0)
```

This would have brought the syntax somewhat closer to the formula (1.1).

However, there still remains the second difference: The symbol k is used as an *argument* of a nameless function `(k => n % k != 0)` in the Scala code, — while the mathematical notation, such as

$$\forall k \in [2, n - 1] : n \neq 0 \pmod k ,$$

does not seem to involve any nameless functions. Instead, the mathematical formula defines the symbol k that “goes over a certain set,” as we might say. The symbol k is then used for writing the predicate $n \neq 0 \pmod k$.

1.1 Translating mathematics into code

However, let us investigate the role of the symbol k more closely.

The symbol k is a mathematical variable that is actually defined *only inside* the expression “ $\forall k : \dots$ ” and makes no sense outside that expression. This becomes clear by looking at Eq. (1.1): The variable k is not present in the left-hand side and could not possibly be used there. The name “ k ” is defined only in the right-hand side, where it is first mentioned as the arbitrary element $k \in [2, n - 1]$ and then used in the sub-expression “ $\dots \bmod k$ ”.

So, the mathematical notation

$$\forall k \in [2, n - 1] : n \neq 0 \bmod k$$

gives two pieces of information: first, we are examining all values from the given range; second, we chose the name k for the values from the given range, and for each of those k we need to evaluate the expression $n \neq 0 \bmod k$, which is a certain given *function of k* that returns a `Boolean` value. Translating the mathematical notation into code, it is therefore natural to use a nameless function of k ,

$$k \Rightarrow n \neq 0 \bmod k \quad ,$$

and to write Scala code that applies this nameless function to each element of the range $[2, n - 1]$ and then requires that all result values be `true`:

```
(2 to n-1).forall(k => n % k != 0)
```

Just as the mathematical notation defines the variable k only in the right-hand side of Eq. (1.1), the argument `k` of the nameless Scala function `k => n % k != 0` is defined only within that function’s body and cannot be used in any code outside the expression `n % k != 0`.

Variables that are defined only inside an expression and are invisible outside are called **bound variables**, or “variables bound in an expression”. Variables that are used in an expression but are defined outside it are called **free variables**, or “variables occurring free in an expression”. These concepts apply equally well to mathematical formulas and to Scala code. For example, in the mathematical expression $k \Rightarrow n \neq 0 \bmod k$ (which is a nameless function), the variable k is bound (because it is named and defined only within that expression) but the variable n is free (it is defined outside that expression).

The main difference between free and bound variables is that bound variables can be *locally renamed* at will, unlike free variables. To see this, consider that we could rename k to x and write instead of Eq. (1.1) an equivalent definition

$$\text{is_prime}(n) = \forall x \in [2, n - 1] : n \neq 0 \bmod x \quad ,$$

1 Values, types, expressions, functions

or in Scala code,

```
(2 to n-1).forall(x => n % x != 0)
```

In the nameless function $k \Rightarrow n \% k \neq 0$, the argument k may be renamed to x or to anything else, without changing the value of the entire program. No code outside this expression needs to be changed after renaming k to x . But the value n is defined outside and thus cannot be renamed locally (i.e. only within the sub-expression). If, for any reason, we wanted to rename n in the sub-expression $k \Rightarrow n \% k \neq 0$, we would also need to change every place in the code that defines and uses n *outside* that expression, or else the program would become incorrect.

Mathematical formulas have bound variables in various constructions such as $\forall k : p(k)$, $\exists k : p(k)$, $\sum_{k=a}^b f(k)$, $\int_0^1 k^2 dk$, $\lim_{n \rightarrow \infty} f(n)$, $\operatorname{argmax}_k f(k)$, and so on. When we translate mathematical expressions into code, we need to recognize the presence of bound variables, which the mathematical notation does not make quite so explicit. For each bound variable, we need to create a nameless function whose argument is that variable, e.g. $k \Rightarrow p(k)$ or $k \Rightarrow f(k)$ for the examples just shown. Only then will our code correctly reproduce the behavior of bound variables in mathematical expressions.

As an example, the mathematical formula

$$\forall k \in [1, n] : p(k) \quad ,$$

has a bound variable k and is translated into Scala code as

```
(1 to n).forall(k => p(k))
```

At this point we can apply the following simplification trick to this code. The nameless function $k \Rightarrow p(k)$ does exactly the same thing as the (named) function p : It takes an argument, which we may call k , and returns $p(k)$. So, we can simplify the Scala code above to

```
(1 to n).forall(p)
```

The simplification of $x \Rightarrow f(x)$ to just f is always possible for functions f of a single argument.³

³Certain features of Scala allow programmers to write code that looks like $\mathfrak{f}(x)$ but actually involves additional implicit or default arguments of the function \mathfrak{f} , or an implicit conversion for its argument x . In those cases, replacing the code $x \Rightarrow \mathfrak{f}(x)$ by just \mathfrak{f} may fail to compile in Scala. But these complications do not arise when working with simple functions.

1.1.4 Aggregating data from lists

Consider the task of finding how many even numbers there are in a given list L of integers. For example, the list $[1, 2, 3, 4, 5]$ contains *two* even numbers: 2 and 4.

A mathematical formula for this task can be written like this,

$$\text{count_even}(L) = \sum_{k \in L} \text{is_even}(k) \quad ,$$

$$\text{is_even}(k) = \begin{cases} 1 & \text{if } k \equiv 0 \pmod{2} \\ 0 & \text{otherwise} \end{cases} \quad .$$

Here we defined a helper function `is_even` in order to write more easily a formula for `count_even`. In mathematics, complicated formulas are often split into simpler parts by defining helper expressions.

We can write the Scala code similarly. We first define the helper function `is_even`; the Scala code can be written in the style quite similar to the mathematical formula:

```
def is_even(k: Int): Int = (k % 2) match {
  case 1 => 1
  case _ => 0
}
```

For such a simple computation, we could also write a shorter code and use a nameless function,

```
val is_even = (k: Int) => if (k % 2 == 0) 1 else 0
```

Given this function, we now need to translate into Scala code the expression $\sum_{k \in L} \text{is_even}(k)$. We can represent the list L using the data type `List[Int]` from the Scala standard library.

To compute $\sum_{k \in L} \text{is_even}(k)$, we must apply the function `is_even` to each element of the list L , which will produce a list of some (integer) results, and then we will need to add all those results together. It is convenient to perform these two steps separately. This can be done with the functions `.map` and `.sum`, defined in the Scala standard library as infix methods for the data type `List`.

The method `.sum` is similar to `.product` and is defined for any `List` of numerical types (`Int`, `Float`, `Double`, etc.). It computes the sum of all numbers in the list:

```
scala> List(1, 2, 3).sum
```

1 Values, types, expressions, functions

```
res0: Int = 6
```

The method `.map` needs more explanation. This method takes a *function* as its second argument, applies that function to each element of the list, and puts all the results into a *new* list, which is then returned as the result value:

```
scala> List(1, 2, 3).map(x => x*x + 100*x)
res1: List[Int] = List(101, 204, 309)
```

In this example, we used the nameless function $x \Rightarrow x^2 + 100x$ as the argument of `.map`. This function is repeatedly applied by `.map` to transform each of the values from a given list, creating a new list as a result.

It is equally possible to define the transforming function separately, give it a name, and then pass it as the argument to `.map`:

```
scala> def func1(x: Int): Int = x*x + 100*x
func1: (x: Int)Int

scala> List(1, 2, 3).map(func1)
res2: List[Int] = List(101, 204, 309)
```

Usually, short and simple functions are defined inline, while longer functions are given a name and defined separately.

An infix method, such as `.map`, can be also used with a “dot-less” syntax:

```
scala> List(1, 2, 3) map func1
res3: List[Int] = List(101, 204, 309)
```

If the transforming function `func1` is used only once, and especially for a simple operation such as $x \Rightarrow x^2 + 100x$, it is easier to work with a nameless function. We can now combine the methods `.map` and `.sum` to define `count_even`:

```
def count_even(s: List[Int]) = s.map(is_even).sum
```

This code can be also written using a nameless function instead of `is_even`:

```
def count_even(s: List[Int]): Int =
  s
    .map { k => if (k % 2 == 0) 1 else 0 }
    .sum
```

1.1 Translating mathematics into code

It is customary in Scala to use infix methods when chaining several operations. For instance `s.map(...).sum` means first apply `s.map(...)`, which returns a *new* list, and then apply `.sum` to that list. To make the code more readable, I put each of the chained methods on a new line.

To test this code, let us run it in the Scala interpreter. In order to let the interpreter work correctly with code entered line by line, the dot character needs to be at the end of the line. The interpreter will automatically insert the visual continuation characters. (In a compiled code, the dots can be at the beginning of the lines since the compiler reads the entire code at once.)

```
scala> def count_even(s: List[Int]): Int =  
    |   s .  
    |     map { k => if (k % 2 == 0) 1 else 0 } .  
    |     sum  
count_even: (s: List[Int])Int  
  
scala> count_even(List(1,2,3,4,5))  
res0: Int = 2  
  
scala> count_even( List(1,2,3,4,5).map(x => x * 2) )  
res1: Int = 5
```

Note that the Scala interpreter prints the types differently for functions defined using `def`. It prints `(s: List[Int])Int` for the function type that one would normally write as `List[Int] => Int`.

1.1.5 Filtering

In addition to the methods `.sum`, `.product`, `.map`, `.forall` that we have already seen, the Scala standard library defines many other useful methods. We will now take a look at using the methods `.max`, `.min`, `.exists`, `.size`, `.filter`, and `.takeWhile`.

The methods `.max`, `.min`, and `.size` are self-explanatory:

```
scala> List(10, 20, 30).max  
res2: Int = 30  
  
scala> List(10, 20, 30).min  
res3: Int = 10
```

1 Values, types, expressions, functions

```
scala> List(10, 20, 30).size  
res4: Int = 3
```

The methods `.forall`, `.exists`, `.filter`, and `.takeWhile` require a predicate as an argument. The `.forall` method returns `true` iff the predicate is true on all values in the list; the `.exists` method returns `true` iff the predicate holds (returns `true`) for at least one value in the list. These methods can be written as mathematical formulas like this:

$$\begin{aligned}\text{forall}(S, p) &= \forall k \in S : p(k) = \text{true} \\ \text{exists}(S, p) &= \exists k \in S : p(k) = \text{true}\end{aligned}$$

However, mathematical notation does not exist for operations such as “finding elements in a list”, so we will focus on the Scala syntax for these operations.

The `.filter` method returns a *new list* that contains only the values for which the predicate returns `true`:

```
scala> List(1, 2, 3, 4, 5).filter(k => k % 3 != 0)  
res5: List[Int] = List(1, 2, 4, 5)
```

The `.takeWhile` method returns a *new list* that contains the initial subset of values from the original list, which are taken until the predicate first returns `false`:

```
scala> List(1, 2, 3, 4, 5).takeWhile(k => k % 3 != 0)  
res6: List[Int] = List(1, 2)
```

In all these cases, the predicate’s argument `k` must be of the same type as the elements in the list. In the examples shown above, the elements are integers (i.e. the lists have type `List[Int]`), therefore `k` must be of type `Int`.

The methods `.max`, `.min`, `.sum`, and `.product` are defined on lists of *numeric types*, such as `Int`, `Double`, and `Long`. The other methods are defined on lists of all types.

Using these methods, we can solve many problems that involve transforming and aggregating data stored in lists (as well as in arrays, sets, or other similar data structures). By **transformation** I here mean a function from a list of values to another list of values; examples of transformation functions are `.filter` and `.map`. By **aggregation** I mean a function from a list of values to a *single* value; examples of aggregation functions are `.max` and `.sum`.

Writing programs by chaining together various operations of transformation and aggregation is known as programming in the **map/reduce style**.

1.2 Worked examples: transformation and aggregation

1. Improve the code for `is_prime` by limiting the search to $k^2 \leq n$,

$$\text{is_prime}(n) = \forall k \in [2, n-1] \text{ such that } k^2 \leq n : n \neq 0 \pmod k .$$

```
def is_prime(n: Int): Boolean =
  (2 to n-1)
    .filter(k => k*k <= n)
    .forall(k => n % k != 0)
```

2. Compute $\prod_{k \in [1,10]} |\sin(k+2)|$.

```
(1 to 10)
  .map(k => math.abs(math.sin(k + 2)))
  .product
```

3. Compute $\sum_{k \in [1,10]; \cos k > 0} \sqrt{\cos k}$.

```
(1 to 10)
  .filter(k => math.cos(k) > 0)
  .map(k => math.sqrt(math.cos(k)))
  .sum
```

It is safe to compute $\sqrt{\cos k}$, because we have first filtered the list by keeping only values k for which $\cos k > 0$:

```
scala> (1 to 10).toList.
  filter(k => math.cos(k) > 0).map(x => math.cos(x))
res0: List[Double] = List(0.5403023058681398,
  0.28366218546322625, 0.9601702866503661, 0.7539022543433046)
```

4. Compute the average of a list of numbers of type `Double` (assuming that the list is not empty).

```
scala> def average(s: List[Double]): Double = s.sum / s.size
average: (s: List[Double])Double
```

1 Values, types, expressions, functions

```
scala> average(List(1.0, 2.0, 3.0))
res4: Double = 2.0
```

5. Given n , compute the **Wallis product** truncated up to $\frac{2n}{2n+1}$:

$$\text{wallis}(n) = \frac{2}{1} \frac{2}{3} \frac{4}{5} \frac{6}{7} \dots \frac{2n}{2n+1} .$$

We use the method `.toDouble` to create `Double` numbers from integers.

```
def wallis_frac(i: Int): Double =
  (2*i).toDouble / ((2*i - 1)*(2*i)/(2*i + 1))

def wallis(n: Int) = (1 to n).map(wallis_frac).product

scala> math.cos(wallis(10000)) // Should be close to 0.
res5: Double = 3.9267453954401036E-5
scala> math.cos(wallis(100000)) // Should be even closer to 0.
res6: Double = 3.926966362362075E-6
```

The limit of the Wallis product is $\frac{\pi}{2}$, so the cosine of `wallis(n)` tends to zero in the limit of large n .

6. Another well-known series related to π is

$$\sum_{k=1}^{\infty} \frac{1}{k^2} = \frac{\pi^2}{6} .$$

Define a function of n that computes a partial sum of this series until $k = n$. Compute the result for a large value of n and compare with the limit value.

```
def euler_series(n: Int): Double =
  (1 to n).map(k => 1.0/k/k).sum

scala> euler_series(100000)
res10: Double = 1.6449240668982423

scala> val pi = 4*math.atan(1)
pi: Double = 3.141592653589793
```

1.2 Worked examples: transformation and aggregation

```
scala> pi*pi/6
res9: Double = 1.6449340668482264
```

7. Check numerically the infinite product formula

$$\prod_{k=1}^{\infty} \left(1 - \frac{x^2}{k^2}\right) = \frac{\sin \pi x}{\pi x} \quad .$$

We will compute the product up to $k = n$ for a fixed $x = 0.1$ and a large value of n , say $n = 10^5$, and compare with the right-hand side:

```
def sine_product(n: Int, x: Double): Double =
  (1 to n).map(k => 1.0 - x*x/k/k).product

scala> sine_product(100000, 0.1)
res11: Double = 0.9836317414461351

scala> math.sin(pi*0.1)/pi/0.1
res12: Double = 0.9836316430834658
```

8. Define a function p that takes a list of integers and a function $f: \text{Int} \Rightarrow \text{Int}$, and returns the largest value of $f(x)$ among all x in the list.

```
def p(s: List[Int], f: Int => Int): Int = s.map(f).max
```

9. Given a list of lists, $s: \text{List}[\text{List}[\text{Int}]]$, select the inner lists of size at least 3. The result must be again of type $\text{List}[\text{List}[\text{Int}]]$.

To “select the inner lists” satisfying a condition means to compute a *new* list containing only the desired inner lists. We use `.filter` on the outer list s . The predicate for the filter is a function that takes an inner list and returns `true` if the size of that list is at least 3. Write the predicate as a nameless function, $t \Rightarrow t.\text{size} \geq 3$. Then the solution is

```
def f(s: List[List[Int]]): List[List[Int]] =
  s.filter(t => t.size >= 3)

scala> f(List( List(1,2), List(1,2,3), List(1,2,3,4) ))
res13: List[List[Int]] = List(List(1, 2, 3), List(1, 2, 3, 4))
```

1 Values, types, expressions, functions

The predicate in the argument of `.filter` is a nameless function $t \Rightarrow t.size \geq 3$ whose argument t is of type `List[Int]`. The Scala compiler deduces the type of t from the code; no other type would work with the way we use `.filter` on a list of lists of integers.

10. Find all integers $k \in [1, 10]$ such that there are at least three different integers j , where $1 \leq j \leq k$, each j satisfying the condition $j^2 > 2k$.

```
scala> (1 to 10).toList.filter(k => (1 to k).  
    |   filter(j => j*j > 2*k).size >= 3)  
res14: List[Int] = List(6, 7, 8, 9, 10)
```

The argument of the outer `.filter` is a nameless function that itself uses another `.filter`. The inner expression,

```
(1 to k).filter(j => j*j > 2*k).size >= 3
```

computes the list of j 's that satisfy the condition $j^2 > 2k$, and then computes the length of that list and imposes the requirement that there should be at least 3 such values of j . We can see how the Scala code closely follows the mathematical formulation of the problem.

1.3 Summary

The following table translates mathematical formulas into code.

<i>Mathematical notation</i>	<i>Scala code</i>
$x \mapsto \sqrt{x^2 + 1}$	<code>x => math.sqrt(x * x + 1)</code>
list $[1, 2, \dots, n]$	<code>(1 to n)</code>
list $[f(1), \dots, f(n)]$	<code>(1 to n).map(k => f(k))</code>
$\sum_{k=1}^n k^2$	<code>(1 to n).map(k => k*k).sum</code>
$\prod_{k=1}^n f(k)$	<code>(1 to n).map(f).product</code>
$\forall k$ such that $1 \leq k \leq n : p(k)$ holds	<code>(1 to n).forall(k => p(k))</code>
$\exists k, 1 \leq k \leq n$ such that $p(k)$ holds	<code>(1 to n).exists(k => p(k))</code>
$\sum_{k \in S: p(k) \text{ holds}} f(k)$	<code>s.filter(p).map(f).sum</code>

What problems can we solve with this knowledge?

- Compute mathematical expressions involving sums, products, and quantifiers, based on integer ranges, such as $\sum_{k=1}^n f(k)$ etc.
- Transform and aggregate data from lists using `.map`, `.filter`, `.sum`, and other methods from the Scala standard library.

What are examples of problems that are not solvable with these tools?

- Example 1: Compute the smallest $n \geq 1$ such that

$$f(f(f(\dots f(0)\dots)) > 1000 \quad ,$$

where the given function f is applied n times.

- Example 2: Given a list s of numbers, compute the list r of running averages: $r_n = \sum_{k=0}^n s_k$.
- Example 3: Perform binary search over a sorted list of integers.

These problems are not yet solvable because the required formulas involve *mathematical induction*, which we cannot yet translate into code in the general case.

Library functions we have seen so far, such as `.map` and `.filter`, implement a restricted class of iterative operations on lists: namely, operations that process each element of a given list independently and accumulate results. For instance, when computing $s.\text{map}(f)$, the number of function applications is given by the size of the initial list. However, Example 1 requires applying a function f repeatedly until a given condition holds—that is, repeating for an *initially unknown* number of times. So it is impossible to write an expression containing `.map`, `.filter`, `.takeWhile`, etc., that solves Example 1. We could use mathematical induction to write the solution of Example 1 as a formula, but we have not yet seen how to implement it in Scala code.

Similarly, Example 2 defines a new list r from s by induction,

$$r_0 = s_0; \quad r_i = s_i + r_{i-1}, \forall i > 0 \quad .$$

However, operations such as `.map` and `.filter` cannot compute r_i depending on the value of r_{i-1} .

Example 3 defines the search result by induction: the list is split in half, and search is performed by inductive hypothesis in the half that contains the required value. This computation requires an initially unknown number of steps.

Chapter 2 explains how to solve these problems by translating mathematical induction into code using recursion.

1.4 Exercises

1. Define a function of type `List[Double] => List[Double]` that “normalizes” the list: finds the element having the largest absolute value and, if that value is nonzero, divides all elements by that factor and returns a new list; otherwise returns the original list.
2. **Machin’s formula** computes π more efficiently than the Wallis fraction:

$$\frac{\pi}{4} = 4 \arctan \frac{1}{5} - \arctan \frac{1}{239} \quad ,$$

$$\arctan \frac{1}{n} = \frac{1}{n} - \frac{1}{3} \frac{1}{n^3} + \frac{1}{5} \frac{1}{n^5} - \dots$$

Implement a function that computes the series for $\arctan \frac{1}{n}$ up to a given number of terms, and compute an approximation of π using this formula. Show that about 12 terms of the series are already sufficient for a full-precision `Double` approximation of π .

3. Using the function `is_prime`, check numerically the **Euler product formula** for the Riemann zeta function $\zeta(s)$ at $s = 4$, since **it is known that** $\zeta(4) = \frac{\pi^4}{90}$,

$$\prod_{k \in \mathbb{N}; k \text{ is prime}} \frac{1}{1 - p^{-4}} = \frac{\pi^4}{90} \quad .$$

4. Define a function `add_20` of type `List[List[Int]] => List[List[Int]]` that adds 20 to every element of every inner list. A test:

```
scala> add_20( List( List(1), List(2, 3) ) )
res0: List[List[Int]] = List(List(21), List(22, 23))
```

5. An integer n is called a “3-factor” if it is divisible by only three different integers j such that $2 \leq j < n$. Compute the set of all “3-factor” integers n among $n \in [1, \dots, 1000]$.
6. Given a function `f: Int => Boolean`, an integer n is called a “3- f ” if there are only three different integers $j \in [1, \dots, n]$ such that `f(j)` returns `true`. Define a function that takes `f` as an argument and returns a sequence of all “3- f ” integers among $n \in [1, \dots, 1000]$. What is the type of that function? Implement Exercise 5 using that function.

7. Define a function `sel_100` of type `List[List[Int]] => List[List[Int]]` that selects only those inner lists whose largest value is at least 100. A test:

```
scala> sel_100( List( List(0, 100), List(60, 80), List(1000) ) )
res0: List[List[Int]] = List(List(0, 100), List(1000))
```

1.5 Discussion

1.5.1 Functional programming as a paradigm

Functional programming (FP) is a **paradigm** of programming,—that is, an approach that guides programmers to write code in specific ways, for a wide range of programming tasks.

The main principle of FP is to write code *as a mathematical expression or formula*. This approach allows programmers to derive code through logical reasoning rather than through guessing, – similarly to how books on mathematics reason about mathematical formulas and derive results systematically, without guessing or “debugging.” Similarly to mathematicians and scientists who reason about formulas, functional programmers can *reason about code* systematically and logically, based on rigorous principles. This is possible only because code is written as a mathematical formula.

Mathematical intuition is backed by the vast experience accumulated while working with data over thousands of years of human history. It took centuries to invent flexible and powerful notation such as $\forall k \in S : p(k)$ and to develop the corresponding rules of reasoning. Functional programmers are fortunate to have at their disposal such a superior reasoning tool.

As we have seen, the Scala code for certain computational tasks corresponds quite closely to mathematical formulas. (Scala conventions and syntax, of course, require programmers to spell out certain things that the mathematical notation leaves out.) Just as in mathematics, large code expressions may be split into parts in a suitable way, so that the parts can be easily reused, flexibly composed together, and written independently from each other. The FP community has developed a toolkit of functions (such as `.map`, `.filter`, etc.) that proved especially useful in real-life programming, although many of them are not standard in mathematical literature.

Mastering FP involves practicing to reason about programs as formulas, building up the specific kind of applied mathematical intuition, familiarizing oneself

with concepts adapted to programming needs, and learning how to translate the mathematics into code in various cases. The FP community has discovered a number of specific design patterns, founded on mathematical principles but driven by practical necessities of programming rather than by the needs of academic mathematics. This book explains the required mathematical principles in detail, developing them through intuition and practical coding tasks.

1.5.2 Functional programming languages

It is possible to apply the FP paradigm while writing code in any programming language. However, some languages lack certain features that make FP techniques much easier to use in practice. For example, in a language such as Python or Ruby, one can productively apply only a small number of the idioms of FP, such as the map/reduce operations. More advanced FP constructions are impractical in these languages because the corresponding code becomes too complicated to read, and mistakes are too easy to make, which negates the advantage of easier reasoning about the FP code.

Some programming languages, such as Haskell and OCaml, were designed specifically for advanced use in the FP paradigm. Other languages, such as ML, F#, Scala, Swift, Elm, PureScript, and Rust, had different design goals but still support enough FP features to be considered FP languages. I will be using Scala in this book, but exactly the same constructions could be implemented in other FP languages in a similar way. At the level of detail needed in this book, the differences between languages such as ML, OCaml, Haskell, F#, Scala, Swift, Elm, PureScript, or Rust will not play a significant role.

1.5.3 The mathematical meaning of variables

The usage of variables in functional programming closely corresponds to how mathematical literature uses variables. In mathematics, **variables** are used first of all as arguments of functions; e.g. the formula

$$f(x) = x^2 + x$$

contains the variable x and defines a function f that takes a number x as its argument (to be definite, let us assume that x is an integer) and computes the value $x^2 + x$. The body of the function is the expression $x^2 + x$.

Mathematics has the convention that a variable, such as x , *does not change* its value within a formula. Indeed, there is no mathematical notation even to talk

about “modifying” the value of x *inside* the formula $x^2 + x$. It would be quite confusing if a mathematics textbook said “before adding the last x in the formula $x^2 + x$, we modify x by adding 4 to it”. If the “last x ” in $x^2 + x$ needs to have a 4 added to it, a mathematics textbook will just write the formula $x^2 + x + 4$.

Arguments of nameless functions are also immutable. Consider, for example,

$$f(n) = \sum_{k=0}^n k^2 + k \quad .$$

Here, n is the argument of the function f , while k is the argument of the nameless function $k \Rightarrow k^2 + k$. Neither n nor k can be “modified” in any sense within the expressions where they are used. The symbols k and n stand for some integer values, and these values are immutable. Indeed, it is meaningless to say that we want to “modify the value 4”. In the same way, we cannot modify k .

So, a variable in mathematics does not actually vary *within the expression* where it is defined; in that expression, a variable is essentially a *named constant* value. A function f can be applied to different values x , to compute a different result $f(x)$ each time. However, a given value of x will remain unmodified within the body of the function f when it is computed.

Functional programming adopts this convention from mathematics: variables are immutable named constants. (Scala also has *mutable* variables, but we will not need to consider them in this book.)

In Scala, function arguments are immutable within the function body:

```
def f(x: Int) = x*x + x // Can't modify x here.
```

The *type* of each mathematical variable (say, integer, string, etc.) is also fixed in advance. In mathematics, each variable is a value from a specific set, known in advance (the set of all integers, the set of all strings, etc.). Mathematical formulas such as $x^2 + x$ do not express any “checking” that x is indeed an integer and not, say, a string, before starting to evaluate $x^2 + x$.

Functional programming adopts the same view: Each argument of each function must have a **type**, which represents *the set of possible allowed values* for that function argument. The programming language’s compiler will automatically check all types of all arguments. A program that calls functions on arguments of incorrect types will not compile.

The second usage of **variables** in mathematics is to denote expressions that will be reused. For example, one writes: let $z = \dots$ and now compute $\cos z + \cos 2z + \cos 3z$. Again, the variable z remains immutable, and its type remains fixed.

1 Values, types, expressions, functions

In Scala, this construction (defining an expression to be reused later) is written with the “`val`” syntax. Each variable defined using “`val`” is a named constant, and its type and value are fixed at the time of definition. Types for “`val`”s are optional in Scala, for instance we could write

```
val x: Int = 123
```

or more concisely,

```
val x = 123
```

because it is clear that this `x` is of type `Int`. However, when types are complicated, it helps to write them out. The compiler will check that the types match correctly everywhere and give an error message if we write

```
scala> val x: Int = "123" // A String instead of an Int.
<console>:11: error: type mismatch;
   found   : String("123")
   required: Int
       val x: Int = "123" // A String instead of an Int.
               ^
```

1.5.4 Iteration without loops

Another distinctive feature of functional programming is the absence of explicit loops or repetition.

Iterative computations are ubiquitous in mathematics; one can see formulas such as

$$\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n s_i s_j - \left(\frac{1}{n} \sum_{i=1}^n s_i \right)^2 .$$

To compute these expressions, we need to iterate over various values of i, j , etc. And yet, no mathematics textbook ever mentions “loops” or says “now repeat this formula ten times”. Indeed, it would be pointless to evaluate a formula such as $x^2 + x$ ten times, or to “repeat” an equation such as

$$(x - 1)(x^2 + x + 1) = x^3 - 1 .$$

Instead of loops, mathematicians write *expressions* such as $\sum_{i=1}^n s_i$, which denote repetitive computations. Mathematically, such expressions are defined using

mathematical induction. The FP paradigm has developed rich tools for translating mathematical induction into code. In this chapter, we have seen methods such as `.map`, `.filter`, and `.sum`, which implement certain kinds of iterative computations. These and other operations can be combined in very flexible ways, which allows programmers to write iterative code *without loops*.

The programmer can avoid writing loops because the iteration is delegated to the library functions `.map`, `.filter`, `.sum`, and so on. It is the job of the library and the compiler to translate these functions into machine code that likely *will* contain loops; but the functional programmer does not need to see that code or to reason about it.

1.5.5 Nameless functions in mathematical notation

Functions in mathematics are mappings from one set to another. A function does not necessarily *need* a name; the mapping just needs to be defined. However, nameless functions have not been widely used in the conventional mathematical notation. It turns out that nameless functions are quite important in functional programming because, in particular, they allow programmers to write code more concisely and use a straightforward, consistent syntax.

Nameless functions have the property that their bound variables are invisible outside their scope. This property is directly reflected by the prevailing mathematical conventions. Compare the formulas

$$f(x) = \int_0^x \frac{dx}{1+x} \quad ; \quad f(x) = \int_0^x \frac{dz}{1+z} \quad .$$

The mathematical convention is that these formulas define the same function f , and that one may rename the integration variable at will.

In programming, the only situation when a variable “may be renamed at will” is when the variable represents an argument of a function. It follows that the notations $\frac{dx}{1+x}$ and $\frac{dz}{1+z}$ correspond to a nameless function whose argument was renamed from x to z . In FP notation, this nameless function would be denoted as $z \Rightarrow \frac{1}{1+z}$, and the integral rewritten as code such as

$$\text{integration}(0, x, g) \text{ where } g = \left(z \Rightarrow \frac{1}{1+z} \right) \quad .$$

1 Values, types, expressions, functions

Now consider the traditional mathematical notations for summation, e.g.

$$\sum_{k=0}^x \frac{1}{1+k} \quad .$$

In sums, the bound variable k is introduced under the \sum symbol; but in integrals, the bound variable follows the special symbol “ d ”. This notational inconsistency could be removed if we were to use nameless functions explicitly, for example:

$$\begin{aligned} \sum_0^x k &\Rightarrow \frac{1}{1+k} \text{ instead of } \sum_{k=0}^x \frac{1}{1+k} \quad , \\ \int_0^x z &\Rightarrow \frac{1}{1+z} \text{ instead of } \int_0^x \frac{dz}{1+z} \quad . \end{aligned}$$

In this notation, the new summation symbol \sum_0^x does not mention the name “ k ” but takes a function as an argument. Similarly, the new integration symbol \int_0^x does not mention “ z ” and does not use the special symbol “ d ” but now takes a function as an argument. Written in this way, the operations of summation and integration become *functions* that take a function as argument. The above summation may be written in a consistent and straightforward manner as a function:

$$\text{summation}(0, x, f) \text{ where } f = (y \Rightarrow \frac{1}{1+y}) \quad .$$

We could implement `summation(a,b,g)` as

```
def summation(a: Int, b: Int, g: Int => Double) =  
  (a to b).map(g).sum  
summation: (a: Int, b: Int, g: Int => Double)Double  
  
scala> summation(1, 10, x => math.sqrt(x))  
res0: Double = 22.4682781862041
```

Numerical integration requires longer code, since the formulas are more com-

plicated. For instance, **Simpson's rule** can be written as

$$\begin{aligned} \text{integration}(a, b, g) &= \frac{\delta}{3} (g(a) + g(b) + 4s_1 + 2s_2) \quad , \\ n &= 2 \left\lceil \frac{b-a}{\varepsilon} \right\rceil, \quad \delta = \frac{b-a}{n} \quad , \\ s_1 &= \sum_{i=1,3,\dots,n-1} g(a+i\delta) \quad , \\ s_2 &= \sum_{i=2,4,\dots,n-2} g(a+i\delta) \quad . \end{aligned}$$

A straightforward translation of this formula into Scala is

```
def integration(a: Double, b: Double,
  g: Double => Double, eps: Double) = {
  // Define some helper values and functions.
  val n: Int = (math.round((b-a)/eps/2)*2).toInt
  val delta_x = (b - a) / n
  val g_i = (i: Int) => g(a + i*delta_x)
  val s1 = (1 to (n-1) by 2).map(g_i).sum
  val s2 = (2 to (n-2) by 2).map(g_i).sum
  // Compute the final result.
  delta_x / 3 * (g(a) + g(b) + 4*s1 + 2*s2)
}

scala> integration(0, 2, x => x*x*x, 0.001) // Exact answer is 4
res0: Double = 4.000000000000003

scala> integration(0, 7, x => x*x*x*x*x*x, 0.001) // Exact answer is
117649
res1: Double = 117649.00000000023
```

The entire code is one large *expression*, with a few sub-expressions defined for convenience as a few helper values and helper functions. In other words, this code is written in the FP paradigm.

1.5.6 Named and nameless expressions and their uses

It is a significant advantage if a programming language supports unnamed (or “nameless”) expressions. To see this, consider a familiar situation where we take the absence of names for granted.

1 Values, types, expressions, functions

In most programming languages today, we can directly write arithmetical expressions such as $(x+123)*y/(2+x)$. Here, x and y are variables with names. Note, however, that the entire expression does not need to have a name. Parts of that expression (such as $x+123$ or $2+x$) also do not need to have separate names. It would be quite inconvenient if we *needed* to assign a name separately to each sub-expression. The code for $(x+123)*y/(2+x)$ would then look like this:

```
r1 = 123
r2 = x + r1
r3 = r2 * y
r4 = 2
r5 = r4 + x
r6 = r3 / r5
return r6
```

This style of programming resembles assembly languages, where *every* sub-expression – that is, every step of every calculation, – must be named separately (and, in the assembly languages, assigned a memory address or a CPU register).

Programmers gain productivity when their programming language supports expressions without names (which I call “nameless expressions”).

This is also common practice in mathematics; names are assigned when needed, but many expressions remain nameless.

It is similarly quite useful if data structures can be declared without a name. For instance, a **dictionary** (also called a “hashmap”) is declared in Scala as

```
Map("a" -> 1, "b" -> 2, "c" -> 3)
```

This is a nameless expression representing a dictionary. Without this construction, programmers have to write cumbersome, repetitive code that creates an initially empty named dictionary and then fills it step by step with values:

```
// Scala code for creating a dictionary:
val myMap = Map("a" -> 1, "b" -> 2, "c" -> 3)
// Java code:
Map<String, Int> myMap = new HashMap<String, Integer>() {{
    put("a", 1);
    put("b", 2);
    put("c", 3);
}}; // The shortest Java code for creating the same dictionary.
```

Nameless functions are useful for the same reason as nameless data values: they allow the programmer to build a larger program from simpler parts in a uniform and concise way.

1.5.7 Nameless functions: historical perspective

Nameless functions were first used in 1936 in a theoretical programming language called “ λ -calculus”. In that language,⁴ all functions are nameless and have a single argument. The letter λ is a syntax separator denoting function arguments in nameless functions. For example, the nameless function $x \Rightarrow x + 1$ could be written as $\lambda x.add\ x\ 1$ in λ -calculus, if it had a function *add* for adding integers (which it does not).

In most programming languages that were in use until around 1990, all functions required names. But by 2015, most languages added support for nameless functions, possibly because programming in the map/reduce style (which invites frequent use of nameless functions) turned out to be quite useful. Table 1.1 shows when nameless functions were introduced in each language.

What I call a “nameless function” is also elsewhere called anonymous function, function expression, function literal, closure, lambda function, lambda expression, or just a “lambda”. I use the term “nameless function” in this book because it is the most descriptive and unambiguous both in speech and in writing.

⁴Although called a “calculus,” it is in reality a (drastically simplified) programming language. It has nothing to do with “calculus” as known in mathematics, such as differential or integral calculus. Also, the letter λ has no particular significance; it plays a purely syntactic role in the λ -calculus. Practitioners of functional programming usually do not need to study any λ -calculus. All practically relevant knowledge related to λ -calculus is explained in Chapter 3 of this book.

1 Values, types, expressions, functions

Language	Year	Code for $k:\text{Int} \Rightarrow k + k$
λ -calculus	1936	$\lambda k. \text{add } k \ k$
typed λ -calculus	1940	$\lambda k : \text{int}. \text{add } k \ k$
LISP	1958	(lambda (k) (+ k k))
Standard ML	1973	fn (k: int) => k+k
OCaml	1985	fun (k: int) -> k+k
Haskell	1990	\ (k: Int) -> k+k
Oz	1991	fun {\$ K} K+K
Ruby	1993	lambda { k k+k }
R	1993	function(k) k+k
Python	1994	lambda k: k+k
JavaScript	1995	function(k) { return k+k; }
Scala	2003	(k: Int) => k+k
F#	2005	fun (k: int) -> k+k
C++ 11	2011	[] (int k) { return k+k; }
Go	2012	func(k int) { return k+k }
Kotlin	2012	{ k: Int -> k+k }
Swift	2014	{(k: int) -> int in return k+k}
Java 8	2014	(int k) -> k+k
Rust	2015	k: i32 k+k

Table 1.1: Nameless functions in various programming languages.

2 Mathematical induction

In this chapter, we will see more functionality that Scala gives us to treat collections in a functional way – that is, in a way that is more suitable for mathematical thinking. In particular, the Scala standard library has methods for performing quite general iterative computations – including those that represent mathematical quantities defined by induction. Translating mathematical induction into code is the main topic of this chapter.

Before we begin, we need to become fluent with some of the types used by the Scala collections library.

2.1 Tuple types

2.1.1 First examples

Many standard library methods in Scala require working with **tuple types**. What are tuple types? The easiest example of a tuple is a *pair* of values, for example, a pair of an integer and a string. The Scala syntax for this type of pair is

```
val a: (Int, String) = (123, "xyz")
```

The type expression `(Int, String)` denotes this tuple type.

We can also have a *triple* of values:

```
val b: (Boolean, Int, Int) = (true, 3, 4)
```

Pairs and triples are examples of tuples. A tuple can contain any fixed number of values, which I call **parts** of a tuple. The parts of a tuple can have different types, but the type of each part is fixed once and for all. It is a type error to use incorrect types in a tuple, or an incorrect number of parts of a tuple:

```
scala> val bad: (Int, String) = (1,2)
<console>:11: error: type mismatch;
found   : Int(2)
```

2 Mathematical induction

```
required: String
    val bad: (Int, String) = (1,2)
                                ^
scala> val bad: (Int, String) = (1,"a",3) <console>:11: error: type
mismatch; found   : (Int, String, Int) required: (Int, String)
    val bad: (Int, String) = (1,"a",3)
                                ^
```

Parts of a tuple can be accessed by number, starting from 1. The Scala syntax for tuple access methods looks like this,

```
scala> val a = (123, "xyz")
a: (Int, String) = (123,xyz)

scala> a._1
res0: Int = 123

scala> a._2
res1: String = xyz
```

It is a type error to access a part that does not exist:

```
scala> a._0
<console>:13: error: value _0 is not a member of (Int, String)
    a._0
      ^

scala> a._5
<console>:13: error: value _5 is not a member of (Int, String)
    a._5
      ^
```

These **type errors** are detected at compile time, – before any computations begin. Tuples can be **nested**: any part of a tuple can itself be of a tuple type.

```
scala> val c: (Boolean, (String, Int), Boolean) = (true, ("abc", 3),
false)
c: (Boolean, (String, Int), Boolean) = (true,(abc,3),false)

scala> c._1
res2: Boolean = true
```

```
scala> c._2
res3: (String, Int) = (abc,3)
```

To define functions whose arguments are tuples, use the syntax

```
def f(p: (Boolean, Int), q: Int): Boolean = p._1 && (p._2 > q)
```

The first argument, *p*, of this function, has a tuple type. The function body uses access methods to compute its result value. Note that the second part of the tuple *p* is of type `Int`, so it is valid to compare it with an integer *q*. It would be a type error to compare the *tuple* *p* with an integer using the expression *p* > *q*. It would be also a type error to apply *f* to an argument *p* that has a wrong type, e.g. type `(Int, Int)` instead of `(Boolean, Int)`.

2.1.2 Pattern matching on tuples

Instead of using access methods when working with tuples, it is more convenient to use **pattern matching**. Scala performs pattern matching in two situations:

- destructuring definition: `val pattern = ...`
- partial function: `case pattern => ...`

An example of a **destructuring definition** is

```
scala> val g = (1, 2, 3)
g: (Int, Int, Int) = (1,2,3)

scala> val (x, y, z) = g
x: Int = 1
y: Int = 2
z: Int = 3
```

The value *g* is a tuple of three integers. After defining *g*, we define the three variables *x*, *y*, *z* *at once* in a single `val` definition. We imagine that this definition “destructures” the data structure contained in *g* and decomposes it into three parts, and then assigns the names *x*, *y*, *z* to these parts. The types of each of the new values are also assigned automatically.

The left-hand side of the destructuring definition contains the tuple pattern `(x, y, z)` that looks like a tuple except that its parts are names *x*, *y*, *z* that are so

2 Mathematical induction

far undefined. These names are called **pattern variables**. The destructuring definition *matches* the three pattern variables to the value of *g*. If *g* does not contain a tuple with exactly three parts, the pattern-matching definition will fail. This computation is called **pattern matching**.

Pattern matching is often used in defining functions on tuples. An example of such a function is

```
def s(p: (Int, Int, Int)): Int = p match { case (x, y, z) => x + y + z }
scala> s(g)
res0: Int = 6
```

The **case expression** `case (x, y, z) => x + y + z` performs pattern matching on the tuple argument *p*. The pattern matching process will “destructure” (i.e. decompose) the tuple and try to match it to the given pattern `(x, y, z)`. In this pattern, *x*, *y*, *z* are as yet undefined new variables, – that is, they are pattern variables. Once the pattern matching succeeds, the pattern variables *x*, *y*, *z* are assigned their values, and the function body can proceed to perform its computation. In this example, *x*, *y*, *z* will be assigned values 1, 2, and 3, so the function returns 6 as its result value.

Pattern matching is especially convenient when working with nested tuples. Here is an example of such code:

```
def t(p: (Int, (String, Int))): String = p match {
  case (x, (str, y)) => str + (x + y).toString
}

scala> t((10, ("result is ", 2)))
res0: String = result is 12
```

The type structure of the argument is visually repeated in the pattern. It is easy to see that *x* and *y* become integers and *str* becomes a string after pattern matching. We could rewrite the same code using tuple access methods instead of pattern matching:

```
def t1(p: (Int, (String, Int))): String = {
  p._2._1 + (p._1 + p._2._2).toString
}
```

This code is more difficult to read and to maintain because, for example, it is not

immediately clear what `p._2._1` refers to.

2.1.3 Using tuples with collections

Tuples can be combined with any other types without restrictions. For instance, we can have a tuple of functions,

```
val q: (Int => Int, Int => Int) = (x => x + 1, x => x - 1)
```

We can have a list of tuples,

```
val r: List[(String, Int)] = List(
  ("apples", 3), ("oranges", 2), ("pears", 0))
```

We can have a tuple of lists of tuples of functions, or any other combination of types.

Here is an example of using the standard method `.map` to transform a list of tuples. The argument of `.map` must be a function taking a tuple as its argument. It is convenient to use pattern matching for writing such functions:

```
scala> val basket: List[(String, Int)] = List(
  ("apples", 3), ("pears", 2), ("lemons", 0))
basket: List[(String, Int)] = List((apples,3), (pears,2), (lemons,0))

scala> basket.map { case (fruit, count) => count *2 }
res0: List[Int] = List(6, 4, 0)

scala> basket.map { case (fruit, count) => count *2 }.sum
res1: Int = 10
```

In this way, we can use the standard methods such as `.map`, `.filter`, `.max`, `.sum` to manipulate sequences of tuples. Note how the chosen names “fruit”, “count” help us remember the meaning of the parts of tuples.

We can easily transform a list of tuples into a list of tuples of different type:

```
scala> basket.map { case (fruit, count) =>
  val isAcidic = fruit == "lemons"
  (fruit, isAcidic)
}
res2: List[(String, Boolean)] = List((apples,false), (pears,false),
  (lemons,true))
```

2 Mathematical induction

In the Scala syntax, a nameless function written with braces { ... } can define local values in its body. The return value of the function is the last expression written in the function body. In this example, the return value of the nameless function is the tuple (fruit, isAcidic).

2.1.4 Using dictionaries (Scala's `Map`s) as collections

In the Scala standard library, tuples are frequently used as types of intermediate values. For instance, the Scala type `Map[K, V]` represents a dictionary with keys of type `K` and values of type `V`. Here `K` and `V` are **type parameters**; they represent unknown types that will be chosen later.

In order to create a dictionary with given keys and values, we can write

```
Map(("apples", 3), ("oranges", 2), ("pears", 0))
```

This is equivalent to taking a sequence of key/value pairs and converting it to a dictionary.

Pairs are used so often that the Scala library defines a special infix syntax for pairs, using the arrow symbol `->`. the expression `x -> y` is equivalent to the pair `(x, y)`:

```
scala> "apples" -> 3
res0: (String, Int) = (apples,3)
```

With this syntax, the code for creating a dictionary looks visually clearer:

```
Map("apples" -> 3, "oranges" -> 2, "pears" -> 0)
```

A list of pairs can be converted to a dictionary using the method `.toMap`. The same method works for other list-like collection types such as `Seq`, `Vector`, `Iterator`, or `Array`.

The method `.toSeq` converts a dictionary into a sequence of pairs:

```
scala> Map("apples" -> 3, "oranges" -> 2, "pears" -> 0).toSeq
res20: Seq[(String, Int)] = ArrayBuffer((apples,3), (oranges,2),
    (pears,0))
```

The `ArrayBuffer` is one of the many list-like data structures defined in the Scala library. All of these data structures are gathered under the common “sequence” type, `Seq`, and the Scala standard library may change the implementation of the

“sequence” for reasons of performance.

The standard library has several useful methods that work with tuple types, such as

- `.map` for dictionaries
- `.filter` for dictionaries
- `.zip` and `.zipWithIndex`
- `.flatten` and `.flatMap`
- `.groupBy`
- `.sliding`

It is important to familiarize yourself with these methods, because it will make programming with sequences and dictionaries much easier.

We have seen in Chapter 1 how the `.map` method works on sequences: it applies a given function to each element of the sequence, and gathers the results into a new sequence. The `.map` method works similarly on dictionaries, except that a dictionary of type `Map[K, V]` is automatically converted to a sequence of pairs, `Seq[(K, V)]`, before applying `.map`. For this reason, the argument of `.map` must be a function operating on tuples of type `(K, V)`. Typically, such functions are written using `case` expressions. Here is an example:

```
scala> val m1 = Map("apples" -> 3, "pears" -> 2, "lemons" -> 0)

scala> m1.map { case (fruit, count) => count*2 }
res0: Seq[Int] = ArrayBuffer(6, 4, 0)
```

If we want to transform a dictionary into another dictionary, we can first create a sequence of pairs and then convert it to dictionary with the `.toMap` method:

```
scala> m1.map { case (fruit, count) => (fruit, count*2) }.toMap
res1: scala.collection.immutable.Map[String,Int] = Map(apples -> 6,
    pears -> 4, lemons -> 0)
```

The `.filter` method works similarly on dictionaries; the predicate needs to be a function of type `(K, V) => Boolean`.

```
scala> m1.filter { case (fruit, count) => count > 0 }.toMap
```

2 Mathematical induction

```
res2: scala.collection.immutable.Map[String,Int] = Map(apples -> 6,
    pears -> 4)
```

The `.zip` method takes two sequences and produces a sequence of pairs, taking one element from each sequence:

```
scala> val s = List(1, 2, 3)
s: List[Int] = List(1, 2, 3)

scala> val t = List(true, false, true)
t: List[Boolean] = List(true, false, true)

scala> s.zip(t)
res3: List[(Int, Boolean)] = List((1,true), (2,false), (3,true))

scala> s zip t
res4: List[(Int, Boolean)] = List((1,true), (2,false), (3,true))
```

In the last line, I used the equivalent infix syntax (`s zip t`) just to illustrate the flexibility of syntax conventions in Scala.

The `.zip` method works equally well on dictionaries, since dictionaries are automatically converted to sequences of pairs.

The `.zipWithIndex` method transforms a sequence into a sequence of pairs, where the second part of the pair is the zero-based index:

```
scala> List("a", "b", "c").zipWithIndex
res5: List[(String, Int)] = List((a,0), (b,1), (c,2))
```

The `.flatten` method converts nested sequences to “flattened” ones:

```
scala> List(List(1, 2), List(2, 3), List(3, 4)).flatten
res6: List[Int] = List(1, 2, 2, 3, 3, 4)
```

Keep in mind that this removes *only one* level of nesting:

```
scala> List(List(List(1), List(2)), List(List(2), List(3))).flatten
res7: List[List[Int]] = List(List(1), List(2), List(2), List(3))
```

The `.flatMap` method is closely related to `.flatten` and can be seen as a shortcut to doing first `.map` and then `.flatten` on a sequence:

```
scala> List(1,2,3,4).flatMap(n => (1 to n).toList)
```

```
res8: List[Int] = List(1, 1, 2, 1, 2, 3, 1, 2, 3, 4)

scala> List(1,2,3,4).map(n => (1 to n).toList)
res9: List[List[Int]] = List(List(1), List(1, 2), List(1, 2, 3),
    List(1, 2, 3, 4))

scala> List(1,2,3,4).map(n => (1 to n).toList).flatten
res10: List[Int] = List(1, 1, 2, 1, 2, 3, 1, 2, 3, 4)
```

It transforms a sequence by mapping each element to a potentially variable number of new elements. At first sight, it is probably not clear why `.flatMap` is useful. However, the use of `.flatMap` (which is related to monads) turns out to be one of the most versatile and powerful design patterns in functional programming.

The `.groupBy` method rearranges data into a dictionary, grouped by a key. The key is computed by a function supplied as the argument to `.groupBy`. For example, we can group together all numbers from a given sequence that have the same remainder after division by 3:

```
scala> List(1,2,3,4).flatMap(n => (1 to n).toList).groupBy(k => k % 3)
res11: scala.collection.immutable.Map[Int,List[Int]] = Map(2 ->
    List(2, 2, 2), 1 -> List(1, 1, 1, 4), 0 -> List(3, 3))
```

The result is a dictionary that maps each key to the sequence of values that have that key. The order of elements in the sequences remains unchanged.

The `.sliding` method creates a sliding window of a given width and returns a nested sequence:

```
scala> (1 to 10).sliding(4).toList
res12: List[scala.collection.immutable.IndexedSeq[Int]] =
    List(Vector(1, 2, 3, 4), Vector(2, 3, 4, 5), Vector(3, 4, 5, 6),
    Vector(4, 5, 6, 7), Vector(5, 6, 7, 8), Vector(6, 7, 8, 9),
    Vector(7, 8, 9, 10))
```

2.1.5 Worked examples: Tuples and collections

Example 2.1.5.1 For a given sequence x_i , compute the sequence of pairs $b_i = (\cos x_i, \sin x_i)$.

Hint: use `.map`, assume `xs: Seq[Double]`.

Solution: We need to produce a sequence that has a pair of values correspond-

2 Mathematical induction

ing to each element of the original sequence. This transformation is precisely what the `.map` method does. So the code is

```
xs.map { x => (math.cos(x), math.sin(x)) }
```

Example 2.1.5.2 In a given sequence x_i , count how many times $\cos x_i > \sin x_i$ occurs.

Hint: use `.count`, assume `xs: Seq[Double]`.

Solution: The method `.count` takes a predicate and returns the number of times the predicate was `true` while evaluated on all elements of the sequence. So the straightforward solution is

```
xs.count { x => math.cos(x) > math.sin(x) }
```

We could also use the solution of Exercise 2.1.5.1 that already computed the cosine and the sine values. The code would then become longer. In particular, doing a `.map` on a sequence of pairs would be normally written using pattern matching for readability:

```
xs.map { x => (math.cos(x), math.sin(x)) }  
  .count { case (cosine, sine) => cosine > sine }
```

Example 2.1.5.3 For given sequences a_i and b_i , compute the sequence of differences $c_i = a_i - b_i$

Hint: use `.zip`, `.map`, assume `as, bs: Seq[Double]`.

Solution: we can use `.zip` on `as` and `bs`, which will give us a sequence of pairs,

```
as.zip(bs) : Seq[(Double, Double)]
```

We then use `.map` on this sequence, to compute the difference between the elements. So the full code is

```
as.zip(bs).map { case (a, b) => a - b }
```

Example 2.1.5.4 In a given sequence a_i , count how many times $a_i > a_{i+1}$ occurs.

Hint: use `.zip` and `.tail`. Make sure `.zip` is applied to sequences of equal length.

Solution:

```
as.take(as.length - 1)
  .zip(as.tail)
  .count { case (a, b) => a > b }
```

Example 2.1.5.5 For a given $k > 0$, compute the sequence $b_i = \max(a_{i-k}, \dots, a_{i+k})$.

Hint: use `.sliding`.

Solution:

```
bs.sliding(k).map(_.max)
```

Example 2.1.5.6 Create a 10×10 multiplication table as a value of type `Map[(Int, Int), Int]`.

Hint: use `.flatMap`.

Solution:

```
(1 to 10).flatMap{ x => (1 to 10)
  .map { i => ((i, x), i*x) }
  .toMap
```

Example 2.1.5.7 For a given sequence x_i , compute the maximum of all of the numbers x_i , $\cos x_i$, $\sin x_i$.

Hint: use `.map`, `.flatMap`, `.max`.

Solution:

```
xs.flatMap{ x => Seq(x, math.cos(x), math.sin(x)) }
  .max
```

Example 2.1.5.8 From a dictionary of type `Map[String, String]` mapping names to addresses, and assuming that the addresses do not repeat, compute a dictionary of type `Map[String, String]` mapping the addresses back to names.

Hint: use `.toMap`, `.map`.

Solution:

```
dict.map { case (name, addr) => (addr, name) }
```

2 Mathematical induction

```
.toMap
```

Example 2.1.5.9 Write the solution of Example 2.1.5.8 as a function with type parameters `Name` and `Address` instead of using the fixed type `String`. Test it with `Name = Boolean` and `Address = Set[String]`.

Solution: The final code is

```
def rev[Name, Address](dict: Map[Name, Address]): Map[Address, Name]
  = {
    dict.map { case (name, addr) => (addr, name) }
      .toMap
  }
```

Example* 2.1.5.10 Given a sequence `words: Seq[String]` of words, compute a sequence of pairs, of type `Seq[(Seq[String], Int)]`, where each inner sequence contains all the words with the same length, and the integer value shows the length. The resulting sequence must be ordered by increasing length of words. For example, the input `Seq("the", "food", "is", "good")` should produce the output

```
Seq( (Seq("is"), 2), (Seq("the"), 3), (Seq("food", "good"), 4) )
```

Hint: use `.groupBy`, `.length`, and `.sortBy`.

Solution: It is clear that we need to group the words by length. The library method `.groupBy` takes a function that computes a grouping key” from each element of a sequence. In our case, we need to group by word length, which is computed with the method `.length` applied to a string. So the first step is

```
words.groupBy{ word => word.length }
```

or, more concisely,

```
words.groupBy(_.length)
```

The result of this expression is a dictionary that maps each length to the words of that length:


```
scala> words.groupBy(_.length)
res0: scala.collection.immutable.Map[Int,Seq[String]] = Map(2 ->
  List(is), 4 -> List(food, good), 3 -> List(the))
```

This is already close to what we want; if we convert this dictionary to a sequence, we will have a list of pairs

```
scala> words.groupBy(_.length).toSeq
res1: Seq[(Int, Seq[String])] = ArrayBuffer((2,List(is)),
  (4,List(food, good)), (3,List(the)))
```

It remains to swap the length and the list of words, and to sort the result by the increasing length. We can do this in any order: first sort, then swap; or first swap, then sort. The final code is

```
words.groupBy(_.length).toSeq
  .sortBy{ case (len, words) => len }
  .map{ case (len, words) => (words, len) }
```

2.1.6 Exercises: Tuples and collections

Exercise 2.1.6.1 Find all pairs i, j within $(0, 1, \dots, 9)$ such that $i + 4 * j > i * j$.

Hint: use `.flatMap`.

Exercise 2.1.6.2 Same task as in Exercise 2.1.6.1, but for i, j, k and the condition $i + 4 * j + 9 * k > i * j * k$

Exercise 2.1.6.3 Given two sequences $p: \text{Seq[String]}$ and $q: \text{Seq[Boolean]}$ of equal length, compute a Seq[String] with those elements of p for which the corresponding element of q is `true`.

Hint: use `.zip`, `.map`, `.filter`.

Exercise 2.1.6.4 Convert a Seq[Int] into a $\text{Seq[(Int, Boolean)]}$ where the `Boolean` value is `true` when the element is followed by a larger value. For example, $\text{Seq}(1, 3, 2, 4)$ is to be converted into $\text{Seq}((1, \text{true}), (3, \text{false}), (2, \text{true}), (4, \text{false}))$. (The last element, 4, has no following element.)

2 Mathematical induction

Exercise 2.1.6.5 Given `p: Seq[String]` and `q: Seq[Int]` of equal length, and assuming that elements of `q` do not repeat, compute a `Map[Int, String]` that maps numbers from `q` to their corresponding strings from `p`.

Exercise 2.1.6.6 Write the solution of Exercise 2.1.6.5 as a function with type parameters `P` and `Q` instead of the fixed types `String` and `Int`; test it with `P = Boolean` and `Q = Set[Int]`.

Exercise 2.1.6.7 Given `p: Seq[String]` and `q: Seq[Int]` of equal length, compute a `Seq[String]` that contains the strings from `p` ordered according to the corresponding numbers from `q`.

Hint: use `.sortBy`.

Exercise 2.1.6.8 Write the solution of Exercise 2.1.6.7 as a function with type parameter `S` instead of the fixed type `String`.

Exercise 2.1.6.9 Given a `Seq[(String, Int)]` showing a list of purchased items (where item names may repeat), compute a `Map[String, Int]` showing the total counts: e.g. for the input

```
Seq(("apple", 2), ("pear", 3), ("apple", 5))
```

the output must be

```
Map("apple" -> 7, "pear" -> 3)
```

Implement this as a function with type parameter `S` instead of `String`.

Hint: use `.groupBy`, `.map`, `.sum`.

Exercise 2.1.6.10 Given a `Seq[List[Int]]`, compute a new `Seq[List[Int]]` where each inner list contains *three* largest elements from the initial inner list (or fewer than three if the initial inner list is shorter).

Hint: use `.sortBy`, `.take`.

Exercise 2.1.6.11 Given two sets `p, q` both of type `Set[Int]`, compute a `Set[(Int, Int)]` representing the Cartesian product of the sets `p` and `q` (that is, the set of all pairs `(x, y)` where `x` is from `p` and `y` is from `q`).

Implement this as a function with type parameters `I, J` instead of `Int`.

Hint: use `.flatMap`.

Exercise* 2.1.6.12 Given a `Seq[Map[Person, Amount]]`, showing the amounts various people paid on each day, compute a `Map[Person, Seq[Amount]]`, showing the sequence of payments for each person (assume `Person` and `Amount` are type parameters).

Hint: use `.flatMap`, `.groupBy`.

2.2 Discussion

2.2.1 Total and partial functions

In Scala, functions can be total or partial. A **total** function will always compute a result value, while a **partial** function may fail to compute its result for certain values of its arguments. A partial function is only well-defined on a subset of possible argument values. Functions that use pattern matching are sometimes partial because pattern matching may fail. In most cases, programs can be written to avoid match errors completely.

If the pattern matching fails, the code will throw an exception and stop running. In functional programming, we usually want to avoid this situation because it makes it much harder to reason about program correctness. An example of a potentially fallible pattern matching is

```
def h(p: (Int, Int)): Int = p match { case (x, 0) => x }

scala> h( (1, 0) )
res0: Int = 1

scala> h( (1, 2) )
scala.MatchError: (1,2) (of class scala.Tuple2mcIIsp)
  at .h(<console>:12)
  ... 32 elided
```

In this case, the pattern contains a pattern variable `x` and a constant `0`. This pattern only matches tuples whose second part is equal to `0`. If the second argument is nonzero, a match error occurs and the program crashes. Therefore, `h` is a partial function.

Pattern matching failures never happen if we match a tuple of correct size with

2 Mathematical induction

a pattern such as (x, y, z) , because the pattern variables will always match whatever values the tuple has. On the other hand, the compiler will generate a compile-time error if a function is applied to a tuple of incorrect type. For this reason, pattern matching with a pattern such as (x, y, z) is **infallible** (never fails).

2.2.2 Scope of pattern matching variables

Pattern matching constructions introduce locally scoped variables. Consider this code,

```
def f(x: (Int, Int)): Int = x match {  
  case (x, y) => x + y  
}
```

3 The formal logic of types

3.1 Types of higher-order functions

3.1.1 Curried functions

Consider a function with type signature `Int => (Int => Int)`. This is a function that takes an integer and returns a *function* that again takes an integer and then returns an integer. So, we obtain an integer result only after we apply the function to *two* integer values, one after another. This is, in a sense, equivalent to a function having two arguments, except the application of the function needs to be done in two steps, applying it to one argument at a time. Functions of this sort are called **curried** functions.

One way of defining such a function is

```
def f0(x: Int): Int => Int = { y => x - y }
```

The function takes an integer argument x and returns the expression $y \Rightarrow x - y$, which is a function of type $\text{Int} \Rightarrow \text{Int}$. Thus, the type of `f0` is $\text{Int} \Rightarrow (\text{Int} \Rightarrow \text{Int})$.

To use `f0`, we apply it to an integer value. The result is a value of function type:

```
scala> val r = f0(20)
```

The value `r` can be now applied to another integer argument:

```
scala> r(4)
```

In Scala, $\text{Int} \Rightarrow \text{Int} \Rightarrow \text{Int}$ means the same as $\text{Int} \Rightarrow (\text{Int} \Rightarrow \text{Int})$ and $x \Rightarrow y \Rightarrow x - y$ means the same as $x \Rightarrow (y \Rightarrow x - y)$. In other words, the function symbol \Rightarrow associates to the right. Thus, the type signature of `f0` may be equivalently written as $\text{Int} \Rightarrow \text{Int} \Rightarrow \text{Int}$.

An equivalent way of defining a function with the same type signature is

```
val f1: Int => Int => Int = x => y => x - y
```

3 The formal logic of types

We can also define a function that takes its two arguments at once. Such a function will have a different type signature, for instance $(\text{Int}, \text{Int}) \Rightarrow \text{Int}$ as in this example:

```
def f2(x: Int, y: Int): Int = x - y
```

The syntax for calling the functions `f1` and `f2` is different:

```
scala> f1(20)(4)
scala> f2(20, 4)
```

The main difference between the usage of `f1` and `f2` is that `f2` must be applied at once to both arguments, while `f1(20)` can be evaluated separately, with only the first argument, 20, the result being a *function* that can be later applied to another argument:

```
scala> val r = f1(20)
scala> r(4)
```

Applying a curried function to some but not all of possible arguments is called **partial application**.

More generally, a curried function may have a type signature of the form $A \Rightarrow B \Rightarrow C \Rightarrow \dots \Rightarrow P \Rightarrow Z$, where A, B, C, \dots, Z are some types. Arguments of types A, B, \dots, P are sometimes called the **curried arguments** of the function. This function is, in a sense, equivalent to an **uncurried** function with type signature $(A, B, C, \dots, P) \Rightarrow Z$. The uncurried function takes all arguments at once. The equivalence of curried and uncurried functions is not equality — these functions are *different*; but one of them can be easily reconstructed from the other if necessary. I call such functions **isomorphic**.

From the point of view of programming language theory, curried functions are simpler because they always have a *single* argument (and may return a function that will consume further arguments). From the point of view of programming practice, curried functions are not often used and may be harder to read.

In the short notation (used e.g. in OCaml and Haskell), a curried function such as `f2` is applied to its arguments as `f2 20 4`. This departs further from the mathematical notation and requires some getting used to. If the two arguments are more complicated than just 20 and 4, the resulting expression may become significantly harder to read, since no commas are used to separate the arguments. For this reason, programmers may prefer to first define short names for complicated expressions and then use these names as curried arguments.

3.1 Types of higher-order functions

In Scala, the choice of whether to use curried or uncurried function signatures is largely a matter of syntactic convenience. Most Scala code tends to be written with uncurried functions, while curried functions are used when they produce more easily readable code. One of the syntactic features for curried functions in Scala is the ability to supply a curried argument using the curly brace syntax. Compare the two definitions of the function `summation` described earlier:

```
def summation1(a: Int, b: Int, g: Int ⇒ Int): Int =  
  (a to b).map(g).sum  
def summation2(a: Int, b: Int)(g: Int ⇒ Int): Int =  
  (a to b).map(g).sum  
summation1(1, 10, x ⇒ x*x*x + 2*x)  
summation2(1, 10) { x ⇒ x*x*x + 2*x }
```

The code for `summation2` may be easier to read because the curried argument is syntactically separated from the rest of the code by curly braces. This is especially useful when the curried argument is itself a function, since the Scala curly braces syntax allows function bodies to contain their own local definitions (`val` or `def`).

3.1.2 Calculations with nameless functions

We now need to gain experience working with nameless functions.

In mathematics, functions are evaluated by substituting their argument values into their body. Nameless functions are evaluated in the same way. For example, applying the nameless function $x \Rightarrow x + 10$ to an integer 2, we substitute 2 instead of x in “ $x + 10$ ” and get “ $2 + 10$ ”, which we then evaluate to 12. The computation is written like this,

$$(x \Rightarrow x + 10)(2) = 2 + 10 = 12 \quad .$$

Nameless functions are *values* and can be used as part of larger expressions, just as any other values. For instance, nameless functions can be arguments of other functions (nameless or not). Here is an example of applying a nameless function $f \Rightarrow f(2)$ to a nameless function $x \Rightarrow x + 4$:

$$(f \Rightarrow f(2))(x \Rightarrow x + 4) = (x \Rightarrow x + 4)(2) = 6 \quad .$$

In the nameless function $f \Rightarrow f(2)$, the argument f has to be itself a function, otherwise the expression $f(2)$ would make no sense. In this example, f must have type $\text{Int} \Rightarrow \text{Int}$.

3 The formal logic of types

There are some standard conventions for reducing the number of parentheses when writing expressions involving nameless functions, especially curried functions:

- Function expressions group everything to the right:
 $x \Rightarrow y \Rightarrow z \Rightarrow e$ means the same as $x \Rightarrow (y \Rightarrow (z \Rightarrow e))$.
- Function applications group everything to the left:
 $f(x)(y)(z)$ means $((f(x))(y))(z)$.
- Function applications group stronger than infix operations:
 $x + f(y)$ means $x + (f(y))$, just like in mathematics.

To specify the type of the argument, I will use the underlined superscript, for example: $x^{\underline{\text{Int}}} \Rightarrow x + 2$.

The convention of grouping functions to the right reduces the number of parentheses for curried function types used most often. It is rare to find use for function types such as $((a \Rightarrow b) \Rightarrow c) \Rightarrow d$ that require many parentheses with this convention.

Here are some more examples of performing function applications symbolically. I will omit types for brevity, since every non-function value is of type `Int` in these examples.

$$\begin{aligned} (x^{\underline{\text{Int}}} \Rightarrow x * 2) (10) &= 10 * 2 = 20 \quad . \\ (p \Rightarrow z \Rightarrow z * p) (t) &= (z \Rightarrow z * t) \quad . \\ (p \Rightarrow z \Rightarrow z * p) (t)(4) &= (z \Rightarrow z * t)(4) = 4 * t \quad . \end{aligned}$$

Some results of these computation are integer values such as 20; in other cases, results are *function values* such as $z \Rightarrow z * t$.

In the following examples, some function arguments are themselves functions:

$$\begin{aligned} (f \Rightarrow p \Rightarrow f(p)) (g \Rightarrow g(2)) &= (p \Rightarrow p(2)) \quad . \\ (f \Rightarrow p \Rightarrow f(p)) (g \Rightarrow g(2)) (x \Rightarrow x + 4) &= (p \Rightarrow p(2)) (x \Rightarrow x + 4) \\ &= 2 + 4 = 6 \quad . \end{aligned}$$

Here I have been performing calculations step by step, as usual in mathematics. A Scala program is evaluated in a similar way at run time.

3.1.3 Short syntax for function applications

In mathematics, function applications are sometimes written without parentheses, for instance $\cos x$ or $\arg z$. There are also cases where formulas such as $\sin 2x = 2 \sin x \cos x$ imply parentheses as $\sin(2x) = 2 \cdot \sin(x) \cdot \cos(x)$.

Many programming languages (such as ML, OCaml, F#, Haskell, Elm, PureScript) have adopted this “short syntax”, in which parentheses are optional for function arguments. The result is a concise notation where $f\ x$ means the same as $f(x)$. Parentheses are still used where necessary to avoid ambiguity or for readability.¹

The conventions for nameless functions in the short syntax become:

- Function expressions group everything to the right:
 $x \Rightarrow y \Rightarrow z \Rightarrow e$ means $x \Rightarrow (y \Rightarrow (z \Rightarrow e))$.
- Function applications group everything to the left:
 $f\ x\ y\ z$ means $((f\ x)\ y)\ z$.
- Function applications group stronger than infix operations: $x + f\ y$ means $x + (f\ y)$, just like in mathematics $x + \cos y$ groups $\cos y$ stronger than the infix “+” operation.

So, $x \Rightarrow y \Rightarrow a\ b\ c + p\ q$ means $x \Rightarrow (y \Rightarrow ((a\ b)\ c) + (p\ q))$. When this notation becomes hard to read correctly, one needs to add parentheses, e.g. to write $f(x \Rightarrow g\ h)$ instead of $f\ x \Rightarrow g\ h$.

In this book, I will sometimes use this “short syntax” when reasoning about code. Scala does not support the short syntax; in Scala, parentheses need to be put around every curried argument. The infix method syntax such as `List(1,2,3).map func1` does not apply to curried functions in Scala.

3.1.4 Higher-order functions

The **order** of a function is the number of function arrows “ \Rightarrow ” contained in the type signature of that function. If a function’s type signature contains more than one function arrow, the function is called a **higher-order** function. A higher-order function takes a function as argument and/or returns a function as its result value.

¹The no-parentheses syntax is also used in Unix shell commands, for example `cp file1 file2`, as well as in the programming language Tcl. In LISP and Scheme, each function application is enclosed in parentheses but the arguments are separated by spaces, for example `(+ 1 2 3)`.

3 The formal logic of types

Examples:

```
def f1(x: Int): Int = x + 10
```

The function `f1` has type signature $\text{Int} \Rightarrow \text{Int}$ and order 1, so it is *not* a higher-order function.

```
def f2(x: Int): Int  $\Rightarrow$  Int = z  $\Rightarrow$  z + x
```

The function `f2` has type signature $\text{Int} \Rightarrow \text{Int} \Rightarrow \text{Int}$ and is a higher-order function, of order 2.

```
def f3(g: Int  $\Rightarrow$  Int): Int = g(123)
```

The function `f3` has type signature $(\text{Int} \Rightarrow \text{Int}) \Rightarrow \text{Int}$ and is a higher-order function of order 2.

Although `f2` is a higher-order function, its higher-orderness comes from the fact that the return value is of function type. An equivalent computation can be performed by an uncurried function that is not higher-order:

```
scala> def f2u(x: Int, z: Int): Int = z + x
```

The Scala library defines methods to transform between curried and uncurried functions:

```
scala> def f2u(x: Int, z: Int): Int = z + x
scala> val f2c = (f2u _).curried
scala> val f2u1 = Function.uncurried(f2c)
```

The syntax `(f2u _)` is used in Scala to convert methods to function values. Recall that Scala has two ways of defining a function: one as a method (using `def`), another as a function value (using `val`).

The methods `.curried` and `.uncurried` can be easily implemented in Scala code, as we will see in the worked examples.

Unlike `f2`, the function `f3` cannot be converted to a non-higher-order function because `f3` has an *argument* of function type, rather than a return value of function type. Converting to an uncurried form cannot eliminate an argument of function type.

3.1.5 Worked examples: higher-order functions

1. Using both `def` and `val`, define a function that...

- a) ...adds 20 to its integer argument.

```
def fa(i: Int): Int = i + 20
val fa_v: (Int ⇒ Int) = k ⇒ k + 20
```

It is not necessary to specify the type of the argument k because we already fully specified the type $(Int \Rightarrow Int)$ of `fa_v`. The parentheses around the type of `fa_v` are optional, I added them for clarity.

- b) ...takes an integer x , and returns a *function* that adds x to its argument.

```
def fb(x: Int): (Int ⇒ Int) = k ⇒ k + x
val fb_v: (Int ⇒ Int ⇒ Int) = x ⇒ k ⇒ k + x
def fb_v2(x: Int)(k: Int): Int = k + x
```

Since functions are values, we can directly return new functions. When defining the right-hand sides as function expressions in `fb` and `fb_v`, it is not necessary to specify the type of the arguments x and k because we already fully specified the type signatures of `fb` and `fb_v`. The last version, `fb_v2`, may be easier to read and is equivalent to `fb_v`.

- c) ...takes an integer x and returns true iff $x + 1$ is a prime. Use the function `is_prime` defined previously.

```
def fc(x: Int): Boolean = is_prime(x + 1)
val fc_v: (Int ⇒ Boolean) = x ⇒ is_prime(x + 1)
```

- d) ...returns its integer argument unchanged. (This is called the **identity function** for integer type.)

```
def fd(i: Int): Int = i
val fd_v: (Int ⇒ Int) = k ⇒ k
```

- e) ...takes x and always returns 123, ignoring its argument x . (This is called a **constant function**.)

```
def fe(x: Int): Int = 123
val fe_v: (Int ⇒ Int) = x ⇒ 123
```

To emphasize the fact that the argument x is ignored, use the special syntax where x is replaced by the underscore:

```
val fe_v1: (Int ⇒ Int) = _ ⇒ 123
```

3 The formal logic of types

- f) ...takes x and returns a constant function that always returns the fixed value x . (This is called the **constant combinator**.)

```
def ff(x: Int): Int ⇒ Int = _ ⇒ x
val ff_v: (Int ⇒ Int ⇒ Int) = x ⇒ _ ⇒ x
def ff_v2(x: Int)(y: Int): Int = x
```

The syntax of `ff_v2` may be easier to read, but then we cannot omit the name `y` for the unused argument.

2. Define a function `comp` that takes two functions $f : \text{Int} \Rightarrow \text{Double}$ and $g : \text{Double} \Rightarrow \text{String}$ as arguments, and returns a new function that computes $g(f(x))$. What is the type of the function `comp`?

```
def comp(f: Int ⇒ Double, g: Double ⇒ String): (Int ⇒ String) =
  x ⇒ g(f(x))
scala> val f: Int ⇒ Double = x ⇒ 5.67 + x
scala> val g: Double ⇒ String = x ⇒ f"x=%3.2f"
scala> val h = comp(f, g)
scala> h(10)
```

The function `comp` has two arguments, of types $\text{Int} \Rightarrow \text{Double}$ and $\text{Double} \Rightarrow \text{String}$. The result value of `comp` is of type $\text{Int} \Rightarrow \text{String}$, because `comp` returns a new function that takes an argument x of type Int and returns a String . So the full type signature of the function `comp` is written as

```
/// (Int ⇒ Double, Double ⇒ String) ⇒ (Int ⇒ String)
```

This is an example of a function that both takes other functions as arguments *and* returns a new function.

3. Define a function `uncurry2` that takes a curried function of type $\text{Int} \Rightarrow \text{Int} \Rightarrow \text{Int}$ and returns an uncurried equivalent function of type $(\text{Int}, \text{Int}) \Rightarrow \text{Int}$.

```
def uncurry2(f: Int ⇒ Int ⇒ Int): (Int, Int) ⇒ Int =
  (x, y) ⇒ f(x)(y)
```

3.2 Discussion

3.2.1 Scope of bound variables

A bound variable is invisible outside the scope of the expression (often called **local scope** whenever it is clear which expression is being considered). This is

why bound variables may be renamed at will: no outside code could possibly use them and depend on their values. However, outside code may define variables that (by chance or by mistake) have the same name as a bound variable inside the scope.

Consider this example from calculus: In the integral

$$f(x) = \int_0^x \frac{dx}{1+x} \quad ,$$

a bound variable named x is defined in *two* local scopes: in the scope of f and in the scope of the nameless function $x \Rightarrow \frac{1}{1+x}$. The convention in mathematics is to treat these two x 's as two *completely different* variables that just happen to have the same name. In sub-expressions where both of these bound variables are visible, priority is given to the bound variable defined in the closest inner scope. The outer definition of x is **shadowed**, i.e. hidden, by the definition of the inner x . For this reason, mathematicians expect that evaluating $f(10)$ will give

$$f(10) = \int_0^{10} \frac{dx}{1+x} \quad ,$$

rather than $\int_0^{10} \frac{dx}{1+10}$, because the outer definition $x = 10$ is shadowed, within the expression $\frac{1}{1+x}$, by the closer definition of x in the local scope of $x \Rightarrow \frac{1}{1+x}$.

Since this is the prevailing mathematical convention, the same convention is adopted in FP. A variable defined in a local scope (i.e. a bound variable) is invisible outside that scope but will shadow any outside definitions of a variable with the same name.

It is better to avoid name shadowing, because it usually decreases the clarity of code and thus invites errors. Consider this function,

$$x \Rightarrow x \Rightarrow x \quad .$$

This code may look confusing; let us decipher this syntax. The symbol \Rightarrow associates to the right, so $x \Rightarrow x \Rightarrow x$ is the same as $x \Rightarrow (x \Rightarrow x)$. So, it is a function that takes x and returns $x \Rightarrow x$. Since the returned nameless function, $(x \Rightarrow x)$, may be renamed to $(y \Rightarrow y)$ without changing its value, we can rewrite the code to

$$x \Rightarrow (y \Rightarrow y) \quad .$$

It is now easier to understand this code and reason about it. For instance, it becomes immediately clear that this function actually ignores its argument x .

3.3 Exercises

1. Define a function of type $\text{Int} \Rightarrow \text{List}[\text{List}[\text{Int}]] \Rightarrow \text{List}[\text{List}[\text{Int}]]$ similar to Exercise 7 except that the hard-coded number 100 must be a *curried* first argument. Implement Exercise 7 using this function.
2. Define a function q that takes a function $f : \text{Int} \Rightarrow \text{Int}$ as its argument, and returns a new function that computes $f(f(f(x)))$. What is the required type of the function q ?
3. Define a function `curry2` that takes an uncurried function of type $(\text{Int}, \text{Int}) \Rightarrow \text{Int}$ and returns a curried equivalent function of type $\text{Int} \Rightarrow \text{Int} \Rightarrow \text{Int}$.

3.4 Disjunctive types

3.4.1 Discussion

3.5 The Curry-Howard correspondence

3.5.1 Discussion

4 Functors

4.1 Discussion

4.2 Practical use

4.3 Laws and structure

5 Type-level functions and type classes

5.1 Combining typeclasses

5.2 Inheritance

5.3 Functional dependencies

5.4 Discussion

6 Filterable functors

6.1 Practical use

6.1.1 Discussion

6.2 Laws and structure

6.2.1 Discussion

7 Semimonads and monads

7.1 Practical use

7.1.1 Discussion

7.2 Laws and structure

7.2.1 Discussion

8 Applicative functors, contrafunctors, and profunctors

8.1 Practical use

8.1.1 Discussion

8.2 Laws and structure

9 Traversable functors and profunctors

9.1 Discussion

10 Free constructions

10.1 Discussion

11 Monad transformers

11.1 Practical use

11.2 Laws and structure

A monad transformer $T_L^{M,A}$ is a type constructor with a type parameter A and a monad parameter M , such that the following five laws hold:

1. **Monad construction law:** $T_L^{M,\bullet}$ is a lawful monad for any monad M . For instance, the transformed monad $T_L^{M,\bullet}$ has methods pu_T and ftn_T that satisfy the monad laws.
2. **Identity law:** $T_L^{\text{Id},\bullet} \cong L^\bullet$ via a monadic isomorphism, where Id is the identity monad, $\text{Id}^A \triangleq A$.
3. **Lifting law:** For any monad M , the function $\text{lift} : M^A \Rightarrow T_L^{M,A}$ is a monadic morphism. (In a shorter notation, $\text{lift} : M \rightsquigarrow T_L^M$.)
4. **Runner law:** For any monads M, N and any monadic morphism $\phi : M \rightsquigarrow N$, the runner $\text{mrun } \phi : T_L^M \rightsquigarrow T_L^N$ is a monadic morphism. Moreover, the function mrun is a lawful lifting of monadic morphisms from $M \rightsquigarrow N$ to $T_L^M \rightsquigarrow T_L^N$:

$$\text{mrun id} = \text{id} \quad , \quad \text{mrun } \phi \circ \text{mrun } \chi = \text{mrun } (\phi \circ \chi) \quad .$$

5. **Base runner law:** For any monadic morphism $\theta : L \rightsquigarrow \text{Id}$ and for any monad M , the base runner $\text{brun } \theta : T_L^M \rightsquigarrow M$ is a monadic morphism. This morphism must satisfy the nondegeneracy law

$$\text{lift}; \text{brun } \theta = \text{id} \quad .$$

(Since it is not possible to transform the base monad L into an arbitrary other monad, there are no lifting laws for brun , unlike mrun .)

11.2.1 Properties of monadic morphisms

Statement 11.2.1.1 If $\phi : M^A \Rightarrow N^A$ is a monadic morphism (here M and N are monads) then ϕ is also a natural transformation between functors M and N .

Statement 11.2.1.2 If L, M, N are monads and $\phi : L^\bullet \rightsquigarrow M^\bullet$ and $\chi : M^\bullet \rightsquigarrow N^\bullet$ are monadic morphisms then the composition $\phi \circ \chi : L^\bullet \rightsquigarrow N^\bullet$ is also a monadic morphism.

Statement 11.2.1.3 For any monad M , the function $\Delta : M^A \Rightarrow M^A \times M^A$ is a monadic morphism between monads M and $M \times M$.

Statement 11.2.1.4 For any monads K, L, M, N and monadic morphisms $\phi : K^\bullet \rightsquigarrow M^\bullet$ and $\chi : L^\bullet \rightsquigarrow N^\bullet$, the componentwise function product $\phi \boxtimes \chi : K^\bullet \times L^\bullet \rightsquigarrow M^\bullet \times N^\bullet$ is a monadic morphism.

Statement 11.2.1.5 For any monads M and N , the function $\nabla_1 : M^\bullet \times N^\bullet \rightsquigarrow M^\bullet$ is a monadic morphism. Same for $\nabla_2 : M^\bullet \times N^\bullet \rightsquigarrow N^\bullet$.

Statement 11.2.1.6 For any monads M and N , the component-swapping function $\sigma : M^\bullet \times N^\bullet \rightsquigarrow N^\bullet \times M^\bullet$ is a monadic morphism.

Proof: The code for σ can be written as $\sigma = \Delta \circ (\nabla_2 \boxtimes \nabla_1)$. Written in this way, the code clearly shows that it is a monadic morphism: it is a composition of Δ , which is a monadic morphism, and a product of monadic morphisms ∇_1 and ∇_2 .

11.3 Functor composition using transformed monads

Suppose we are working with a base monad L and a foreign monad M , and we have constructed the transformed monad T_L^M . In this section, let us denote this monad simply by T .

A useful property of monad transformers is that the monad T adequately describes the effects of both monads L and M at the same time. Suppose we are

11.4 Stacking composition of monad transformers

working with a deeply nested type constructor involving many functor layers of monads L , M , and T such as

$$T^{M^{T^L M^L A}}.$$

The properties of the transformer allow us to convert this type to a single layer of the transformed monad T . In this example, we will have a natural transformation

$$T^{M^{T^L M^L A}} \Rightarrow T^A.$$

To achieve this, we first use the methods `blift` and `lift` to convert each layer of L or M to a layer of T , raising to functors as necessary. The result will be a number of nested layers of T . Second, we use `fnT` as many times as necessary to flatten all layers of T into a single layer.

11.4 Stacking composition of monad transformers

11.4.1 Stacking two monads

Suppose we know the transformers T_P and T_Q for some given monads P and Q . We can transform Q with P and obtain a monad $R^A \triangleq T_P^{Q,A}$. What would be the monad transformer T_R for the monad R ?

A simple solution is to first transform the foreign monad M with T_Q , obtaining a new monad $T_Q^{M,\bullet}$, and then to transform that new monad with T_P . So the formula for the transformer T_R is

$$T_R^{M,A} = T_P^{T_Q^{M,\bullet},A}.$$

Here the monad $T_Q^{M,\bullet}$ was substituted into $T_P^{M,A}$ as the foreign monad M (not as the type parameter A). This way of composition is called **stacking** the monad transformers.

In Scala code, this “stacking” composition is written as

```
type RT[M, A] = PT[QT[M, ?], A]
```

The resulting monad is a **stack** of three monads P , Q , and M . The order of monads in the stack is significant since, in general, there will be no monadic isomorphism between monads stacked in a different order.

11 Monad transformers

We will now show that the transformer T_R is lawful (satisfies all five monad transformer laws), as long as the same five laws are satisfied by both T_P and T_Q . To shorten the notation, I will talk about a “monad T_P^M ” meaning the monad defined as $T_P^{M,\bullet}$ or, more verbosely, the monad $G^A \triangleq T_P^{M,A}$.

Monad construction law We need to show that $T_P^{T_Q^M}$ is a monad for any monad M . The monad law for T_Q says that T_Q^M is a monad. Denote this monad temporarily by $S \triangleq T_Q^M$. The monad law for T_P says that T_P^M is a monad for any monad M ; in particular, for $M = S$. Therefore, T_P^S is a monad, as required.

Identity law We need to show that $T_P^{T_Q^{\text{Id}}} \cong T_P^Q$ via a monadic isomorphism. The identity law for T_Q says that $T_Q^{\text{Id}} \cong Q$ via a monadic isomorphism. So, we have a monadic morphism $\phi : Q \rightsquigarrow T_Q^{\text{Id}}$ and its inverse, $\chi : T_Q^{\text{Id}} \rightsquigarrow Q$. The runner law for T_P can be applied to both ϕ and χ since they are monadic morphisms. So we obtain two new monadic morphisms,

$$\text{mr} \text{run}_P \phi : T_P^Q \rightsquigarrow T_P^{T_Q^{\text{Id}}} \quad ; \quad \text{mr} \text{run}_P \chi : T_P^{T_Q^{\text{Id}}} \rightsquigarrow T_P^Q \quad .$$

Are these two monadic morphisms inverses of each other? Since $\text{mr} \text{run}_P$ is a lifting of monadic morphisms, it lifts identity into identity and composition into composition. So the composition $\phi \circ \chi = \text{id}$ is lifted into

$$\text{mr} \text{run} \phi \circ \text{mr} \text{run} \chi = \text{id} \quad ,$$

and similarly for the composition $\chi \circ \phi$. Thus we obtain a monadic isomorphism between monads T_P^Q and $T_P^{T_Q^{\text{Id}}}$.

Lifting law We need to show that there exists a monadic morphism $M \rightsquigarrow T_P^{T_Q^M}$ for any monad M . The lifting law for T_Q gives a monadic morphism $\text{lift}_Q : M \rightsquigarrow T_Q^M$. The lifting law for T_P can be applied to the monad T_Q^M , which gives a monadic morphism

$$\text{lift}_P : T_Q^M \rightsquigarrow T_P^{T_Q^M} \quad .$$

The composition of this with lift_Q is a monadic morphism of the required type $M \rightsquigarrow T_P^{T_Q^M}$. (We use the fact that a composition of monadic morphisms is again a monadic morphism.)

11.4 Stacking composition of monad transformers

Runner law We need to show that there exists a lawful lifting

$$\text{mrun}_R : (M \rightsquigarrow N) \Rightarrow T_P^{T_Q^M} \rightsquigarrow T_P^{T_Q^N} .$$

First, we have to define $\text{mrun}_R \phi$ for any given $\phi : M \rightsquigarrow N$. We use the lifting law for T_Q to get a monadic morphism

$$\text{lift}_Q \phi : T_Q^M \rightsquigarrow T_Q^N .$$

Now we can apply the lifting law for T_P to this monadic morphism and obtain

$$\text{lift}_P (\text{lift}_Q \phi) : T_P^{T_Q^M} \rightsquigarrow T_P^{T_Q^N} .$$

This function has the correct type signature. So we can define

$$\text{lift}_R \triangleq \text{lift}_P \circ \text{lift}_Q .$$

It remains to prove that lift_R is a lawful lifting. We use the fact that both lift_P and lift_Q are lawful liftings; we need to show that their composition is also a lawful lifting. To verify the identity law of lifting, apply lift_R to an identity function $\text{id} : M \rightsquigarrow M$,

$$\begin{aligned} \text{lift}_R (\text{id}^{M \rightsquigarrow M}) &= \text{lift}_P (\text{lift}_Q \text{id}^{M \rightsquigarrow M}) \\ \text{identity law for } \text{lift}_Q : &= \text{lift}_P (\text{id}^{T_Q^M \rightsquigarrow T_Q^M}) \\ \text{identity law for } \text{lift}_P : &= \text{id} . \end{aligned}$$

To verify the composition law of lifting, apply lift_R to a composition of two monadic morphisms $\phi : L \rightsquigarrow M$ and $\chi : M \rightsquigarrow N$,

$$\begin{aligned} \text{lift}_R (\phi \circ \chi) &= \text{lift}_P (\text{lift}_Q (\phi \circ \chi)) \\ \text{composition law for } \text{lift}_Q : &= \text{lift}_P (\text{lift}_Q \phi \circ \text{lift}_Q \chi) \\ \text{composition law for } \text{lift}_P : &= \text{lift}_P (\text{lift}_Q \phi) \circ \text{lift}_P (\text{lift}_Q \chi) \\ &= \text{lift}_R \phi \circ \text{lift}_R \chi . \end{aligned}$$

Base runner law We need to show that for any monadic morphism $\theta : T_P^Q \rightsquigarrow \text{Id}$ and for any monad M , there exists a monadic morphism $\text{brun}_R \theta : T_P^{T_Q^M} \rightsquigarrow M$. To

11 Monad transformers

define this morphism for a given θ , we clearly need to use the base runners for T_P and T_Q . The base runner for T_Q has the type signature

$$\text{brun}_Q : (Q \rightsquigarrow \text{Id}) \Rightarrow T_Q^M \rightsquigarrow M \quad .$$

We can apply the base runner for T_P to T_Q^M as the foreign monad,

$$\text{brun}_P : (P \rightsquigarrow \text{Id}) \Rightarrow T_P^{T_Q^M} \rightsquigarrow T_Q^M \quad .$$

It is now clear that we could obtain a monadic morphism $T_P^{T_Q^M} \rightsquigarrow M$ if we had some monadic morphisms $\phi : P \rightsquigarrow \text{Id}$ and $\chi : Q \rightsquigarrow \text{Id}$,

$$\text{brun}_P \phi \circ \text{brun}_Q \chi : T_P^{T_Q^M} \rightsquigarrow M \quad .$$

However, we are only given a single monadic morphism $\theta : T_P^Q \rightsquigarrow \text{Id}$. How can we compute ϕ and χ out of θ ? We can use the liftings $\text{blift}_P : P \rightsquigarrow T_P^Q$ and $\text{lift}_P : Q \rightsquigarrow T_P^Q$, which are both monadic morphisms, and compose them with θ :

$$\text{blift}_P \circ \theta : P \rightsquigarrow \text{Id} \quad ; \quad \text{lift}_P \circ \theta : Q \rightsquigarrow \text{Id} \quad .$$

So it follows that we can define a monadic morphism

$$\text{brun}_R \theta = \text{brun}_P (\text{blift}_P \circ \theta) \circ \text{brun}_Q (\text{lift}_P \circ \theta) : T_P^{T_Q^M} \rightsquigarrow M \quad .$$

Since we have defined $\text{brun}_R \theta$ as a composition of monadic morphisms, it is itself a monadic morphism.

11.4.2 Stacking multiple monads

The monad transformer for T_P^Q can be applied to another monad K ; the result is the transformed monad

$$S^A \triangleq T_P^{T_Q^K, A}.$$

What is the monad transformer for the monad S ? Assuming that we know the monad transformer T_K , we could stack the transformers one level higher:

$$T_S^{M, A} \triangleq T_P^{T_Q^{T_K^M}, A} \quad .$$

11.5 Monad transformers that use functor composition

This looks like a stack of four monads P , Q , K , and M . Note that the type parameter A is used as $T_P^{(\dots),A}$, that is, it belongs to the outer transformer T_P .

We can define a transformer stack for any number of monads P, Q, \dots, Z in the same way,

$$T_S^{M,A} \triangleq T_P^{T_Q \cdots T_Z^M, A} . \quad (11.1)$$

The type parameter A will always remain at the outer transformer level, while the foreign monad M will be in the innermost nested position.

It appears plausible that T_S is a lawful monad transformer for *any* number of stacked monads. We can prove this by induction on the number of monads. In the previous section, we have derived the transformer laws for any *three* stacked monads (two monads P, Q within the transformer and one foreign monad M). Now we need to derive the same laws for a general transformer stack, such as that in Eq. (11.1). Let us temporarily denote by J the monad

$$J \triangleq T_Q \cdots T_Z^{\text{Id}} ,$$

where we replaced the foreign monad M by the identity monad. The monad J is a shorter transformer stack than S , so the inductive assumption tells us that the transformer laws already hold for the transformer T_J defined by

$$T_J^M \triangleq T_Q \cdots T_Z^M .$$

Since both T_P and T_J are lawful transformers, their stacking composition $T_P^{T_J^M, A}$ is also a lawful transformer (this was shown in the previous section). In our notation, $T_S^{M,A} = T_P^{T_J^M, A}$, and so we have shown that $T_S^{M,A}$ is a lawful transformer.

11.5 Monad transformers that use functor composition

We have seen examples of monad transformers that work via functor composition, either from inside or from outside. The simplest such examples are the `OptionT` transformer,

$$L^A \triangleq 1 + A, \quad T_L^{M,A} \triangleq M^{L^A} = M^{1+A} ,$$

11 Monad transformers

which puts the base monad L *inside* the monad M , and the ReaderT transformer,

$$L^A \triangleq R \Rightarrow A, \quad T_L^{M,A} \triangleq L^{M^A} = R \Rightarrow M^A, \quad ,$$

which puts the base monad L *outside* the foreign monad M .

We can obtain many properties of both kinds of monad transformers at once from a single derivation if we temporarily drop the distinction between the base monad and the foreign monad. We simply assume that we are given two different monads, L and M , whose functor composition $T^\bullet \triangleq L^{M^\bullet}$ also happens to be a monad. Since the assumptions on the monads L and M are the same, the resulting properties of the composed monad T will apply equally to both kinds of monad transformers.

To interpret the results, we will assume that L is the base monad for the composed-outside transformers, and that M is the base monad for the composed-inside transformers. For instance, we will be able to prove the laws of liftings $L \rightsquigarrow T$ and $M \rightsquigarrow T$ regardless of the choice of the base monad.

What properties of monad transformers will *not* be derivable in this way? Monad transformers depend on the structure on the base monad, but not on the structure of the foreign monad; the transformer's methods `pure` and `flatten` are generic in the foreign monad. This is expressed via the monad transformer laws for the runners `mrn` and `brn`, which we will need to derive separately for each of the two kinds of transformers.

11.5.1 Motivation for the swap function

The first task is to show that the composed monad $T^\bullet \triangleq L^{M^\bullet}$ obeys the monad laws. For this, we need to define the methods for the monad T , namely `pure` (short notation "`puT`") and `flatten` (short notation "`ftnT`"), with the type signatures

$$\text{pu}_T : A \Rightarrow L^{M^A}, \quad \text{ftn}_T : L^{M^{L^{M^A}}} \Rightarrow L^{M^A}.$$

How can we implement these methods? *All we know* about L and M is that they are monads with their own methods `puL`, `ftnL`, `puM`, and `ftnM`. We can easily implement

$$\text{pu}_T \triangleq \text{pu}_M \circ \text{pu}_L. \quad (11.2)$$

11.5 Monad transformers that use functor composition

$$\begin{array}{ccc}
 A & \xrightarrow{\text{pu}_M} & M^A \\
 & \searrow \text{pu}_T \triangleq & \downarrow \text{pu}_L \\
 & & L^{M^A}
 \end{array}$$

It remains to implement ftn_T . In the type $L^{M^{L^{M^A}}}$, we have two layers of the functor L and two layers of the functor M . We could use the available method ftn_L to flatten the two layers of L if we could *somehow* bring these nested layers together. However, these layers are separated by a layer of the functor M . To show this layered structure in a more visual way, let us employ another notation for the functor composition,

$$L \circ M \triangleq L^{M^\bullet} \quad .$$

In this notation, the type signature for flatten is written as

$$\text{ftn}_T : L \circ M \circ L \circ M \rightsquigarrow L \circ M \quad .$$

If we had $L \circ L \circ \dots$ here, we would have applied ftn_L and flattened the two layers of the functor L . Then we would have flattened the remaining two layers of the functor M . How can we achieve this? The trick is to *assume* that we can interchange the order of the layers in the middle, using a function called *swap* (short notation “sw”), with the type signature

$$\text{sw} : M \circ L \rightsquigarrow L \circ M \quad ,$$

which is equivalently written in the short type notation as

$$\text{sw} : M^{L^A} \Rightarrow L^{M^A} \quad .$$

If this operation were *somehow* defined for the two monads L and M , we could implement ftn_T by first swapping the order of the inner layers M and L as

$$L \circ M \circ L \circ M \rightsquigarrow L \circ L \circ M \circ M$$

and then applying the flatten methods of the monads L and M . The resulting code for the function ftn_T and the corresponding type diagram are

$$\text{ftn}_T \triangleq \text{sw}^{\uparrow L} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \quad . \quad (11.3)$$

11 Monad transformers

$$\begin{array}{ccccc}
 L^{M^L M^A} & \xrightarrow{\text{sw}^\uparrow L} & L^{L M^M A} & \xrightarrow{\text{ftn}_L} & L^{M^M A} \\
 & \searrow \text{ftn}_T \triangleq & & \swarrow \text{ftn}_M^\uparrow L & \\
 & & L^{M^A} & &
 \end{array}$$

It turns out that in *both* cases (the composed-inside and the composed-outside transformers), the new monad's `flatten` method can be defined through the `swap` operation. For the two kinds of transformers, the type signatures of these functions are

$$\begin{aligned}
 \text{composed-inside : } \quad & \text{ftn}_T : M^{L^{M^L A}} \Rightarrow M^{L^A} \quad , \quad \text{sw} : L^{M^A} \Rightarrow M^{L^A} \quad , \\
 \text{composed-outside : } \quad & \text{ftn}_T : L^{M^{L M^A}} \Rightarrow L^{M^A} \quad , \quad \text{sw} : M^{L^A} \Rightarrow L^{M^A} \quad .
 \end{aligned}$$

The operations `swap` and `sequence` There is a certain similarity between the `swap` operation introduced here and the `sequence` operation introduced in Chapter 9 for traversable functors. Indeed, the type signature of the `sequence` operation is

$$\text{seq} : L^{F^A} \Rightarrow F^{L^A} \quad ,$$

where F is an arbitrary applicative functor (which could be M , since monads are applicative functors) and L is a traversable functor. However, the similarity ends here. The laws required for the `swap` operation to yield a monad T are much stronger than the laws of traversable functors. In particular, if we wish M^{L^A} to be a monad, it is insufficient to require the monad L to be a traversable functor. A simple counterexample is found with $L^A \triangleq A \times A$ and $M^A \triangleq 1 + A$. Both L and M are traversable (since they are polynomial); but their composition $Q^A \triangleq 1 + A \times A$ is not a monad.

Another difference between `swap` and `sequence` is that the `swap` operation needs to be generic in the foreign monad, which may be either L or M according to the type of the monad transformer; whereas `sequence` is always generic in F . To denote more clearly the monad with respect to which `swap` is generic, we may write

$$\begin{aligned}
 \text{sw}_L^M : L^{M^A} &\Rightarrow M^{L^A} && \text{for the composed-inside transformers,} \\
 \text{sw}_L^M : M^{L^A} &\Rightarrow L^{M^A} && \text{for the composed-outside transformers.}
 \end{aligned}$$

11.5 Monad transformers that use functor composition

The superscript M in sw_L^M shows that M is a *type parameter* in swap ; that is, swap is generic in the monad M . The subscript L in sw_L^M shows that the implementations of swap may need to use completely different code for different monads L .

To avoid confusion, I use the name “ swap ” rather than “sequence” for the function $M^{L^\bullet} \leadsto L^{M^\bullet}$ in the context of monad transformers. Let us now find out what laws are required for the swap operation.¹

11.5.2 Laws for swap

The first law is that swap must be a natural transformation. Since swap has only one type parameter, there is one naturality law: for any function $f : A \Rightarrow B$,

$$f \uparrow^{L \uparrow M} \circ \text{sw} = \text{sw} \circ f \uparrow^{M \uparrow L} . \quad (11.4)$$

$$\begin{array}{ccc} M^{L^A} & \xrightarrow{f \uparrow^{L \uparrow M}} & M^{L^B} \\ \text{sw} \downarrow & & \downarrow \text{sw} \\ L^{M^A} & \xrightarrow{f \uparrow^{M \uparrow L}} & L^{M^B} \end{array}$$

To derive further laws for swap , consider the requirement that the transformed monad T should satisfy the monad laws:

$$\begin{aligned} \text{pu}_T \circ \text{ftn}_T &= \text{id} , & \text{pu}_T^{\uparrow T} \circ \text{ftn}_T &= \text{id} , \\ \text{ftn}_T^{\uparrow T} \circ \text{ftn}_T &= \text{ftn}_T \circ \text{ftn}_T . \end{aligned}$$

Additionally, T must satisfy the laws of a monad transformer. We will now discover the laws for swap that make the laws for ftn_T hold automatically, as long as ftn_T is derived from swap using Eq. (11.3).

¹The “ swap ” operation was first introduced in a [1993 paper](#) “Composing monads” by M. P. Jones and L. Duponcheel. They studied various ways of composing monads and also gave some arguments to show that no generic transformer could compose all monads L, M . The impossibility of a generic monad composition is demonstrated by the `State` monad that, as I show in this chapter, does not compose with arbitrary other monads M – either from inside or from outside.

11 Monad transformers

We substitute Eq. (11.3) into the left identity law for ftn_T and simplify:

$$\begin{aligned}
 \text{id} &= \text{pu}_T \circ \text{ftn}_T \\
 \text{replace } \text{ftn}_T \text{ using Eq. (11.3)} : &= \text{pu}_T \circ \text{sw}^{\uparrow L} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \\
 \text{replace } \text{pu}_T \text{ using Eq. (11.2)} : &= \text{pu}_M \circ \underline{\text{pu}_L \circ \text{sw}^{\uparrow L}} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \\
 \text{naturality of } \text{pu}_L : &= \text{pu}_M \circ \text{sw} \circ \underline{\text{pu}_L \circ \text{ftn}_L} \circ \text{ftn}_M^{\uparrow L} \\
 \text{left identity law for } L : &= \text{pu}_M \circ \text{sw} \circ \text{ftn}_M^{\uparrow L} \quad . \quad (11.5)
 \end{aligned}$$

How could the last expression in Eq. (11.5) be equal to id ? We know nothing about the `pure` and `flatten` methods of the monads L and M , except that these methods satisfy their monad laws. We could satisfy the law in Eq. (11.5) if we somehow reduce that expression to

$$\text{pu}_M^{\uparrow L} \circ \text{ftn}_M^{\uparrow L} = (\text{pu}_M \circ \text{ftn}_M)^{\uparrow L} = \text{id} \quad .$$

This will be possible only if we are able to interchange the order of function compositions with `sw` and eliminate `swap` from the expression. So, we must require the “ M -identity law” for `swap`,

$$\text{pu}_M \circ \text{sw} = \text{pu}_M^{\uparrow L} \quad . \quad (11.6)$$

$$\begin{array}{ccc}
 L^A & \xrightarrow{\text{pu}_M} & M^{L^A} \\
 & \searrow \text{pu}_M^{\uparrow L} & \downarrow \text{sw} \\
 & & L^{M^A}
 \end{array}$$

Intuitively, this law says that a pure layer of the monad M remains pure after interchanging the order of layers with `swap`.

With this law, we can finish the derivation in Eq. (11.5) as

$$\begin{aligned}
 &\text{pu}_M \circ \text{sw} \circ \text{ftn}_M^{\uparrow L} \\
 M\text{-identity law for } \text{sw} : &= \text{pu}_M^{\uparrow L} \circ \text{ftn}_M^{\uparrow L} \\
 \text{functor composition law for } L : &= (\text{pu}_M \circ \text{ftn}_M)^{\uparrow L} \\
 \text{left identity law for } M : &= \text{id}^{\uparrow L} \\
 \text{functor identity law for } L : &= \text{id} \quad .
 \end{aligned}$$

11.5 Monad transformers that use functor composition

So, the M -identity law for `swap` entails the left identity law for T .

In the same way, we motivate the “ L -identity” law for `swap`,

$$\text{pu}_L^{\uparrow M} \circ \text{sw} = \text{pu}_L \quad . \quad (11.7)$$

$$\begin{array}{ccc} M^A & \xrightarrow{\text{pu}_L^{\uparrow M}} & M^{L^A} \\ & \searrow \text{pu}_L & \downarrow \text{sw} \\ & & L^{M^A} \end{array}$$

This law expresses the idea that a pure layer of the functor L remains pure after swapping the order of layers.

Assuming this law, we can derive the right identity law for T :

$$\begin{aligned} & \text{pu}_T^{\uparrow T} \circ \text{ftn}_T \\ \text{(note that } f^{\uparrow T} \triangleq f^{\uparrow M \uparrow L} \text{)} : &= (\text{pu}_T)^{\uparrow M \uparrow L} \circ \text{ftn}_T \\ \text{definitions of } \text{pu}_T \text{ and } \text{ftn}_T : &= \text{pu}_M^{\uparrow M \uparrow L} \circ \underline{\text{pu}_L^{\uparrow M \uparrow L} \circ \text{sw}^{\uparrow L}} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \\ L\text{-identity law for } \text{sw}, \text{ under } \uparrow^L : &= \text{pu}_M^{\uparrow M \uparrow L} \circ \underline{\text{pu}_L^{\uparrow L} \circ \text{ftn}_L} \circ \text{ftn}_M^{\uparrow L} \\ \text{right identity law for } L : &= \text{pu}_M^{\uparrow M \uparrow L} \circ \text{ftn}_M^{\uparrow L} = \underline{(\text{pu}_M^{\uparrow M} \circ \text{ftn}_M)^{\uparrow L}} \\ \text{right identity law for } M : &= \text{id}^{\uparrow L} = \text{id} \quad . \end{aligned}$$

Deriving the monad associativity law for T ,

$$\text{ftn}_T^{\uparrow T} \circ \text{ftn}_T = \text{ftn}_T \circ \text{ftn}_T \quad ,$$

turns out to require *two* further laws for `swap`. Substituting the definition of `ftnT` into the associativity law, we get

$$\begin{aligned} & (\text{sw}^{\uparrow L} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L})^{\uparrow M \uparrow L} \circ \text{sw}^{\uparrow L} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \\ &= \text{sw}^{\uparrow L} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \circ \text{sw}^{\uparrow L} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \quad . \end{aligned} \quad (11.8)$$

The only hope of proving this law is being able to interchange `ftnL` as well as `ftnM` with `sw`. In other words, the `swap` function should be in some way adapted to the `flatten` methods of *both* monads L and M .

Let us look for such interchange laws. One possibility is to have a law involving `ftnM ∘ sw`, which is a function of type $M^{M^A} \Rightarrow L^{M^A}$ or, in another notation,

11 Monad transformers

$M \circ M \circ L \rightsquigarrow L \circ M$. This function first flattens the two adjacent layers of M , obtaining $M \circ L$, and then swaps the two remaining layers, moving the L layer outside. Let us think about what law could exist for this kind of transformation. It is plausible that we may obtain the same result if we first swap the layers *twice*, so that the L layer moves to the outside, obtaining $L \circ M \circ M$, and then flatten the two inner M layers. Writing this assumption in code, we obtain the “ M -interchange” law

$$\text{ftn}_M \circ \text{sw} = \text{sw}^{\uparrow M} \circ \text{sw} \circ \text{ftn}_M^{\uparrow L} . \quad (11.9)$$

$$\begin{array}{ccccc}
 & & M^{M^L A} & \xrightarrow{\text{ftn}_M} & M^{L A} \\
 & \swarrow \text{sw}^{\uparrow M} & & & \downarrow \text{sw} \\
 M^{L M^A} & \xrightarrow{\text{sw}} & L^{M M^A} & \xrightarrow{\text{ftn}_M^{\uparrow L}} & L^{M^A}
 \end{array}$$

The analogous “ L -interchange” law involving two layers of L and a transformation $M \circ L \circ L \rightsquigarrow L \circ M$ is written as

$$\text{ftn}_L^{\uparrow M} \circ \text{sw} = \text{sw} \circ \text{sw}^{\uparrow L} \circ \text{ftn}_L . \quad (11.10)$$

$$\begin{array}{ccccc}
 & & M^{L L A} & \xrightarrow{\text{ftn}_L^{\uparrow M}} & M^{L A} \\
 & \swarrow \text{sw} & & & \downarrow \text{sw} \\
 L^{M L A} & \xrightarrow{\text{sw}^{\uparrow L}} & L^{L M^A} & \xrightarrow{\text{ftn}_L} & L^{M^A}
 \end{array}$$

At this point, we have simply written down these two interchange laws, hoping that they will help us derive the associativity law for T . It remains to prove that this is indeed so.

Both sides of the law in Eq. (11.8) involve compositions of several `flattens` and `swaps`. The heuristic idea of the proof is to use various laws to move all `flattens` to right of the composition, while moving all `swaps` to the left. In this way we will transform both sides of Eq. (11.8) into a similar form, hoping to prove that they are equal.

We begin with the right-hand side of Eq. (11.8) since it is simpler than the left-hand side, and look for ways of using the interchange laws. At every step of the

11.5 Monad transformers that use functor composition

calculation, there happens to be only one place where some law can be applied:

$$\begin{aligned}
& \text{sw}^{\uparrow L} \circ \text{ftn}_L \circ \underline{\text{ftn}_M^{\uparrow L} \circ \text{sw}^{\uparrow L}} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \\
\text{composition for } L : &= \text{sw}^{\uparrow L} \circ \text{ftn}_L \circ \underline{(\text{ftn}_M \circ \text{sw})^{\uparrow L}} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \\
M\text{-interchange for sw} : &= \text{sw}^{\uparrow L} \circ \text{ftn}_L \circ \underline{(\text{sw}^{\uparrow M} \circ \text{sw} \circ \text{ftn}_M^{\uparrow L})^{\uparrow L}} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \\
\text{composition for } L : &= \text{sw}^{\uparrow L} \circ \text{ftn}_L \circ \underline{\text{sw}^{\uparrow M \uparrow L} \circ \text{sw}^{\uparrow L}} \circ \underline{\text{ftn}_M^{\uparrow L} \circ \text{ftn}_L} \circ \text{ftn}_M^{\uparrow L} \\
\text{naturality of } \text{ftn}_L : &= \text{sw}^{\uparrow L} \circ \text{ftn}_L \circ \underline{(\text{sw}^{\uparrow M} \circ \text{sw})^{\uparrow L}} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \circ \text{ftn}_M^{\uparrow L} \\
\text{naturality of } \text{ftn}_L : &= \text{sw}^{\uparrow L} \circ \underline{(\text{sw}^{\uparrow M} \circ \text{sw})^{\uparrow L \uparrow L}} \circ \text{ftn}_L \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \circ \text{ftn}_M^{\uparrow L} .
\end{aligned}$$

Now all `swaps` are on the left and all `flattens` on the right of the expression.

Transform the right-hand side of Eq. (11.8) in the same way as

$$\begin{aligned}
& \left(\text{sw}^{\uparrow L} \circ \text{ftn}_L \circ \underline{\text{ftn}_M^{\uparrow L}} \right)^{\uparrow M \uparrow L} \circ \underline{\text{sw}^{\uparrow L} \circ \text{ftn}_L} \circ \text{ftn}_M^{\uparrow L} \\
\text{functor composition} : &= \left(\text{sw}^{\uparrow L} \circ \text{ftn}_L \right)^{\uparrow M \uparrow L} \circ \underline{(\text{ftn}_M^{\uparrow L \uparrow M} \circ \text{sw})^{\uparrow L}} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \\
\text{naturality of sw} : &= \left(\text{sw}^{\uparrow L} \circ \text{ftn}_L \right)^{\uparrow M \uparrow L} \circ \underline{(\text{sw} \circ \text{ftn}_M^{\uparrow M \uparrow L})^{\uparrow L}} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \\
\text{naturality of } \text{ftn}_L : &= \text{sw}^{\uparrow L \uparrow M \uparrow L} \circ \text{ftn}_L^{\uparrow M \uparrow L} \circ \underline{\text{sw}^{\uparrow L} \circ \text{ftn}_L} \circ \underline{\text{ftn}_M^{\uparrow M \uparrow L} \circ \text{ftn}_M^{\uparrow L}} \\
\text{associativity of } \text{ftn}_M : &= \text{sw}^{\uparrow L \uparrow M \uparrow L} \circ \underline{(\text{ftn}_L^{\uparrow M} \circ \text{sw})^{\uparrow L}} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \circ \text{ftn}_M^{\uparrow L} \\
L\text{-interchange for sw} : &= \text{sw}^{\uparrow L \uparrow M \uparrow L} \circ (\text{sw} \circ \text{sw}^{\uparrow L} \circ \underline{\text{ftn}_L})^{\uparrow L} \circ \underline{\text{ftn}_L} \circ \text{ftn}_M^{\uparrow L} \circ \text{ftn}_M^{\uparrow L} \\
\text{associativity of } \text{ftn}_L : &= (\text{sw}^{\uparrow L \uparrow M} \circ \text{sw} \circ \text{sw}^{\uparrow L})^{\uparrow L} \circ \underline{\text{ftn}_L \circ \text{ftn}_L} \circ \text{ftn}_M^{\uparrow L} \circ \text{ftn}_M^{\uparrow L} .
\end{aligned}$$

We have again managed to move all `swaps` to the left and all `flattens` to the right of the expression.

Comparing now the two sides of the associativity law, we see that all the `flattens` occur in the same combination: $\text{ftn}_L \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \circ \text{ftn}_M^{\uparrow L}$. It remains to show that

$$\text{sw}^{\uparrow L} \circ (\text{sw}^{\uparrow M} \circ \text{sw})^{\uparrow L \uparrow L} = (\text{sw}^{\uparrow L \uparrow M} \circ \text{sw} \circ \text{sw}^{\uparrow L})^{\uparrow L} .$$

or equivalently

$$(\text{sw} \circ \text{sw}^{\uparrow M \uparrow L} \circ \text{sw}^{\uparrow L})^{\uparrow L} = (\text{sw}^{\uparrow L \uparrow M} \circ \text{sw} \circ \text{sw}^{\uparrow L})^{\uparrow L} .$$

11 Monad transformers

The two sides are equal due to the naturality law of `swap`,

$$\text{sw} \circ \text{sw}^{\uparrow M \uparrow L} = \text{sw}^{\uparrow L \uparrow M} \circ \text{sw}.$$

This completes the proof of the following theorem:

Theorem 11.5.2.1 If two monads L and M are such that there exists a natural transformation

$$\text{sw} : M^{L^A} \Rightarrow L^{M^A}$$

(called “`swap`”) satisfying the four additional laws,

$$\begin{aligned} L\text{-identity} : \quad & \text{pu}_L^{\uparrow M} \circ \text{sw} = \text{pu}_L \quad , \\ M\text{-identity} : \quad & \text{pu}_M \circ \text{sw} = \text{pu}_M^{\uparrow L} \quad , \\ L\text{-interchange} : \quad & \text{ftn}_L^{\uparrow M} \circ \text{sw} = \text{sw} \circ \text{sw}^{\uparrow L} \circ \text{ftn}_L \quad , \\ M\text{-interchange} : \quad & \text{ftn}_M \circ \text{sw} = \text{sw}^{\uparrow M} \circ \text{sw} \circ \text{ftn}_M^{\uparrow L} \quad , \end{aligned}$$

then the functor composition

$$T^A \triangleq L^{M^A}$$

is a monad with the methods `pure` and `flatten` defined by

$$\text{pu}_T \triangleq \text{pu}_M \circ \text{pu}_L \quad , \tag{11.11}$$

$$\text{ftn}_T \triangleq \text{sw}^{\uparrow L} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \quad . \tag{11.12}$$

11.5.3 Intuition behind the laws of `swap`

The interchange laws for `swap` guarantee that in any functor composition built up from L and M , e.g. like this,

$$M \circ M \circ L \circ M \circ L \circ L \circ M \circ M \circ L \quad ,$$

we can flatten the layers using `ftnL`, or `ftnM`, `ftnT`, or interchange the layers with `swap`, in any order. We will always get the same final value, which we can transform to the monad type $T = L \circ M$.

In other words, the monadic effects of the monads L and M can be arbitrarily interleaved, swapped, and flattened in any order, with no change to the final results. The programmer is free to refactor a monadic program, say, by computing

11.5 Monad transformers that use functor composition

some L -effects in a separate functor block of `L-flatMaps` and only then combining the result with the rest of the computation in the monad T . Regardless of the refactoring, the monad T computes all the effects correctly. This is what programmers would certainly expect of the monad T , if it is to be a useful monad transformer.

We will now derive the properties of ftn_T that correspond to the interchange laws. We will find that it is easier to formulate these laws in terms of `swap` than in terms of ftn_T . In practice, all known examples of compositional monad transformers (the “linear” and the “rigid” monads) are defined via `swap`.

11.5.4 Deriving `swap` from `flatten`

We have seen that the `flatten` method for the monad $T^\bullet = L^{M^\bullet}$ can be defined via the `swap` method. However, we have seen examples of some composable monads (such as `Reader` and `Option`) where we already know the definitions of the `flatten` method for T . How do we know whether a suitable `swap` function exist for these examples? In other words, if a `flatten` function for the monad $T^\bullet = L^{M^\bullet}$ is already given, can we establish whether a `swap` function exists such that the given `flatten` function is expressed via Eq. (11.3)?

To answer this question, let us look at the type signature of `flatten` for T :

$$\text{ftn}_T : L \circ M \circ L \circ M \rightsquigarrow L \circ M \quad .$$

This type signature is different from $\text{sw} : M \circ L \rightsquigarrow L \circ M$ only because the argument of ftn_T has extra layers of the functors L and M that are placed outside the $M \circ L$ composition. We can use the pure methods of M and L to add these extra layers to a value of type $M \circ L$, without modifying any monadic effects present in $M \circ L$. This will allow us to apply ftn_T and to obtain a value of type $L \circ M$. The resulting code for the function ftn_T and the corresponding type diagram are

$$\text{sw} = \text{pu}_M^{\uparrow L \uparrow M} \circ \text{pu}_L \circ \text{ftn}_T \quad . \quad (11.13)$$

We have expressed ftn_T and `sw` through each other. Are these functions always equivalent? To decide this, we need to answer two questions:

11 Monad transformers

1. If we first define ftn_T using Eq. (11.3) through a given implementation of sw , and then substitute that ftn_T into Eq. (11.13), will we always recover the initially given function sw ? (Yes.)
2. If we first define sw using Eq. (11.13) through a given implementation of ftn_T , and then substitute that sw into Eq. (11.3), will we always recover the initially given function ftn_T ? (No, not without additional assumptions for ftn_T .)

To answer the first question, substitute ftn_T from Eq. (11.3) into Eq. (11.13):

$$\begin{aligned}
 & \text{pu}_M^{\uparrow L \uparrow M} \circ \text{pu}_L \circ \text{ftn}_T \\
 \text{substitute } \text{ftn}_T : &= \text{pu}_M^{\uparrow L \uparrow M} \circ \text{pu}_L \circ \text{sw}^{\uparrow L} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \\
 \text{naturality of } \text{pu}_L : &= \text{pu}_M^{\uparrow L \uparrow M} \circ \text{sw} \circ \text{pu}_L \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \\
 \text{left identity law for } L : &= \text{pu}_M^{\uparrow L \uparrow M} \circ \text{sw} \circ \text{ftn}_M^{\uparrow L} \\
 \text{naturality of } \text{sw} : &= \text{sw} \circ \text{pu}_M^{\uparrow M \uparrow L} \circ \text{ftn}_M^{\uparrow L} \\
 \text{functor composition for } L : &= \text{sw} \circ (\text{pu}_M^{\uparrow M} \circ \text{ftn}_M)^{\uparrow L} \\
 \text{right identity law for } M : &= \text{sw} \quad .
 \end{aligned}$$

So, yes, we do always recover the initial swap function.

To answer the second question, substitute sw from Eq. (11.13) into Eq. (11.3):

$$\begin{aligned}
 & \text{sw}^{\uparrow L} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \\
 \text{substitute } \text{sw} : &= (\text{pu}_M^{\uparrow L \uparrow M} \circ \text{pu}_L \circ \text{ftn}_T)^{\uparrow L} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \\
 \text{functor composition} : &= \text{pu}_M^{\uparrow L \uparrow M \uparrow L} \circ \text{pu}_L^{\uparrow L} \circ \text{ftn}_T^{\uparrow L} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \quad . \quad (11.14)
 \end{aligned}$$

At this point, we are stuck: no laws can be applied to transform the last expression any further. Without additional assumptions, it does not follow that this expression is equal to ftn_T . Let us derive the necessary additional assumptions.

A simplification of Eq. (11.14) clearly requires a law involving the function

$$\text{ftn}_T^{\uparrow L} \circ \text{ftn}_L : L \circ L \circ M \circ L \circ M \rightsquigarrow L \circ M \quad .$$

This function flattens the two layers of $(L \circ M)$ and then flattens the remaining two layers of L . An alternative transformation with the same type signature is

11.5 Monad transformers that use functor composition

to first flatten the two *outside* layers of L and then to flatten the two remaining layers of $(L \circ M)$:

$$\text{ftn}_L \circ \text{ftn}_T : L \circ L \circ M \circ L \circ M \rightsquigarrow L \circ M \quad .$$

So we conjecture that a possibly useful additional law for ftn_T is

$$\text{ftn}_L \circ \text{ftn}_T = \text{ftn}_T^{\uparrow L} \circ \text{ftn}_L \quad .$$

$$\begin{array}{ccc} L^{L^M L^{M^A}} & \xrightarrow{\text{ftn}_L} & L^{M^L M^A} \\ \text{ftn}_T^{\uparrow L} \downarrow & & \downarrow \text{ftn}_T \\ L^{L^{M^A}} & \xrightarrow{\text{ftn}_L} & L^{M^A} \end{array}$$

This law shows a kind of “compatibility” between the monads L and T .

With this law, the right-hand side of Eq. (11.14) becomes

$$\begin{aligned} & \text{pu}_M^{\uparrow L \uparrow M \uparrow L} \circ \text{pu}_L^{\uparrow L} \circ \text{ftn}_L \circ \text{ftn}_T \circ \text{ftn}_M^{\uparrow L} \\ \text{right identity law of } L : & = \text{pu}_M^{\uparrow L \uparrow M \uparrow L} \circ \text{ftn}_T \circ \text{ftn}_M^{\uparrow L} \quad . \end{aligned}$$

Again, we cannot proceed unless we have a law involving the function

$$\text{ftn}_T \circ \text{ftn}_M^{\uparrow L} : L \circ M \circ L \circ M \circ M \rightsquigarrow L \circ M \quad .$$

This function first flattens the two layers of $(L \circ M)$ and then flattens the remaining two layers of M . An alternative order of flattenings is to first flatten the *innermost* two layers of M :

$$\text{ftn}_M^{\uparrow L \uparrow M \uparrow L} \circ \text{ftn}_T : L \circ M \circ L \circ M \circ M \rightsquigarrow L \circ M \quad .$$

The second conjectured law is therefore

$$\text{ftn}_T \circ \text{ftn}_M^{\uparrow L} = \text{ftn}_M^{\uparrow L \uparrow M \uparrow L} \circ \text{ftn}_T \quad .$$

$$\begin{array}{ccc} L^{M^L M^{M^A}} & \xrightarrow{\text{ftn}_T} & L^{M^{M^A}} \\ \text{ftn}_M^{\uparrow L \uparrow M \uparrow L} \downarrow & & \downarrow \text{ftn}_M^{\uparrow L} \\ L^{M^L M^A} & \xrightarrow{\text{ftn}_T} & L^{M^A} \end{array}$$

11 Monad transformers

This law expresses “compatibility” between the `flatten` methods of M and T .

With this law, we can finally complete the derivation:

$$\begin{aligned}
 & \text{pu}_M^{\uparrow L \uparrow M \uparrow L} \circ \text{ftn}_T \circ \text{ftn}_M^{\uparrow L} \\
 \text{substitute the second conjecture : } &= \text{pu}_M^{\uparrow L \uparrow M \uparrow L} \circ \text{ftn}_M^{\uparrow L \uparrow M \uparrow L} \circ \text{ftn}_T \\
 \text{functor composition : } &= (\text{pu}_M \circ \text{ftn}_M)^{\uparrow L \uparrow M \uparrow L} \circ \text{ftn}_T \\
 \text{left identity law of } M : &= \text{ftn}_T .
 \end{aligned}$$

We recovered the initial ftn_T by assuming two additional laws.

It turns out that these additional laws will always hold when ftn_T is defined via `swap` (see Exercise 11.5.4.1).

Given some monads L and M , it is cumbersome to verify the associativity and the interchange laws for T directly because of the deeply nested type constructors, such as $L \circ M \circ L \circ M \circ L \circ M$. If the monad T is defined via a suitable `swap` method (in practice, this is always the case), it is simpler to verify the laws of T by checking the laws of `swap`. This is what we will do next.

Exercise 11.5.4.1 Assuming that

- L and M are monads,
- the method `swap` is a natural transformation $M \circ L \rightsquigarrow L \circ M$,
- the method ftn_T of the monad $T = L \circ M$ is defined via `swap` by Eq. (11.3),

show that the two interchange laws hold for T :

$$\begin{aligned}
 \text{\textit{L}-interchange : } & \text{ftn}_L \circ \text{ftn}_T = \text{ftn}_T^{\uparrow L} \circ \text{ftn}_L , \\
 \text{\textit{M}-interchange : } & \text{ftn}_T \circ \text{ftn}_M^{\uparrow L} = \text{ftn}_M^{\uparrow L \uparrow M \uparrow L} \circ \text{ftn}_T .
 \end{aligned}$$

Exercise 11.5.4.2 With the same assumptions as Exercise 11.5.4.1 and additionally assuming the L -identity and the M -identity laws for `swap` (see Theorem 11.5.2.1), show that the monad $T^\bullet \triangleq L^{M^\bullet}$ satisfies two “pure compatibility” laws,,

$$\begin{aligned}
 \text{\textit{L-pure-compatibility : } } & \text{ftn}_L = \text{pu}_M^{\uparrow L} \circ \text{ftn}_T : L^{L^{M^A}} \Rightarrow L^{M^A} , \\
 \text{\textit{M-pure-compatibility : } } & \text{ftn}_M^{\uparrow L} = \text{pu}_L^{\uparrow T} \circ \text{ftn}_T : L^{M^{M^A}} \Rightarrow L^{M^A} ,
 \end{aligned}$$

11.5 Monad transformers that use functor composition

or, expressed equivalently through the flm methods instead of ftn ,

$$\begin{aligned}\text{flm}_L f : A \Rightarrow L^{M^B} &= \text{pure}_M^{\uparrow L} \circ \text{flm}_T f : A \Rightarrow L^{M^B} , \\ (\text{flm}_M f : A \Rightarrow M^B)^{\uparrow L} &= \text{pure}_L^{\uparrow T} \circ \text{flm}_T (f^{\uparrow L}) .\end{aligned}$$

11.5.5 Laws of monad transformer liftings

We will now derive the laws of liftings from the laws of swap , using Eqs. (11.11)–(11.12) as definitions of the methods of T .

To be specific, let us assume that L is the base monad of the transformer. Only names will need to change for the other choice of the base monad.

The lifting morphisms of a compositional monad transformer are defined by

$$\begin{aligned}\text{lift} &= \text{pu}_L : M^A \Rightarrow L^{M^A} , \\ \text{blift} &= \text{pu}_M^{\uparrow L} : L^A \Rightarrow L^{M^A} .\end{aligned}$$

Their laws of liftings (the identity and the composition laws) are

$$\begin{aligned}\text{pu}_M \circ \text{lift} &= \text{pu}_T , & \text{pu}_L \circ \text{blift} &= \text{pu}_T , \\ \text{ftn}_M \circ \text{lift} &= \text{lift}^{\uparrow M} \circ \text{lift} \circ \text{ftn}_T , & \text{ftn}_L \circ \text{blift} &= \text{blift}^{\uparrow L} \circ \text{blift} \circ \text{ftn}_T .\end{aligned}$$

The identity laws are verified quickly,

$$\begin{aligned}\text{expect to equal } \text{pu}_T : & \text{pu}_M \circ \text{lift} = \text{pu}_M \circ \text{pu}_L \\ \text{definition of } \text{pu}_T : & = \text{pu}_T , \\ \text{expect to equal } \text{pu}_T : & \text{pu}_L \circ \text{blift} = \text{pu}_L \circ \text{pu}_M^{\uparrow L} \\ \text{naturality of } \text{pu}_L : & = \text{pu}_M \circ \text{pu}_L = \text{pu}_T .\end{aligned}$$

To verify the composition laws, we need to start from their right-hand sides because the left-hand sides cannot be simplified, and we substitute the definition

11 Monad transformers

of ftn_T in terms of swap . The composition law for lift :

$$\begin{aligned}
 \text{expect to equal } \text{ftn}_M \circ \text{pu}_L &: \text{lift}^{\uparrow M} \circ \text{lift} \circ \text{ftn}_T \\
 \text{definitions of lift and } \text{ftn}_T &: = \text{pu}_L^{\uparrow M} \circ \text{pu}_L \circ \text{sw}^{\uparrow L} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \\
 \text{naturality of } \text{pu}_L &: = \text{pu}_L^{\uparrow M} \circ \text{sw} \circ \text{pu}_L \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \\
 \text{left identity law of } L &: = \text{pu}_L^{\uparrow M} \circ \text{sw} \circ \text{ftn}_M^{\uparrow L} \\
 L\text{-identity law of } \text{sw} &: = \text{pu}_L \circ \text{ftn}_M^{\uparrow L} \\
 \text{naturality of } \text{pu}_L &: = \text{ftn}_M \circ \text{pu}_L \quad .
 \end{aligned}$$

The composition law for blift :

$$\begin{aligned}
 \text{expect to equal } \text{ftn}_L \circ \text{pu}_M^{\uparrow L} &: \text{blift}^{\uparrow L} \circ \text{blift} \circ \text{ftn}_T \\
 \text{definitions of blift and } \text{ftn}_T &: = \text{pu}_M^{\uparrow L \uparrow L} \circ \text{pu}_M^{\uparrow L} \circ \text{sw}^{\uparrow L} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \\
 \text{functor composition in } L &: = \text{pu}_M^{\uparrow L \uparrow L} \circ (\text{pu}_M \circ \text{sw})^{\uparrow L} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \\
 M\text{-identity law of } \text{sw} &: = (\text{pu}_M^{\uparrow L \uparrow L} \circ \text{pu}_M^{\uparrow L \uparrow L}) \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \\
 \text{naturality of } \text{ftn}_L &: = \text{ftn}_L \circ (\text{pu}_M^{\uparrow L} \circ \text{pu}_M^{\uparrow L}) \circ \text{ftn}_M^{\uparrow L} \\
 \text{right identity law of } M &: = \text{ftn}_L \circ \text{pu}_M^{\uparrow L} \quad .
 \end{aligned}$$

So, the lifting laws for T follow from the laws of swap .

11.5.6 Laws of monad transformer runners

The laws of runners are not symmetric with respect to the base monad and the foreign monad: the runner, mrun , is generic in the foreign monad (but not in the base monad). In each case, the swap function must be monadically natural in the *foreign* monad. This law needs to be written differently, depending on the choice of the base monad. So we need to consider separately the cases when either L or M is the base monad.

If the base monad is L , the runners are

$$\begin{aligned}
 \text{mrun } \phi^{M^\bullet \rightsquigarrow N^\bullet} : L^{M^\bullet} &\rightsquigarrow L^{N^\bullet} \quad , \quad \text{mrun } \phi = \phi^{\uparrow L} \quad , \\
 \text{brun } \theta^{L^\bullet \rightsquigarrow \bullet} : L^{M^\bullet} &\rightsquigarrow M^\bullet \quad , \quad \text{brun } \theta = \theta \quad .
 \end{aligned}$$

11.5 Monad transformers that use functor composition

The monadic naturality of the `swap` function is written as

$$\text{sw}_L^M \circ \phi^{\uparrow L} = \phi \circ \text{sw}_L^N \quad .$$

$$\begin{array}{ccc} M^{L^A} & \xrightarrow{\text{sw}_L^M} & L^{M^A} \\ \phi \downarrow & & \downarrow \phi^{\uparrow L} \\ N^{L^A} & \xrightarrow{\text{sw}_L^N} & L^{N^A} \end{array}$$

The laws of runners require that `mrrun ϕ` and `brun θ` must be monadic morphisms, i.e. the identity and composition laws of monadic morphisms must hold for `mrrun ϕ` and `brun θ` :

$$\begin{aligned} \text{pu}_{L \circ M} \circ \text{mr_{run} } \phi &= \text{pu}_{L \circ N} \quad , \\ \text{ftn}_{L \circ M} \circ \text{mr_{run} } \phi &= (\text{mr_{run} } \phi)^{\uparrow M \uparrow L} \circ \text{mr_{run} } \phi \circ \text{ftn}_{L \circ N} \quad , \\ \text{pu}_{L \circ M} \circ \text{br_{un} } \theta &= \text{pu}_M \quad , \\ \text{ftn}_{L \circ M} \circ \text{br_{un} } \theta &= (\text{br_{un} } \theta)^{\uparrow M \uparrow L} \circ \text{br_{un} } \theta \circ \text{ftn}_M \quad . \end{aligned}$$

To derive these laws, we may use the identity and composition laws of monadic morphisms for ϕ and θ . We also use Eqs. (11.11)–(11.12) as definitions of the methods of T .

The first law to be shown is the identity law for `mrrun ϕ` :

$$\begin{aligned} \text{expect this to equal } \text{pu}_{L \circ N} : & \quad \text{pu}_{L \circ M} \circ \text{mr_{run} } \phi \\ \text{definitions of mr_{run} and pu}_{L \circ M} : & \quad = \text{pu}_M \circ \underline{\text{pu}_L \circ \phi^{\uparrow L}} \\ \text{naturality of pu}_L : & \quad = \underline{\text{pu}_M \circ \phi} \circ \text{pu}_L \\ \text{identity law for } \phi : & \quad = \text{pu}_N \circ \text{pu}_L \\ \text{definition of pu}_{L \circ N} : & \quad = \text{pu}_{L \circ N} \quad . \end{aligned}$$

11 Monad transformers

The next law to be shown is the composition law for $\text{mrun } \phi$:

$$\begin{aligned}
 \text{expect this to equal } \text{ftn}_T \circ \phi^{\uparrow L} : & \quad (\text{mrun } \phi)^{\uparrow M \uparrow L} \circ \text{mrun } \phi \circ \text{ftn}_{L \circ N} \\
 \text{definitions of mrun and } \text{ftn}_{L \circ N} : & \quad = \phi^{\uparrow L \uparrow M \uparrow L} \circ \phi^{\uparrow L} \circ \text{sw}_L^{N \uparrow L} \circ \text{ftn}_L \circ \text{ftn}_N^{\uparrow L} \\
 \text{monadic naturality of } \text{sw}_L : & \quad = \phi^{\uparrow L \uparrow M \uparrow L} \circ \text{sw}_L^{M \uparrow L} \circ \phi^{\uparrow L \uparrow L} \circ \text{ftn}_L \circ \text{ftn}_N^{\uparrow L} \\
 \text{naturality of } \text{sw}_L : & \quad = \text{sw}_L^{M \uparrow L} \circ \phi^{\uparrow M \uparrow L \uparrow L} \circ \phi^{\uparrow L \uparrow L} \circ \text{ftn}_L \circ \text{ftn}_N^{\uparrow L} \\
 \text{naturality of } \text{ftn}_L : & \quad = \text{sw}_L^{M \uparrow L} \circ \text{ftn}_L \circ \left(\phi^{\uparrow M} \circ \phi \circ \text{ftn}_N \right)^{\uparrow L} \\
 \text{composition law for } \phi : & \quad = \text{sw}_L^{M \uparrow L} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \circ \phi^{\uparrow L} \\
 \text{definition of } \text{ftn}_T : & \quad = \text{ftn}_T \circ \phi^{\uparrow L} \quad .
 \end{aligned}$$

11.6 Composed-inside transformers: Linear monads

11.7 Composed-outside transformers: Rigid monads

A monad R is **rigid** if its monad transformer can be defined by $T_R^{M, \bullet} \triangleq R^{M \bullet}$ for a foreign monad M .

***I do not know whether one could derive the properties 1-3 from each other. However, I will show that these properties hold for the four known rigid monad constructions, namely

$$\begin{aligned}
 \text{contrafunctor choice : } & \quad R^A \triangleq H^A \Rightarrow A \\
 \text{generalized selector : } & \quad S^A \triangleq F^{A \Rightarrow R^Q} \Rightarrow R^A \text{ if } R \text{ is rigid} \\
 \text{product : } & \quad S^A \triangleq P^A \times Q^A \text{ if both } P, Q \text{ are rigid} \\
 \text{composition : } & \quad S^A \triangleq P^{Q^A} \text{ if both } P, Q \text{ are rigid}
 \end{aligned}$$

11.7.1 Rigid monad construction 1: contrafunctor choice

The construction I call **contrafunctor choice**, $R^A \triangleq H^A \Rightarrow A$, defines a rigid monad R for *any* given contrafunctor H .

This monad's effect is to choose and return a value of type A , given a contrafunctor H that *consumes* values of type A (and presumably can check some conditions on those values). The contrafunctor H could be a constant contrafunctor $H^A \triangleq Q$, a function such as $H^A \triangleq A \Rightarrow Q$, or a more complicated contrafunctor.

Different choices of the contrafunctor H will yield simple rigid monad examples such as $R^A \triangleq 1$ (the unit monad), $R^A \triangleq A$ (the identity monad), $R^A \triangleq Z \Rightarrow A$ (the reader monad), as well as the **search² monad** $R^A \triangleq (A \Rightarrow Q) \Rightarrow A$.

The search monad represents the effect of searching for a value of type A that satisfies a condition expressed through a function of type $A \Rightarrow Q$. The simplest example of a search monad is found by setting $Q \triangleq \text{Bool}$. One may implement a function of type $(A \Rightarrow \text{Bool}) \Rightarrow A$ that *somehow* finds a value of type A that might satisfy the given predicate of type $A \Rightarrow \text{Bool}$. The intention is to return a value that satisfies the predicate if possible. If no such value can be found, some value of type A is still returned.

A closely related monad is the search-with-failure $R^A \triangleq (A \Rightarrow \text{Bool}) \Rightarrow 1 + A$, which is of the “selector” type. This (non-rigid) monad will return an empty value $1 + 0$ if no value of type A can be found that satisfies the predicate. Given a search monad, one can compute a search-with-failure monad by checking whether the value returned by the search monad does actually satisfy the predicate.

Assume that H is a contrafunctor and M is a monad, and denote for brevity

$$T^A \triangleq R^{M^A} \triangleq H^{M^A} \Rightarrow M^A \quad .$$

We need to prove that (1) a non-degenerate fi method is defined for R^A ; (2) T^A is a monad, and (3) T satisfies the compatibility laws.

Statement 2. $T^\bullet \triangleq R^{M^\bullet} \triangleq H^{M^\bullet} \Rightarrow M^\bullet$ is a monad if M is any monad and H is any contrafunctor. (If we set $M^A \triangleq A$, this will also prove that R itself is a monad.)

Proof. We need to define the monad instance for T and prove the identity and the associativity laws for T , assuming that the monad M satisfies these laws.

²See <http://math.andrej.com/2008/11/21/>

11 Monad transformers

To define the monad instance for T , it is convenient to use the Kleisli category formulation of the monad. In this formulation, we consider Kleisli morphisms of type $A \Rightarrow T^B$ and then define the Kleisli identity morphism, $\text{pure}_T : A \Rightarrow T^A$, and the Kleisli product operation \diamond_T ,

$$f:A \Rightarrow T^B \diamond_T g:B \Rightarrow T^C : A \Rightarrow T^C \quad .$$

We are then required to define the operation \diamond_T and to prove identity and associativity laws for it.

We notice that since the type constructor R is itself a function type $H^A \Rightarrow A$, the type of the Kleisli morphism $A \Rightarrow T^B$ is actually $A \Rightarrow T^B \triangleq A \Rightarrow H^{M^B} \Rightarrow M^B$. While proving the monad laws for T , we will need to use the monad laws for M (since M is an arbitrary, unknown monad). In order to use the monad laws for M , it would be helpful if we had the Kleisli morphisms for M of type $A \Rightarrow M^B$ more easily available. If we flip the curried arguments of the Kleisli morphism type $A \Rightarrow H^{M^B} \Rightarrow M^B$ and instead consider the **flipped Kleisli** morphisms of type $H^{M^B} \Rightarrow A \Rightarrow M^B$, the type $A \Rightarrow M^B$ will be easier to reason about. Since the type $A \Rightarrow H^{M^B} \Rightarrow M^B$ is equivalent to $A \Rightarrow H^{M^B} \Rightarrow M^B$, any laws we prove for the flipped Kleisli morphisms will yield the corresponding laws for the standard Kleisli morphisms. The use of flipped Kleisli morphisms makes the proof significantly shorter.

We temporarily denote by püre_T and $\tilde{\diamond}_T$ the flipped Kleisli operations:

$$\begin{aligned} \text{püre}_T : H^{M^A} \Rightarrow A \Rightarrow M^A \\ f:H^{M^B} \Rightarrow A \Rightarrow M^B \tilde{\diamond}_T g:H^{M^C} \Rightarrow B \Rightarrow M^C : H^{M^C} \Rightarrow A \Rightarrow M^C \quad . \end{aligned}$$

To define the operations püre_T and $\tilde{\diamond}_T$, we may use the methods pure_M and flm_M as well as the Kleisli product \diamond_M for the given monad M . The definitions are

$$\begin{aligned} \text{püre}_T q &= q \Rightarrow \text{pure}_M \quad (\text{the argument } q \text{ is unused}), \\ f \tilde{\diamond}_T g &= q \Rightarrow (f p) \diamond_M (g q) \quad \text{where} \\ p:H^{M^B} &= (\text{flm}_M (g q))^{\downarrow H} q \quad . \end{aligned}$$

This definition works by using the Kleisli product \diamond_M on values $f p : A \Rightarrow M^B$ and $g q : B \Rightarrow M^C$. To obtain a value $p : H^{M^B}$, we use the function $\text{flm}_M (g q) : M^B \Rightarrow M^C$ to H -contramap $q : H^{M^C}$ into $p : H^{M^B}$.

11.7 Composed-outside transformers: Rigid monads

Written as a single expression, the definition of $\tilde{\otimes}_T$ is

$$f\tilde{\otimes}_T g = q \Rightarrow f \left((\text{flm}_M (g q))^{\downarrow H} q \right) \diamond_M (g q) \quad . \quad (11.15)$$

Checking the left identity law:

$$\begin{aligned} & \text{p}\tilde{\text{u}}_T \tilde{\otimes}_T g \\ \text{definition of } \tilde{\otimes}_T : &= q \Rightarrow \text{p}\tilde{\text{u}}_T \left((\text{flm}_M (g q))^{\downarrow H} q \right) \diamond_M (g q) \\ \text{definition of } \text{p}\tilde{\text{u}}_T : &= q \Rightarrow \text{pu}_M \diamond_M g q \\ \text{left identity law for } M : &= q \Rightarrow g q \\ \text{function expansion :} &= g \end{aligned}$$

Checking the right identity law:

$$\begin{aligned} & f\tilde{\otimes}_T \text{p}\tilde{\text{u}}_T \\ \text{definition of } \tilde{\otimes}_T : &= q \Rightarrow f \left((\text{flm}_M (\text{p}\tilde{\text{u}}_T q))^{\downarrow H} q \right) \diamond_M (\text{p}\tilde{\text{u}}_T q) \\ \text{definition of } \text{p}\tilde{\text{u}}_T : &= q \Rightarrow f \left((\text{flm}_M (\text{pure}_M))^{\downarrow H} q \right) \diamond_M \text{pure}_M \\ \text{right identity law for } M : &= q \Rightarrow f \left((\text{id})^{\downarrow H} q \right) \\ \text{identity law for } H : &= q \Rightarrow f q \\ \text{function expansion :} &= f \end{aligned}$$

Checking the associativity law: $(f\tilde{\otimes}_T g) \tilde{\otimes}_T h$ must equal $f\tilde{\otimes}_T (g\tilde{\otimes}_T h)$. We have

$$\begin{aligned} & (f\tilde{\otimes}_T g) \tilde{\otimes}_T h \\ &= \left(r \Rightarrow f \left((\text{flm}_M (g r))^{\downarrow H} r \right) \diamond_M (g r) \right) \tilde{\otimes}_T h \\ &= q \Rightarrow f \left((\text{flm}_M (g r))^{\downarrow H} r \right) \diamond_M (g r) \diamond_M (h q) \quad \text{where} \\ & \quad r \triangleq (\text{flm}_M (h q))^{\downarrow H} q \quad ; \end{aligned}$$

while

$$\begin{aligned} & f\tilde{\otimes}_T (g\tilde{\otimes}_T h) \\ &= f\tilde{\otimes}_T \left(q \Rightarrow g \left((\text{flm}_M (h q))^{\downarrow H} q \right) \diamond_M (h q) \right) \\ &= q \Rightarrow f \left((\text{flm}_M k)^{\downarrow H} q \right) \diamond_M k \quad \text{where} \\ & \quad r \triangleq (\text{flm}_M (h q))^{\downarrow H} q \quad \text{and} \\ & \quad k \triangleq (g r) \diamond_M (h q) \quad . \end{aligned}$$

11 Monad transformers

It remains to show that the following two expressions are equal,

$$\begin{aligned} & f \left((\text{flm}_M (g r))^{\downarrow H} r \right) \diamond_M (g r) \diamond_M (h q) \quad \text{and} \\ & f \left((\text{flm}_M ((g r) \diamond_M (h q)))^{\downarrow H} q \right) \diamond_M (g r) \diamond_M (h q), \quad \text{where} \\ & r \triangleq (\text{flm}_M (h q))^{\downarrow H} q. \end{aligned}$$

These long expressions differ only by the following two sub-expressions, namely

$$(\text{flm}_M (g r))^{\downarrow H} r$$

and

$$(\text{flm}_M ((g r) \diamond_M (h q)))^{\downarrow H} q,$$

where $r \triangleq (\text{flm}_M (h q))^{\downarrow H} q$. Writing out the value r in the last argument of $(\text{flm}_M (g r))^{\downarrow H} r$ but leaving r unexpanded everywhere else, we now rewrite the differing sub-expressions as

$$\begin{aligned} & (\text{flm}_M (g r))^{\downarrow H} (\text{flm}_M (h q))^{\downarrow H} q \quad \text{and} \\ & (\text{flm}_M ((g r) \diamond_M (h q)))^{\downarrow H} q. \end{aligned}$$

Now it becomes apparent that we need to put the two “ flm_M ”s closer together and to combine them by using the associativity law of the monad M . Then we can rewrite the first sub-expression and transform it into the second one:

$$\begin{aligned} & (\text{flm}_M (g r))^{\downarrow H} (\text{flm}_M (h q))^{\downarrow H} q \\ \text{composition law for } H : & = (\text{flm}_M (g r) \circledast \text{flm}_M (h q))^{\downarrow H} q \\ \text{associativity law for } M : & = (\text{flm}_M ((g r) \circledast \text{flm}_M (h q)))^{\downarrow H} q \\ \text{definition of } \diamond_M \text{ via } \text{flm}_M : & = (\text{flm}_M ((g r) \diamond_M (h q)))^{\downarrow H} q. \end{aligned}$$

This proves the associativity law for \circledast_T .

11.7 Composed-outside transformers: Rigid monads

Statement 3. The monad instance for T defined in Statement 3 can be defined equivalently as

$$\begin{aligned}
 \text{pure}_T (a^{\cdot A}) &: H^{M^A} \Rightarrow M^A \\
 \text{pure}_T a &\triangleq - \Rightarrow \text{pure}_M a \\
 \text{flm}_T (f^{\cdot A \Rightarrow H^{M^B} \Rightarrow M^B}) &: (H^{M^A} \Rightarrow M^A) \Rightarrow H^{M^B} \Rightarrow M^B \\
 \text{flm}_T f &\triangleq t^{\cdot R^{M^A}} \Rightarrow q^{\cdot H^{M^B}} \Rightarrow p^{\uparrow R} t q \quad \text{where} \\
 p^{\cdot M^A \Rightarrow M^B} &\triangleq \text{flm}_M (x^{\cdot A} \Rightarrow f x q) \quad .
 \end{aligned}$$

Proof. The definition of $\tilde{\diamond}_T$ in Statement 3 used the flipped types of Kleisli morphisms, which is not the standard way of defining the methods of a monad. To restore the standard types, we need to unflip the arguments:

$$\begin{aligned}
 f^{\cdot A \Rightarrow H^{M^B} \Rightarrow M^B} \diamond_T g^{\cdot B \Rightarrow H^{M^C} \Rightarrow M^C} &: A \Rightarrow H^{M^C} \Rightarrow M^C \quad ; \\
 f \diamond_T g = t \Rightarrow q \Rightarrow &\left(\tilde{f} \left((\text{flm}_M (b \Rightarrow g b q))^{\downarrow H} q \right) \diamond_M (b \Rightarrow g b q) \right) t \quad ,
 \end{aligned}$$

where $\tilde{f} \triangleq h \Rightarrow k \Rightarrow f k h$ is the flipped version of f . To replace \diamond_M by flm_M , express $x \diamond_M y = x \circ \text{flm}_M y$ to find

$$f \diamond_T g = t \Rightarrow q \Rightarrow \left(\tilde{f} \left(p^{\downarrow H} q \right) \circ p \right) t \quad \text{where } p = \text{flm}_M (x \Rightarrow g x q) \quad .$$

To obtain an implementation of flm_T , express flm_T through \diamond_T as

$$\text{flm}_T g^{\cdot A \Rightarrow T^B} = \text{id}^{\cdot T^A \Rightarrow T^A} \diamond_T g \quad .$$

Now we need to substitute $f^{\cdot T^A \Rightarrow T^A} = \text{id}$ into $f \diamond_T g$. Noting that \tilde{f} will then become

$$\tilde{f} = (h \Rightarrow k \Rightarrow \text{id } k h) = (h \Rightarrow k \Rightarrow k h) \quad ,$$

11 Monad transformers

we get

$$\begin{aligned}
 \text{flm}_T g^{A \Rightarrow T^B} &= \text{id} \circ \text{flm}_T g \\
 \text{definition of } \diamond_T : \quad &= t^{T^A} \Rightarrow q^{H^{M^B}} \Rightarrow (\tilde{f} (p^{\downarrow H} q) \circ p) t \\
 &\quad \text{where } p = \text{flm}_M (x \Rightarrow g \ x \ q) \\
 \text{substitute } f = \text{id} : \quad &= t \Rightarrow q \Rightarrow ((h \Rightarrow k \Rightarrow k \ h) (p^{\downarrow H} q) \circ p) t \\
 \text{apply } k \text{ to } p^{\downarrow H} q : \quad &= t \Rightarrow q \Rightarrow ((k \Rightarrow k (p^{\downarrow H} q)) \circ p) t \\
 \text{definition of } \circ : \quad &= t \Rightarrow q \Rightarrow p (t (p^{\downarrow H} q)) \quad .
 \end{aligned}$$

By definition of the functor $R^A \triangleq H^A \Rightarrow A$, we lift any function $p^{A \Rightarrow B}$ as

$$\begin{aligned}
 p^{\uparrow R} : (H^A \Rightarrow A) &\Rightarrow H^B \Rightarrow B \quad , \\
 p^{\uparrow R} r^{H^A \Rightarrow A} &\triangleq p^{\downarrow H} \circ r \circ p \\
 &= q^{H^B} \Rightarrow p (r (p^{\downarrow H} q)) \quad .
 \end{aligned}$$

Finally, renaming g to f , we obtain the desired code,

$$\text{flm}_T f = t \Rightarrow q \Rightarrow p^{\uparrow R} t \ q \quad \text{where } p \triangleq \text{flm}_M (x \Rightarrow f \ x \ q) \quad .$$

Statement 4. The rigid monad $R^A \triangleq H^A \Rightarrow A$ satisfies the two compatibility laws with respect to any monad M .

Proof. Denote for brevity $T^\bullet \triangleq R^{M^\bullet}$ and rewrite the compatibility laws in terms of flm instead of ftn (since our expressions for the monad T always need to involve flm_M). Then the two laws that we need to prove are

$$\begin{aligned}
 \text{flm}_R f^{A \Rightarrow R^{M^B}} &= \text{pure}_M^{\uparrow R} \circ \text{flm}_T f \quad \text{as functions } R^A \Rightarrow R^{M^B}, \\
 (\text{flm}_M f^{A \Rightarrow M^B})^{\uparrow R} &= \text{pure}_R^{\uparrow T} \circ \text{flm}_T (f^{\uparrow R}) \quad \text{as functions } R^{M^A} \Rightarrow R^{M^B}.
 \end{aligned}$$

A definition of flm_R is obtained from flm_T by choosing the identity monad $M^A \triangleq A$ instead of an arbitrary monad M . This replaces flm_M by id :

$$\begin{aligned}
 \text{flm}_T f^{A \Rightarrow R^{M^B}} &= t^{R^{M^A}} \Rightarrow q^{H^{M^B}} \Rightarrow (\text{flm}_M (x^A \Rightarrow f \ x \ q))^{\uparrow R} t \ q \quad ; \\
 \text{flm}_R f^{A \Rightarrow R^B} &= r^{R^A} \Rightarrow q^{H^B} \Rightarrow (x^A \Rightarrow f \ x \ q)^{\uparrow R} r \ q \quad .
 \end{aligned}$$

11.7 Composed-outside transformers: Rigid monads

To prove the first compatibility law, rewrite its right-hand side as

$$\begin{aligned}
 & \text{pure}_M^{\uparrow R} \circ \text{flm}_T f \\
 \text{definition of } \text{flm}_T : &= (r^{R^A} \Rightarrow \text{pure}_M^{\uparrow R} r) \circ (r \Rightarrow q \Rightarrow (\text{flm}_M (x \Rightarrow f x q))^{\uparrow R} r q) \\
 \text{expand } \circ \text{ and simplify : } &= r \Rightarrow q \Rightarrow (\text{flm}_M (x \Rightarrow f x q))^{\uparrow R} (\text{pure}_M^{\uparrow R} r) q \\
 \text{composition law for } R : &= r \Rightarrow q \Rightarrow (\text{pure}_M^{\uparrow R} \circ \text{flm}_M (x \Rightarrow f x q))^{\uparrow R} r q \\
 \text{left identity law for } M : &= r \Rightarrow q \Rightarrow (x \Rightarrow f x q)^{\uparrow R} r q \\
 \text{definition of } \text{flm}_R : &= \text{flm}_R f \quad .
 \end{aligned}$$

To prove the second compatibility law, rewrite its right-hand side as

$$\begin{aligned}
 & \text{pure}_R^{\uparrow T} \circ \text{flm}_T (f^{\uparrow R}) \\
 \text{definition of } \text{flm}_T (f^{\uparrow R}) : &= (t^{T^A} \Rightarrow \text{pure}_R^{\uparrow T} t) \circ (t \Rightarrow q \Rightarrow (\text{flm}_M (x \Rightarrow f^{\uparrow R} x q))^{\uparrow R} t q) \\
 \text{expand } \circ \text{ and simplify : } &= t^{T^A} \Rightarrow q \Rightarrow (\text{flm}_M (x^{R^A} \Rightarrow f^{\uparrow R} x q))^{\uparrow R} (\text{pure}_R^{\uparrow M \uparrow R} t) q \\
 \text{composition law for } R : &= t \Rightarrow q \Rightarrow (\text{pure}_R^{\uparrow M} \circ \text{flm}_M (x^{R^A} \Rightarrow f^{\uparrow R} x q))^{\uparrow R} t q \\
 \text{left naturality of } \text{flm}_M : &= t \Rightarrow q \Rightarrow (\text{flm}_M (x^A \Rightarrow (\text{pure}_R^{\uparrow R} f^{\uparrow R}) x q))^{\uparrow R} t q \\
 \text{naturality of } \text{pure}_R : &= t \Rightarrow q \Rightarrow (\text{flm}_M (x^A \Rightarrow \text{pure}_R (f x) q))^{\uparrow R} t q \\
 \text{definition of } \text{pure}_R : &= t \Rightarrow q \Rightarrow (\text{flm}_M (x^A \Rightarrow f x))^{\uparrow R} t q \\
 \text{unapply } f \text{ and } \text{flm}_M f : &= t \Rightarrow q \Rightarrow (\text{flm}_M f)^{\uparrow R} t q = (\text{flm}_M f)^{\uparrow R} \quad .
 \end{aligned}$$

Statement 5. Assume that R and M are arbitrary (not necessarily rigid) monads such that the functor composition $T^\bullet \triangleq R^{M^\bullet}$ is also a monad whose method pure_T is defined by

$$\text{pure}_{RM} = \text{pure}_M \circ \text{pure}_R \quad .$$

Then the two “liftings” and the two “runners” defined by

$$\begin{aligned}
 \text{lift} : M^\bullet &\rightsquigarrow R^{M^\bullet} \triangleq \text{pure}_R \quad , \\
 \text{blift} : R^\bullet &\rightsquigarrow R^{M^\bullet} \triangleq \text{pure}_M^{\uparrow R} \quad , \\
 \text{mrun} (\phi^{M^\bullet \rightsquigarrow N^\bullet}) : R^{M^\bullet} &\Rightarrow R^{N^\bullet} \triangleq \phi^{\uparrow R} \quad , \\
 \text{brun} \theta^{R^\bullet \rightsquigarrow \bullet} : R^{M^\bullet} &\Rightarrow M^\bullet \triangleq \theta^{R^{M^\bullet} \rightsquigarrow M^\bullet} \quad ,
 \end{aligned}$$

11 Monad transformers

will satisfy the identity laws

$$\begin{aligned}\text{pure}_M \circ \text{lift} &= \text{pure}_T, \\ \text{pure}_R \circ \text{blift} &= \text{pure}_T, \\ \text{pure}_{RM} \circ (\text{mrun } \phi) &= \text{pure}_N \circ \text{pure}_R, \\ \text{pure}_{RM} \circ (\text{brun } \theta) &= \text{pure}_M,\end{aligned}$$

where we assume that $\phi : M^\bullet \leadsto N^\bullet$ and $\theta : M^\bullet \leadsto \bullet$ are monadic morphisms, and N is an arbitrary monad.

Proof. The first law is a direct consequence of the definition of lift . The second law follows from the naturality of pure_R ,

$$\text{pure}_R \circ \text{pure}_M^{\uparrow R} = \text{pure}_M \circ \text{pure}_R.$$

The third law follows from the identity law for ϕ ,

$$\text{pure}_M \circ \phi = \text{pure}_N,$$

if we write

$$\begin{aligned}\text{definitions of } \text{pure}_{RM} \text{ and } \text{mrun} : & \quad \text{pure}_{RM} \circ (\text{mrun } \phi) \\ \text{naturality of } \text{pure}_R : & \quad = \text{pure}_M \circ \text{pure}_R \circ \phi^{\uparrow R} \\ \text{identity law for } \phi : & \quad = \text{pure}_M \circ \phi \circ \text{pure}_R \\ \text{identity law for } \phi : & \quad = \text{pure}_N \circ \text{pure}_R.\end{aligned}$$

The fourth law follows from the identity law for θ ,

$$\text{pure}_R \circ \theta = \text{id},$$

if we write

$$\begin{aligned}\text{definitions of } \text{pure}_{RM} \text{ and } \text{brun} : & \quad \text{pure}_{RM} \circ (\text{brun } \theta) \\ \text{identity law for } \theta : & \quad = \text{pure}_M \circ \text{pure}_R \circ \theta \\ \text{identity law for } \theta : & \quad = \text{pure}_M.\end{aligned}$$

In this proof, we have not actually used the assumption that R , M , R^{M^\bullet} , and N are monads. We only used the naturality of pure_R . So, it is sufficient to assume that these functors are pointed, and that ϕ and θ are morphisms of pointed functors. However, the four laws we just derived are mainly useful as identity laws for monad transformer methods.

11.7 Composed-outside transformers: Rigid monads

Statement 6. A monad transformer for the rigid monad $R^A \triangleq H^A \Rightarrow A$ is $T_R^{M,A} \triangleq R^{M^A}$, with the four required methods defined as

$$\begin{aligned} \text{lift}^{M^\bullet \rightsquigarrow T^\bullet} &\triangleq \text{pure}_R \quad , \\ \text{blift}^{R^\bullet \rightsquigarrow T^\bullet} &\triangleq \text{pure}_M^{\uparrow R} \quad , \\ \text{mrun } \phi^{M^\bullet \rightsquigarrow N^\bullet} &\triangleq \phi^{\uparrow R} \quad , \\ \text{brun } \theta^{R^\bullet \rightsquigarrow \bullet} &\triangleq \theta^{R^{M^\bullet} \rightsquigarrow M^\bullet} \quad . \end{aligned}$$

Each of these four methods is a monadic morphism for any monadic morphism $\phi : M^\bullet \rightsquigarrow N^\bullet$, where M and N are arbitrary (not necessarily rigid) monads, and for any monadic morphism $\theta : R^A \rightsquigarrow A$.

Proof. Since blift is a special case of mrun, we only need to prove the monadic morphism laws for the three methods: lift, mrun, and brun.

The identity laws follow from Statement 5, since $T_R^{M,\bullet} = R^{M^\bullet}$ is a functor composition. It remains to verify the composition laws,

$$\begin{aligned} (f^{A \Rightarrow M^B} \circ \text{lift}) \diamond_{RM} (g^{B \Rightarrow M^C} \circ \text{lift}) &= (f \diamond_M g) \circ \text{lift} \quad , \\ (f^{A \Rightarrow R^{M^B}} \circ \text{mrun } \phi) \diamond_{RN} (g^{B \Rightarrow R^{M^C}} \circ \text{mrun } \phi) &= (f \diamond_{RM} g) \circ \text{mrun } \phi \quad , \\ (f^{A \Rightarrow R^{M^B}} \circ \text{brun } \theta) \diamond_M (g^{B \Rightarrow R^{M^C}} \circ \text{brun } \theta) &= (f \diamond_{RM} g) \circ \text{brun } \theta \quad . \end{aligned}$$

For the first law, we need to use the definition of lift $\triangleq \text{pure}_R$, which is

$$\text{pure}_R x^{A} \triangleq _ \cdot H^A \Rightarrow x \quad .$$

So the definition of lift can be written as

$$f \circ \text{lift} = x^{A} \Rightarrow _ \cdot H^{M^B} \Rightarrow f x$$

as a function that ignores its second argument (of type H^{M^B}). It is convenient to

11 Monad transformers

use the flipped Kleisli product $\tilde{\circ}_{RM}$ defined before in eq. (11.15). We compute

$$\begin{aligned}
 & (f:A \Rightarrow M^B \circ \text{lift}) \tilde{\circ}_{RM} (g:B \Rightarrow M^C \circ \text{lift}) \\
 \text{definition of lift : } & = (_ \Rightarrow f) \tilde{\circ}_{RM} (_ \Rightarrow g) \\
 \text{eq. (11.15) : } & = q^{H^{M^B}} \Rightarrow (_ \Rightarrow f) \left((\text{flm}_M ((_ \Rightarrow g) \ q)) \downarrow^H q \right) \diamond_M ((_ \Rightarrow g) \) \\
 \text{expanding the functions : } & = q \Rightarrow f \diamond_M g \\
 \text{unflipping the definition of lift : } & = q \Rightarrow x \Rightarrow ((f \diamond_M g) \circ \text{lift}) \ a \ x \ .
 \end{aligned}$$

Note that the last function, of type $H^{M^B} \Rightarrow A \Rightarrow M^B$, ignores its first argument. Therefore, after unflipping this Kleisli function, it will ignore its second argument. This is precisely what the function $(f \diamond_M g) \circ \text{lift}$ must do.

For the second and third laws, we need to use the composition laws for ϕ and θ , which can be written as

$$\begin{aligned}
 (f:A \Rightarrow M^B \circ \phi) \diamond_N (g:B \Rightarrow M^C \circ \phi) &= (f \diamond_M g) \circ \phi \ , \\
 (f:A \Rightarrow M^B \circ \theta) \circ (g:B \Rightarrow M^C \circ \theta) &= (f \diamond_M g) \circ \theta \ .
 \end{aligned}$$

11.7.2 Rigid monad construction 2: composition

Statement 1. The composition $R_1^{R_2^\bullet}$ is a rigid functor if both R_1 and R_2 are rigid functors.

Proof. We will show that the non-degeneracy law $\text{fi}_T \circ \text{fo}_T = \text{id}$ holds for the rigid functor $T^\bullet \triangleq R^{M^\bullet}$ as long as M is rigid.

Since it is given that M is rigid, we may use its method fi_M satisfying the non-degeneracy law $\text{fo}_M (\text{fi}_M f) = f$.

Flip the curried arguments of the function type $A \Rightarrow T^B \triangleq A \Rightarrow H^{M^B} \Rightarrow M^B$, to obtain $H^{M^B} \Rightarrow A \Rightarrow M^B$, and note that $A \Rightarrow M^B$ can be mapped to $M^{A \Rightarrow B}$ using fi_M . So we can implement $\tilde{\text{fi}}_T$ using fi_M :

$$\begin{aligned}
 \tilde{\text{fi}}_T : (H^{M^B} \Rightarrow A \Rightarrow M^B) &\Rightarrow H^{M^{A \Rightarrow B}} \Rightarrow M^{A \Rightarrow B} \\
 \tilde{\text{fi}}_T &= f \Rightarrow h \Rightarrow \text{fi}_M \left(f \left((b \Rightarrow _ \Rightarrow b) \uparrow^{M \downarrow H} h \right) \right) \\
 \tilde{\text{fo}}_T : (H^{M^{A \Rightarrow B}} \Rightarrow M^{A \Rightarrow B}) &\Rightarrow H^{M^B} \Rightarrow A \Rightarrow M^B \\
 \tilde{\text{fo}}_T &= g \Rightarrow h \Rightarrow a \Rightarrow \text{fo}_M \left(g \left((p:A \Rightarrow B \Rightarrow p \ a) \uparrow^{M \downarrow H} h \right) \right) a
 \end{aligned}$$

11.7 Composed-outside transformers: Rigid monads

To show the non-degeneracy law for T , compute

$$\begin{aligned}
 & \tilde{f}o_T \left(\tilde{f}i_T f \right) h^{H^{M^B}} a^A \\
 \text{insert the definition of } \tilde{f}o_T : &= fo_M \left(\left(\tilde{f}i_T f \right) \left((p \Rightarrow p a)^{\uparrow M \downarrow H} h \right) \right) a \\
 \text{insert the definition of } \tilde{f}i_T : &= fo_M \left(fi_M \left(f \left((b \Rightarrow _ \Rightarrow b)^{\uparrow M \downarrow H} (p \Rightarrow p a)^{\uparrow M \downarrow H} h \right) \right) \right) a \\
 \text{nondegeneracy law for } fi_M : &= f \left((b \Rightarrow _ \Rightarrow b)^{\uparrow M \downarrow H} (p \Rightarrow p a)^{\uparrow M \downarrow H} h \right) a \\
 \text{composition laws for } M, H : &= f \left((b \Rightarrow _ \Rightarrow b \circ p \Rightarrow p a)^{\uparrow M \downarrow H} h \right) a \\
 \text{simplify :} &= f \left((b \Rightarrow b)^{\uparrow M \downarrow H} h \right) a \\
 \text{identity laws for } M, H : &= f h a \quad .
 \end{aligned}$$

We obtained $\tilde{f}o_T \left(\tilde{f}i_T f \right) h a = f h a$. Therefore the non-degeneracy law $\tilde{f}i_T \circ \tilde{f}o_T = \text{id}$ holds.

Statement 2. The composition $R_1^{R_2^\bullet}$ has property 2 of a rigid monad if both R_1 and R_2 have that property.

Proof. To prove property 2 for $R_1^{R_2^\bullet}$, first consider that $R_1^{R_2^\bullet}$ is itself a monad, which is property 2 of the rigid monad R_1 . Now consider an arbitrary monad M and the type constructor $Q^\bullet \triangleq R_1^{R_2^{M^\bullet}}$. Property 2 of rigid monads requires that Q be a monad. Since R_2 is rigid, $R_2^{M^\bullet}$ is a monad, and since R_1 is rigid, Q is a monad.

To prove property 3 for $R_1^{R_2^\bullet}$, assume that this property already holds for R_1 and R_2 separately, and consider an arbitrary monad M and the type constructor $Q^\bullet \triangleq R_1^{R_2^{M^\bullet}}$. ???

11.7.3 Rigid monad construction 3: product

The product of rigid monads, $R_1^A \times R_2^A$, is a rigid monad.

11.7.4 Rigid monad construction 4: selector monad

The **selector monad** $S^A \triangleq F^A \Rightarrow R^Q \Rightarrow R^A$ is rigid if R^\bullet is a rigid monad, F^\bullet is any functor, and Q is any fixed type.

11.8 Irregular monad transformers

11.8.1 The state monad transformer StateT

11.8.2 The continuation monad transformer ContT

11.8.3 The codensity monad transformer CodT

The **codensity monad** over a functor F is defined as

$$\text{Cod}^{F,A} \triangleq \forall B. (A \Rightarrow F^B) \Rightarrow F^B$$

Properties:

- $\text{Cod}^{F,\bullet}$ is a monad for any functor F
- If F^\bullet is itself a monad then we have monadic morphisms $\text{inC} : F^\bullet \leadsto \text{Cod}^{F,\bullet}$ and $\text{outC} : \text{Cod}^{F,\bullet} \leadsto F^\bullet$ such that $\text{inC} \circ \text{outC} = \text{id}$
- A monad transformer for the codensity monad is

$$T_{\text{Cod}}^{M,A} = \forall B. (A \Rightarrow M^{F^B}) \Rightarrow M^{F^B}$$

However, this transformer does not have the base lifting morphism

$$\text{blift} : (\forall B. (A \Rightarrow F^B) \Rightarrow F^B) \Rightarrow \forall B. (A \Rightarrow M^{F^B}) \Rightarrow M^{F^B}$$

since this type signature cannot be implemented. The codensity transformer also does not have any of the required “runner” transformations mrun and brun ,

$$\begin{aligned} \text{mrun} : (M^\bullet \leadsto N^\bullet) &\Rightarrow (\forall B. (A \Rightarrow M^{F^B}) \Rightarrow M^{F^B}) \Rightarrow \forall B. (A \Rightarrow N^{F^B}) \Rightarrow N^{F^B} , \\ \text{brun} : ((\forall B. (A \Rightarrow F^B) \Rightarrow F^B) \Rightarrow A) &\Rightarrow (\forall B. (A \Rightarrow M^{F^B}) \Rightarrow M^{F^B}) \Rightarrow M^A . \end{aligned}$$

11.9 Discussion

11.10 Rigid functors

The properties of rigid monads can be extended to a slightly larger class of rigid functors.

R is a **rigid functor** if there exists a natural transformation `fuseIn` (short notation “ fi_R ”) with the type signature

$$\text{fi}_R : (A \Rightarrow R^B) \Rightarrow R^{A \Rightarrow B} \quad .$$

Note that any functor admits the natural transformation `fuseOut` (short notation fo_R), defined by

$$\begin{aligned} \text{fo} : R^{A \Rightarrow B} &\Rightarrow A \Rightarrow R^B \quad , \\ \text{fo}(r) &= a \Rightarrow (f^{A \Rightarrow B} \Rightarrow f(a))^{\uparrow R} r \quad . \end{aligned}$$

The method `fuseIn` must satisfy the nondegeneracy law

$$\text{fi}_R \circ \text{fo}_R = \text{id} \quad .$$

Statement 1. The nondegeneracy law $\text{fi} \circ \text{fo} = \text{id}$ holds for the rigid functor $R^A \triangleq H^A \Rightarrow A$ where H^A is any contrafunctor.

Proof. The `fi` method for R is defined by

$$\begin{aligned} \text{fi} : (A \Rightarrow H^B \Rightarrow B) &\Rightarrow H^{A \Rightarrow B} \Rightarrow A \Rightarrow B \\ \text{fi} &= f^{A \Rightarrow H^B \Rightarrow B} \Rightarrow h^{H^{A \Rightarrow B}} \Rightarrow a \Rightarrow f(a)((b \Rightarrow _ \Rightarrow b)^{\downarrow H} h) \end{aligned}$$

It is easier to manipulate the methods `fi` and `fo` if we flip the curried arguments of the function type $A \Rightarrow R^B \triangleq A \Rightarrow H^B \Rightarrow B$ and instead consider the equivalent methods $\tilde{\text{fi}}$ and $\tilde{\text{fo}}$ defined by

$$\begin{aligned} \tilde{\text{fi}} : (H^B \Rightarrow A \Rightarrow B) &\Rightarrow H^{A \Rightarrow B} \Rightarrow A \Rightarrow B \\ \tilde{\text{fi}} &= f^{H^B \Rightarrow A \Rightarrow B} \Rightarrow h^{H^{A \Rightarrow B}} \Rightarrow f((b \Rightarrow _ \Rightarrow b)^{\downarrow H} h) \\ \tilde{\text{fo}} : (H^{A \Rightarrow B} \Rightarrow A \Rightarrow B) &\Rightarrow H^B \Rightarrow A \Rightarrow B \\ \tilde{\text{fo}} &= g^{H^{A \Rightarrow B} \Rightarrow A \Rightarrow B} \Rightarrow h^{H^B} \Rightarrow a^A \Rightarrow g((p^{A \Rightarrow B} \Rightarrow p a)^{\downarrow H} h) a \end{aligned}$$

11 Monad transformers

To show the non-degeneracy law for R , compute

$$\begin{aligned}
 & \tilde{f}o \left(\tilde{f}i f \right) h \cdot^{H^B} a \cdot^A \\
 \text{definition of } \tilde{f}o: &= \left(\tilde{f}i f \right) \left((p \Rightarrow p a)^{\downarrow H} h \right) a \\
 \text{definition of } \tilde{f}i: &= f \left((b \Rightarrow _ \Rightarrow b)^{\downarrow H} (p \Rightarrow p a)^{\downarrow H} h \right) a \\
 \text{composition law for } H: &= f \left((b \Rightarrow _ \Rightarrow b \circ p \Rightarrow p a)^{\downarrow H} h \right) a \\
 \text{simplify } b \Rightarrow _ \Rightarrow b \circ p \Rightarrow p a \text{ to :} &= f \left((b \Rightarrow b)^{\downarrow H} h \right) a \\
 \text{identity law for } H: &= f h a \quad .
 \end{aligned}$$

We obtained $\tilde{f}o \left(\tilde{f}i f \right) h a = f h a$, so the non-degeneracy law $\tilde{f}i \circ \tilde{f}o = \text{id}$ holds.

12 Recursive types

12.1 Fixpoints and type recursion schemes

12.2 Row polymorphism and OO programming

12.3 Column polymorphism

12.4 Discussion

13 Co-inductive typeclasses. Comonads

13.1 Practical use

13.2 Laws and structure

13.3 Co-free constructions

13.4 Co-free comonads

13.5 Comonad transformers

13.6 Discussion

14 Irregular type classes*

14.1 Distributive functors

14.2 Lenses and prisms

14.3 Discussion

15 Value-dependent types*

16 Conclusions and discussion

17 Essay: Software engineers and software artisans

Published on 2018-07-23

Let us look at the differences between the kind of activities we ordinarily call engineering, as opposed to artisanship or craftsmanship. It will then become apparent that today's computer programmers are better understood as "software artisans" rather than software engineers.

17.1 Engineering disciplines

Consider what kinds of process a mechanical engineer, a chemical engineer, or an electrical engineer follows in their work, and what kind of studies they require for proficiency in their work.

A mechanical engineer **studies** calculus, linear algebra, differential geometry, and several areas of physics such as theoretical mechanics, thermodynamics, and elasticity theory, and then uses calculations to guide the design of a bridge, say. A chemical engineer **studies** chemistry, thermodynamics, calculus, linear algebra, differential equations, some areas of physics such as thermodynamics and kinetic theory, and uses calculations to guide the design of a chemical process, say. An electrical engineer **studies** advanced calculus, linear algebra, as well as several areas of physics such as electrodynamics and quantum physics, and uses calculations to guide the design of an antenna or a microchip.

The pattern here is that an engineer uses mathematics and natural sciences in order to design new devices. Mathematical calculations and scientific reasoning are required *before* drawing a design, let alone building a real device or machine.

Some of the studies required for engineers include arcane abstract concepts such as a "**rank-4 elasticity tensor**" (used in calculations of elasticity of materials), "**Lagrangian with non-holonomic constraints**" (used in robotics), the "Gibbs free energy" (for **chemical reactor design**), or the "**Fourier transform of the delta function**" and the "**inverse Z-transform**" (for digital signal processing).

To be sure, a significant part of what engineers do is not covered by any theory: the *know-how*, the informal reasoning, the traditional knowledge passed on from expert to novice, — all those skills that are hard to formalize. Nevertheless, engineering is crucially based on natural science and mathematics for some of its decision-making about new designs.

17.2 Artisanship: Trades and crafts

Now consider what kinds of things shoemakers, plumbers, or home painters do, and what they have to learn in order to become proficient in their profession.

A novice shoemaker, for example, would begin by **copying some drawings** and then cutting leather in a home workshop. Apprenticeships proceed via learning by doing while listening to comments and instructions from an expert. After a few years of apprenticeship (for example, a **painter apprenticeship in California** can be as short as 2 years), a new specialist is ready to start productive work.

All these trades operate entirely from tradition and practical experience. The trades do not require any academic study because there is no formal theory from which to proceed. To be sure, there is *a lot* to learn in the crafts, and it takes a large amount of effort to become a good artisan in any profession. But there are no rank-4 tensors to calculate, nor any differential equations to solve; no Fourier transforms to apply to delta functions, and no Lagrangians to check for non-holonomic constraints.

Artisans do not study any formal science or mathematics because their professions do not make use of any *formal computation* for guiding their designs or processes.

17.3 Programmers today are artisans, not engineers

Now I will argue that programmers are *not engineers* in the sense we normally see the engineering professions.

17.3.1 No requirement of formal study

According to this recent Stack Overflow survey, **about half of the programmers do not have a degree in Computer Science**. I am one myself; my degrees are in

17.3 Programmers today are artisans, not engineers

physics, and I have never formally studied computer science. I took no academic courses in algorithms, data structures, computer networks, compilers, programming languages, or any other topics ordinarily included in the academic study of “computer science”. None of the courses I took at university or at graduate school were geared towards programming. I am a completely self-taught software developer.

There is a large number of successful programmers who *never* studied at a college, or perhaps never studied formally in any sense. They acquired all their knowledge and skills through self-study and practical work. Robert C. Martin is one such prominent example; an outspoken guru in the arts of programming who has seen it all, he routinely refers to programmers as artisans and uses the appropriate imagery: novices, trade and craft, the “honor of the guild”, etc. He compares programmers to plumbers, electricians, lawyers, and surgeons, but not to mathematicians, physicists, or engineers of any kind. According to one of his blog posts, he started working at age 17 as a self-taught programmer, and then went on to more jobs in the software industry; he never mentions going to college. It is clear that R. C. Martin is an expert craftsman, and that he did not need academic study to master his craft.

Here is another opinion (emphasis is theirs):

Software Engineering is unique among the STEM careers in that it absolutely does *not* require a college degree to be successful. It most certainly does not require licensing or certification. *It requires experience.*

This is a description that fits a career in crafts—but certainly not a career, say, in electrical engineering.

The high demand for software developers gave rise to “developer boot camps”—vocational schools that prepare new programmers very quickly, with no formal theory or mathematics involved, through purely practical training. These vocational schools are successful in job placement. But it is unimaginable that a 6-month crash course or even a 2-year vocational school could prepare an engineer to work successfully on designing, say, quantum analog computers, without ever having studied quantum physics or calculus.

17.3.2 No mathematical formalism to guide software development

Most books on software engineering contain no formulas or equations, no mathematical derivations of any results, and no precise definitions of the various technical terms they are using (such as “object-oriented” or “software architecture”). Some books on software engineering even have no program code in them—just words and illustrative diagrams. These books talk about how programmers should approach their job, how to organize the work flow and the code architecture, in vague and general terms: “code is about detail”, “you must never abandon the big picture”, “you should avoid tight coupling in your modules”, “a class must serve a single responsibility”, and so on. Practitioners such as R. C. Martin never studied any formalisms and do not think in terms of formalisms; instead they think in **vaguely formulated, heuristic “principles”**.

In contrast, every textbook on mechanical engineering or electrical engineering has a significant amount of advanced mathematics in it. The design of a microwave antenna **is guided** not by the principle of “serving a single responsibility” but by calculations of wave propagation, based on electrodynamics.

Donald Knuth’s classic textbook is called “*The Art of Programming*”. It is full of tips and tricks about how to program; but it does not provide any formal theory that could guide programmers while actually *writing* programs. There is nothing in that book that would be similar to the way mathematical formalism guides designs in electrical or mechanical engineering. If Knuth’s books were based on such formalism, they would have looked quite differently: some theory would be first explained and then applied to help us write code.

Knuth’s books provide many algorithms, including mathematical ones. But algorithms are similar to patented inventions: They can be used immediately without further study. Understanding an algorithm is not similar to understanding a mathematical theory. Knowing one algorithm does not make it easier to develop another algorithm in an unrelated domain. In comparison, knowing how to solve differential equations will be applicable to thousands of different areas of science and engineering.

A book exists with the title “**Science of Programming**”, but the title is misleading. The author does not propose a science, similar to physics, at the foundation of the process of designing programs, similarly to how calculations in quantum physics predict the properties of a quantum device. The book claims to give precise methods that guide programmers in writing code, but the scope of proposed methods is narrow: the design of simple algorithms for iterative manipulation

17.3 Programmers today are artisans, not engineers

of data. The procedure suggested in that book is far from a formal mathematical *derivation* of programs from specification. (A book with that title also exists, and similarly disappoints.) Programmers today are mostly oblivious to these books and do not use the methods explained there.

Standard computer science courses today do not teach a true *engineering* approach to software construction. They do teach analysis of programs using formal mathematical methods; the main such methods are **complexity analysis** (the so-called “big O notation”), and **formal verification**. But programs are analyzed only *after* they are complete. Theory does not guide the actual *process* of writing code, does not suggest good ways of organizing the code (e.g. which classes or functions or modules should be defined), and does not tell programmers which data structures or APIs would be best to implement. Programmers make these design decisions purely on the basis of experience and intuition, trial-and-error, copy-paste, and guesswork.

The theory of program analysis and verification is analogous to writing a mathematical equation for the surface of a shoe made by a fashion designer. True, the “shoe surface equations” are mathematically unambiguous and can be “analyzed” or “verified”; but the equations are written after the fact and do not guide the fashion designers in actually making shoes. It is understandable that fashion designers do not study the mathematical theory of surfaces.

17.3.3 Programmers avoid academic terminology

Programmers appear to be taken aback by such terminology as “*functor*”, “*monad*”, or “*lambda-functions*”.

Those fancy words used by functional programmers purists really annoy me. Monads, Functors... Nonsense!!!

In my experience, only a tiny minority of programmers actually complain about this. The vast majority has never heard these words and are unaware of functors or monads.

However, chemical engineers do not wince at “phase diagram” or “Gibbs free energy”, and apparently accept the need for studying differential equations. Electrical engineers do not complain that the word “Fourier” is foreign and difficult to spell, or that “delta-function” is such a weird thing to say. Mechanical engineers take it for granted that they need to calculate with “tensors” and “Lagrangians” and “non-holonomic constraints”. Actually, it seems that the arcane

terminology is the least of their difficulties! Their textbooks are full of complicated equations and long, difficult derivations.

Similarly, software engineers would not complain about the word “functor”, or about having to study the derivation of the algebraic laws for “monads,”—if they were actually *engineers*. True software engineers’ textbooks would be full of equations and derivations, which would be used to perform calculations required *before* starting to write code.

17.4 Towards software engineering

It is now clear that we do not presently have true software engineering. The people employed under that job title are actually artisans. They employ artisanal methods for their work, and their culture and work process are that of a crafts guild.

One could point out that numerical simulations required for physics or the matrix calculations required for machine learning are mathematical. True, these programming *tasks* are mathematical in nature and require formal theory to be *formulated*. However, mathematical subject matter (aerospace control, physics or astronomy experiments, mathematical statistics) does not automatically make the *process of programming* into engineering. Data scientists, aerospace engineers, and natural scientists all write code nowadays – and they are all working as artisans when they write code.

True software engineering would be achieved if we had theory that guides and informs our process of creating programs,—not theory that describes or analyzes programs *after* they are somehow written.

We expect that software engineers’ textbooks should be full of equations. What theory should those equations represent?

I believe this theory already exists, and I call it **functional type theory**. It is the algebraic foundation of the modern practice of functional programming, as implemented in languages such as OCaml, Haskell, and Scala. This theory is a blend of type theory, category theory, and logical proof theory. It has been in development since late 1990s and is still being actively worked on by a community of academic computer scientists and advanced software practitioners.

To appreciate that functional programming, unlike any other programming paradigm, *has a theory that guides coding*, we can look at some recent software engineering conferences such as [Scala By the Bay](#) or [BayHac](#), or at the numerous FP-related online tutorials and blogs. We cannot fail to notice that much time

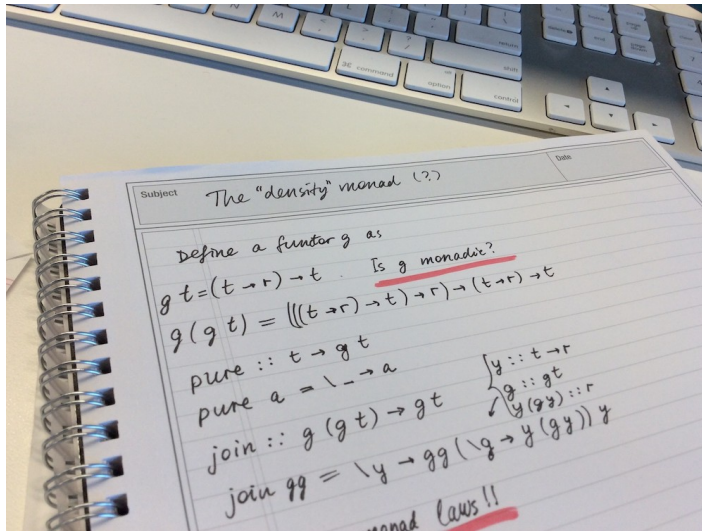


Figure 17.1: Example calculation in functional type theory.

is devoted not to showing code but to a peculiar kind of mathematical reasoning. Rather than focusing on one or another API or algorithm, as it is often the case with other software engineering blogs or presentations, an FP speaker describes a *mathematical structure*—such as the “**applicative functor**” or the “**free monad**”—and illustrates its use for practical coding.

These people are not some graduate students showing off their theoretical research; they are practitioners, software engineers who use FP on their jobs. It is just the nature of FP that certain mathematical tools—coming from formal logic and category theory—are now directly applicable to practical programming tasks.

These mathematical tools are not mere tricks for a specific programming language; they apply equally to all FP languages. Before starting to write code, the programmer can jot down certain calculations in a mathematical notation (see Fig. 17.1). The results of those calculations will help design the code fragment the programmer is about to write. This activity is quite similar to that of an engineer who first performs some mathematical calculations and only then embarks on a real-life design project.

A recent example of the hand-in-hand development of the functional type theory and its applications is seen in the “free applicative functor” construction. It was first described in a 2014 paper; a couple of years later, a combined free applicative + free monad data type was designed and its implementation proposed in Scala as well as in Haskell. This technique allows programmers to work with declarative side-effect computations where some parts are sequential but other parts can be computed in parallel, and to achieve the parallelism *automatically* while maintaining the composability of the resulting programs. The new technique has distinct advantages over using monad transformers, which was the previous method of composing declarative side-effects.

The “free applicative + free monad” combination was designed and implemented by true software engineers. They first wrote down the types and derived the necessary algebraic properties; the obtained results directly guided them about how to proceed writing the library API.

Another example of a development in functional type theory is the so called “tagless final” encoding of data types, first described in 2009. This technique, developed from category theory and type theory motivations, has several advantages over the free monad technique and can improve upon it in a number of cases—just as the free monad itself was designed to cure certain problems with monad transformers. The new technique is also not a trick in a specific programming language; rather, it is a theoretical development that is available to programmers in any language (even in Java).

This example shows that we may need several more years of work before the practical aspects of using “functional type theory” are sufficiently well understood by the FP community. The theory is in active development, and its design patterns—as well as the exact scope of the requisite theoretical material—are still being figured out. If the 40-year gap hypothesis holds, we should expect functional type theory (perhaps under a different name) to become mainstream by 2030. This book is a step towards a clear designation of the scope of that theory.

17.5 Do we need software *engineering*, or are artisans good enough?

The demand for programmers is growing. “Software developer” was #1 best job in the US in 2018. But is there a demand for engineers, or just for artisans?

We do not seem to be able to train enough software artisans. Therefore, it

17.5 Do we need software engineering, or are artisans good enough?

is probably impossible to train as many software engineers in the true sense of the word. Modern courses in Computer Science do not actually train engineers in that sense; at best, they train academics who act as software artisans when writing code. The few existing true software *engineers* are all self-taught. Recalling the situation in construction business, with a few architects and hundreds of construction workers, we might also conclude that, perhaps, only a few software engineers are required per hundred software artisans.

What is the price of *not* having engineers, of replacing them with artisans?

Software practitioners have long bemoaned the mysterious difficulty of software development. Code “becomes rotten with time”, programs grow in size “out of control”, and operating systems have been notorious for ever-appearing **security flaws** despite many thousands of programmers and testers employed. I think this shows we are overestimating the artisanal creative capacity of the human brain.

It is precisely in designing very large and robust software systems that we would benefit from true engineering. Consider that humanity has been using chemical reactions and building bridges by trial, error, and adherence to tradition, long before mechanical or chemical engineering disciplines were developed and founded upon rigorous theory. Once the theory became available, humanity proceeded to create unimaginably more complicated and powerful structures and devices than ever before.

For building large and reliable software, such as new mobile or embedded operating systems or distributed peer-to-peer trust architectures, we will most likely need the qualitative increase in productivity and reliability that can only come from transforming artisanal programming into a proper engineering discipline. Functional type theory and functional programming are first steps in that direction.

18 Essay: Towards functional data engineering with Scala

Published on 2017-09-29

Data engineering is among **the most in-demand** novel occupations in the IT world today. Data engineers create software pipelines that process large volumes of data efficiently. Why did the Scala programming language **emerge as a premier tool** for crafting the foundational data engineering technologies such as Spark or Akka? Why is **Scala in such demand** within the world of big data software?

There are reasons to believe that the choice of Scala was quite far from pure happenstance.

18.1 Data is math

Humanity has been working with data at least since **Babylonian tax tables** and the **ancient Chinese number books**. Mathematics summarizes several millennia's worth of data processing experience into a few tenets:

- Data is *immutable*, because facts are immutable.
- Each *type* of values—population count, land area, distances, prices, dates, times,—needs to be handled separately; e.g. it is meaningless to add a distance to a population count.
- Data processing is to be codified by *mathematical formulas*.

Violating these tenets produces nonsense (see Fig. 18.1).

The power of the basic principles of mathematics extends over all epochs and all cultures; they are the same in Rio de Janeiro, in Kuala-Lumpur, and even in Pyongyang (see Fig. 18.2).



Figure 18.1: A nonsensical calculation arises when mixing incompatible data types.

18.2 Functional programming is math

The functional programming paradigm is based on similar principles: values are immutable, data processing is coded through formula-like expressions, and each type of data is required to match correctly during the computations. A flexible system of data types helps programmers automatically prevent many kind of coding errors. In addition, modern programming languages such as Scala and Haskell have a set of features adapted to building powerful abstractions and domain-specific languages. This power of abstraction is not accidental. Since mathematics is the ultimate art of building abstractions, math-based functional programming languages capitalize on the advantage of several millennia of mathematical experience.

A prominent example of how mathematics informs the design of programming languages is the connection between **constructive logic** and the programming language's type system, called the **Curry-Howard (CH) correspondence**. The main idea of the CH correspondence is to think of programs as mathematical formulas that compute a value of a certain type A . The CH correspondence is between programs and logical propositions. To any program that computes a

value of type A , there corresponds a proposition stating that “a value of type A can be computed”.

This may sound rather theoretical so far. To see the real value of the CH correspondence, recall that formal logic has operations “*and*”, “*or*”, and “*implies*”. For any two propositions A, B , we can construct the propositions “ A *and* B ”, “ A *or* B ”, “ A *implies* B ”. These three logical operations are foundational; without one of them, the logic is *incomplete* (you cannot derive some theorems).

A programming language **obeys the CH correspondence** to the logic if for any two types A, B , the language also contains composite types corresponding to the logical formulas “ A *or* B ”, “ A *and* B ”, “ A *implies* B ”. In Scala, these composite types are `Either[A,B]`, the tuple `(A,B)`, and the function type, $A \Rightarrow B$. All modern functional languages such as OCaml, Haskell, Scala, F#, Swift, Rust, Elm, PureScript support these three type constructions and thus are faithful to the CH correspondence. Having a *complete* logic in a language’s type system enables **declarative domain-driven code design**.

It is interesting to note that most older programming languages (C/C++, Java, JavaScript, Python) do not support some of these composite types. In other words, these programming languages have type systems based on an incomplete logic. As a result, users of these languages have to implement burdensome workarounds that make for error-prone code. Failure to follow mathematical principles has real costs.

18.3 The power of abstraction

Data engineering at scale poses problems of such complexity that many software companies adopt functional programming languages as their main implementation tool. Netflix, LinkedIn, Twitter started using Scala early on and were able to reap the benefits of the powerful abstractions Scala affords, such as asynchronous streams and parallelized functional collections. In this way, Scala enabled these businesses to engineer and scale up their massively concurrent computations. What exactly makes Scala suitable for big data processing?

The only way to manage massively concurrent code is to use sufficiently high-level abstractions that make application code declarative. The two most important such abstractions are the “resilient distributed dataset” (RDD) of Apache Spark and the “reactive stream” used in systems such as Kafka, Apache Storm, Akka Streams, and Apache Flink. While these abstractions are certainly implementable in Java or Python, true declarative and type-safe usage is possible only



Figure 18.2: The Pyongyang method of error-free programming.

in a programming language with a sufficiently sophisticated functional type system. Among the currently available mature functional languages, only Scala and Haskell would be technically adequate for that task, due to their support for type classes and higher-order generic collections.

It remains to see why Scala became the *lingua franca* of big data and not, say, Haskell.

18.4 Scala is Java on math

The recently invented general-purpose functional programming languages can be grouped into “industrial” (F#, Scala, Swift) and “academic” (OCaml, Haskell).

The “academic” languages are clean-room implementations of well-researched mathematical principles of programming language design (the CH correspondence being one such principle). These languages are unencumbered by requirements of compatibility with any existing platform or libraries. Because of this, the “academic” languages are perfect playgrounds for taking various mathematical ideas to their logical conclusion. At the same time, software practitioners strug-

gle to adopt these languages due to a steep learning curve, a lack of enterprise-grade libraries and tool support, and immature package management.

The languages from the “industrial” group are based on existing and mature software ecosystems: F# on .NET, Scala on JVM, and Swift on Apple’s Cocoa/iOS platform. One of the important design requirements for these languages is 100% binary compatibility with their “parent” platforms and languages (F# with C#, Scala with Java, and Swift with Objective-C). Because of this, developers can immediately take advantage of the existing tooling, package management, and industry-strength libraries, while slowly ramping up the idiomatic usage of new language features. However, the same compatibility requirement necessitated certain limitations in the languages, making their design less than fully satisfactory from the functional programming viewpoint.

It is now easy to see why the adoption rate of the “industrial” group of languages is **much higher** than that of the “academic” languages. The transition to the functional paradigm is also made smoother for software developers because F#, Scala, and Swift seamlessly support the familiar object-oriented programming paradigm. At the same time, these new languages still have logically complete type systems, which gives developers an important benefit of type-safe domain modeling.

Nevertheless, the type systems of these languages are not equally powerful. For instance, F# and Swift are similar to OCaml in many ways but omit OCaml’s parameterized modules and some other features. Of all mentioned languages, Scala and Haskell are the only ones that directly support type classes and higher-order generics, which are necessary for expressing abstractions such as an automatically parallelized data set or an asynchronous data stream.

To see the impact of these advanced features of Scala and Haskell, consider LINQ, a domain-specific language for database queries on .NET, implemented in C# and F# through a special syntax supported by Microsoft’s compilers. Analogous functionality is provided in Scala as a *library*, without need to modify the Scala compiler, by several open-source projects such as Slick, Squeryl, or Quill. Similar libraries exist for Haskell—but are impossible to implement in languages with less powerful type systems.

18.5 Conclusion

The decisive advantages of Scala over other contenders (such as OCaml, Haskell, F#, Swift, or Rust) are

18 Essay: Towards functional data engineering with Scala

1. functional collections in the standard library;
2. a highly sophisticated type system, with support for type classes and higher-order generics;
3. seamless compatibility with a mature software ecosystem (JVM).

Based on this assessment, we may be confident in Scala's great future as a main implementation language for big data engineering.

A Notations

I chose certain notations in this book to be different from the notations currently used in the functional programming community. The proposed notation seems to be well adapted to reasoning about types and code, and especially for designing data types and proving the equational laws of typeclasses.

A.1 Summary table

F^A type constructor F with type argument A

$x:A$ value x has type A

$A + B$ a disjunctive type; in Scala, this type is `Either[A, B]`

$A \times B$ a product type; in Scala, this type is `(A,B)`

$A \Rightarrow B$ the function type, mapping from A to B

$x:A \Rightarrow f$ a nameless function (as a value); in Scala, `{x:A => f}`

id the identity function

1 the unit type or its value; in Scala, `Unit` and `()`

0 the void type; in Scala, `Nothing`

\triangleq “equal by definition”

\cong “equivalent” according to an established isomorphism of types

A^{F^B} type annotation, used for defining unfunctors (GADTs)

fmap_F the standard method `fmap` pertaining to a functor F

pu_F the standard method `pure` of a monad F

A Notations

F^\bullet the type constructor F understood as a type-level function; in Scala, `F[_]`

$F^\bullet \rightsquigarrow G^\bullet$ or $F^A \rightsquigarrow G^A$ a natural transformation between functors F and G

$\forall A. P^A$ a universally quantified type expression

$\exists A. P^A$ an existentially quantified type expression

\circ the forward composition of functions: $f \circ g$ is $x \Rightarrow g(f(x))$.

\circ the backward composition of functions: $f \circ g$ is $x \Rightarrow f(g(x))$.

\circ functor composition: $F \circ G$ is F^{G^\bullet}

$f^{\uparrow G}$ a function f lifted to a functor G ; same as `fmapG f`

$f^{\uparrow G \uparrow H}$ a function lifted first to G and then to H ; In Scala, `h.map(_ . map(f))`

$f^{\downarrow H}$ a function f lifted to a contrafunctor

\diamond_M the Kleisli product operation for the monad M

A.2 Explanations

F^A means a type constructor F with a type parameter A . In Scala, this is `F[A]`.

$x^{:A}$ means a value x that has type A . The colon, `:`, in the superscript is to show that it is not a type argument. In Scala, this is `x : A`. The notation $x : A$ can be used as well, but $x^{:A}$ is easier to read when x is inside a larger code expression.

$A + B$ means the disjunctive type made of types A and B (or, a disjunction of A and B). In Scala, this is the type `Either[A, B]`.

$A \times B$ means the product type made of types A and B . In Scala, this is the tuple type `(A,B)`.

$A \Rightarrow B$ means a function type from A to B . In Scala, this is the function type `A => B`.

$x^{:A} \Rightarrow f$ means a nameless function (a literal function value) with argument x of type A and function body f .

id means the identity function. The type of its argument should be either specified as id^A or $\text{id}^{:A \Rightarrow A}$, or else should be unambiguous from the context.

1 means the unit type, and also the value of the unit type. In Scala, the unit type is `Unit`, and its value is denoted by `()`. Example of using the unit type is $1 + A$, which in Scala corresponds to `Option[A]`.

0 means the void type (the type with no values). In Scala, this is the type `Nothing`. Example of using the void type is to denote the empty part of a disjunction. For example, in the disjunction $1 + A$ the non-empty option type is $0 + A$, which in Scala corresponds to `Some[A]`. The empty option type $1 + 0$ corresponds to `None`. Similarly, $A + 0$ may denote the left part of the type $A + B$ (in Scala, `Left[A]`), while $0 + B$ denotes its right part (in Scala, `Right[A]`).

\triangleq means “equal by definition”. Examples:

- $f \triangleq x:\text{Int} \Rightarrow x + 1$ is a definition of the function f . In Scala, this is `val f = (x: Int) => x + 1`.
- $F^A \triangleq 1 + A$ is a definition of a type constructor F . In Scala, this is `type F[A] = Option[A]`.

\cong means “equivalent” according to an equivalence relation that needs to be established in the text. For example, if we have established the equivalence that allows nested tuples to be reordered whenever needed, we can write $(a \times b) \times c \cong a \times (b \times c)$.

A^{F^B} in type expressions means that the type constructor F^\bullet assigns the type F^B to the type expression A . This notation is used for defining unfunctors (GADTs). For example, the Scala code

```
sealed trait F[A]
case class F1() extends F[Int]
case class F2[A](a: A) extends F[(A, String)]
```

defines an unfunctor, which is denoted by $F^A \triangleq 1^{F^{\text{Int}}} + A^{F^{A \times \text{String}}}$.

fmap_F means the standard method `fmap` pertaining to the functor F . In Scala, this is typically implemented as `Functor[F].fmap`. Since each functor F has its own specific definition of `fmap`, the subscript “ F ” is not a type parameter of fmap_F . The method fmap_F actually has *two* type parameters, and its type signature written in full is $\text{fmap}_F^{A,B} : (A \Rightarrow B) \Rightarrow F^A \Rightarrow F^B$. For clarity, one may sometimes write explicitly the type parameters A, B in the expression $\text{fmap}_F^{A,B}$, but in most cases these type parameters A, B can be omitted without loss of clarity. As another example, a monad’s standard method `pure` is denoted by pu_F , where the subscript refers to the monad F . This function has type signature $A \Rightarrow F^A$ that contains a type parameter A . In the short notation, this is denoted by pu_F^A . If we are using the `pure` method with a complicated type, e.g. $1 + P^A$, instead of the type parameter A , we might want to indicate this type for clarity as a type

A Notations

parameter value and write $\text{pu}_F^{1+P^A}$. The type signature of that function is then

$$\text{pu}_F^{1+P^A} : 1 + P^A \Rightarrow F^{1+P^A} \quad .$$

F^\bullet means the type constructor F understood as a type-level function, – that is, with a type argument unspecified. In Scala, this is `F[_]`. The bullet symbol, \bullet , is used as a placeholder for the missing type parameter. I also simply write F when no type argument is needed, and it means the same as F^\bullet . (For example, “a functor F ” and “a functor F^\bullet ” mean the same thing.) However, it is useful for clarity to be able to indicate the place where the type argument would appear. For instance, functor composition is clearly denoted as F^{G^\bullet} ; in Scala, this is `F[G[?]]` when using the “kind projector” plugin.¹ As another example, $T_L^{M,\bullet}$ denotes a monad transformer for the base monad L and the foreign monad M . The foreign monad M is a type parameter in $T_L^{M,\bullet}$, and so is the missing type parameter denoted by the placeholder symbol \bullet . (However, the base monad L is not a type parameter in $T_L^{M,\bullet}$ because the construction of the monad transformer depends sensitively on the internal details of L .)

$F^\bullet \rightsquigarrow G^\bullet$ or $F^A \rightsquigarrow G^A$ means a natural transformation between two functors F and G . In some Scala libraries, this is denoted by `F ~> G`.

$\forall A. P^A$ is a universally quantified type expression, in which A is a bound type parameter.

$\exists A. P^A$ is an existentially quantified type expression, in which A is a bound type parameter.

\circledast means the forward composition of functions: $f \circledast g$ (reads “ f before g ”) is the function defined as $x \Rightarrow g(f(x))$.

\circ means the backward composition of functions: $f \circ g$ (reads “ f after g ”) is the function defined as $x \Rightarrow f(g(x))$.

\circ with type constructors means their functor composition, for example $F \circ G$ denotes the functor F^{G^\bullet} . In Scala, this is `F[G[A]]`.

$f^{\uparrow G}$ means a function f lifted to a functor G . For a function $f:A \Rightarrow B$, the application of $f^{\uparrow G}$ to a value $g : G^A$ is written as $f^{\uparrow G} g$. In Scala, this is `g.map(f)`. Nested liftings can be written as $f^{\uparrow G \uparrow H}$, which means $(f^{\uparrow G})^{\uparrow H}$ and produces a function of type $H^{G^A} \Rightarrow H^{G^B}$. Applying a nested lifting to a value h of type H^{G^A} is written as $f^{\uparrow G \uparrow H} h$. In Scala, this is `h.map(_ . map(f))`.

$f^{\downarrow H}$ means a function f lifted to a contrafunctor H . For a function $f:A \Rightarrow B$, the application of $f^{\downarrow H}$ to a value $h : H^B$ is written as $f^{\downarrow H} h$ and yields a value of type

¹<https://github.com/typelevel/kind-projector>

H^A . In Scala, this is `h.contramap(f)`.

\diamond_M means the Kleisli product operation for the monad M . This is a binary operation working on two Kleisli functions of types $A \Rightarrow M^B$ and $B \Rightarrow M^C$ and yields a new function of type $A \Rightarrow M^C$.

B Glossary of terms

I chose certain terms in this book to be different from the terms currently used in the functional programming community. My proposed terminology is designed to help readers understand and remember the concepts behind the terms.

Nameless function An expression of function type, representing a function. For example, $x \mapsto x * 2$. Also known as function expression, function literal, anonymous function, closure, lambda-function, lambda-expression, or simply a “lambda”.

Product type A type representing several values given at once. In Scala, product types are the tuple types, for example `(Int, String)`, and case classes. Also known as **tuple type**, **struct** (in C and C++), and **record**.

Disjunctive type A type representing one of several distinct possibilities. In Scala, this is usually implemented as a sealed trait extended by several case classes. The standard Scala disjunction types are `Option[A]` and `Either[A, B]`. Also known as **sum type**, **tagged union type**, **co-product type**, and variant type (in OCaml and in Object Pascal). It is shorter to say “sum type,” but the English word “sum” is more ambiguous to the ear than “disjunction”.

Polynomial functor A type constructor built using disjunctions (sums), products (tuples), type parameters and fixed types. For example, in Scala, `type F[A] = Either[(Int, A), A]` is a polynomial functor with respect to the type parameter `A`, while `Int` is a fixed type (not a type parameter). Polynomial functors are also called **algebraic data types**. A polynomial type constructor is always a functor with respect to any of its type parameters, hence I use the shorter name “polynomial functor” instead of “polynomial type constructor”.

Unfunctor A type constructor that cannot possibly be a functor, nor a contrafunctor, nor a profunctor. An example is a type constructor with explic-

B Glossary of terms

itly indexed type parameters, such as $F^A \triangleq (A \times A)^{F^{\text{Int}}} + (\text{Int} \times A)^{F^1}$. Also called a **GADT** (generalized algebraic data type).

Functor block A short syntax for composing various `map`, `flatMap`, and `filter` operations on a functor value. The Functor typeclass instance must be fixed for the entire functor block. For example, in Scala the code

```
for { x <- List(1,2,3); y <- List(10, x); if y > 2 }  
  yield 2 * y
```

is equivalent to the code

```
List(1, 2, 3).flatMap(x => List(10, x))  
  .filter(y => y > 1).map(y => 2 * y)
```

and computes the value `List(20, 20, 20, 6)`. Similar syntax exists in a number of languages and is called a **for-comprehension** or **list comprehension** in Python, **do-notation** or **do-block** in Haskell, and **computation expressions** in F#. I use the name “functor block” in this book because it is shorter and more descriptive.

Method 1) A function defined as a member of a typeclass. For example, `flatMap` is a method defined in the `Monad` typeclass. 2) In Scala, a function defined as a member of a data type declared as a Java-compatible `class` or `trait`. Trait methods are necessary in Scala when implementing functions whose arguments have type parameters (because ordinary Scala function values cannot have type parameters). So, many type classes such as `Functor` or `Monad`, whose methods₁ require type parameters, will use Scala traits with methods₂ for their implementation. The same applies to type constructions with quantified types, such as the Church encoding.

Kleisli function Also called a Kleisli morphism or a Kleisli arrow. A function with type signature $A \Rightarrow M^B$ for some fixed monad M . More verbosely, “a morphism from the Kleisli category corresponding to the monad M ”. The standard monadic method $\text{pure}_M : A \Rightarrow M^A$ has the type signature of a Kleisli function. The Kleisli product operation, \diamond_M , is a binary operation that combines two Kleisli functions (of types $A \Rightarrow M^B$ and $B \Rightarrow M^C$) into a new Kleisli function (of type $A \Rightarrow M^C$).

B.1 On the current misuse of the term “algebra”

Exponential-polynomial type A type constructor built using disjunctions (sums), products, and function types, as well as type parameters or fixed types. For brevity, I call them “exp-poly” types. For example, in Scala, `type F[A] = Either[(A, A), Int > A]` is an exp-poly type constructor. Such type constructors can be functors, contrafunctors, or profunctors.

Short type notation A mathematical notation for type expressions developed in this book for the purpose of quicker and more practical reasoning about types in functional programs. Disjunction types are denoted by $+$, product types by \times , and function types by \Rightarrow . The unit type is denoted by 1 , and the void type by 0 . The function arrow \Rightarrow has weaker precedence than $+$, which is in turn weaker than \times . Type parameters are denoted by superscripts. As an example of using these conventions, the Scala definition

```
type F[A] = Either[(A, A > Option[Int]), String > List[A]]
```

is written in the short type notation as

$$F^A \triangleq A \times (A \Rightarrow 1 + \text{Int}) + (\text{String} \Rightarrow \text{List}^A) \quad .$$

B.1 On the current misuse of the term “algebra”

In this book, I do not use the terms “algebra” or “algebraic”, because these terms are too ambiguous. In the current practice, the functional programming community is using the word “algebra” in at least *four* incompatible ways.

Definition 0. In mathematics, an “algebra” is a vector space with multiplication and certain standard properties. For example, you need $1 * x = x$, the addition must be commutative, the multiplication must be distributive over addition, and so on. As an example, the set of all 10×10 matrices with real coefficients is an “algebra” in this sense. These matrices form a 100-dimensional vector space, and can be multiplied and added. This standard definition of “algebra” is not actually used in functional programming.

Definition 1. An “algebra” is a function with type signature $F^A \Rightarrow A$, where F^A is some fixed functor. This definition comes from category theory, where such types are called *F-algebras*. There is no direct connection between this “algebra”

and Definition 0, except when the functor F is defined by $F^A \triangleq A \times A$, and then a function of type $A \times A \Rightarrow A$ may be interpreted as a “multiplication” operation (but, in any case, A is a type and not a vector space, and there are no distributivity or commutativity laws). I prefer to call such functions “ F -algebras”, emphasizing that they characterize and depend on a chosen functor F . However, F -algebras are not mentioned in this book: knowing how to reason about their properties does not give much help in practical work.

Definition 2. Polynomial functors are often called “algebraic data types”. However, they are not “algebraic” in the sense of Definition 0 or 1. For example, consider the “algebraic data type” `Either[Option[A], Int]`, which is $F^A \triangleq 1 + A + \text{Int}$ in the short type notation. The set of all values of the type F^A does not admit the addition and multiplication operations required by the mathematical definition of “algebra”. The type F^A may admit some binary or unary operations (e.g. that of a monoid), but these operations will not be commutative or distributive. Also, there cannot be a function with type $F^A \Rightarrow A$, as required for Definition 1. It seems that the usage of the word “algebra” here is to refer to “school-level algebra” with polynomials; these data types are built from sums and products of types. In this book, I call such types “polynomial”. However, if the data type contains a function type, e.g. `Option[Int => A]`, the type is no longer polynomial. So I use the more precise terms “polynomial type” and “exponential-polynomial type”.

Definition 3. People talk about the “algebra” of properties of functions such as `map` or `flatMap`, referring to the fact that these functions must satisfy certain equational laws (e.g. the identity, composition, or associativity laws). But these laws do not form an “algebra” in the sense of Definition 0, nor do the functions such as `map` or `flatMap` themselves (there are no binary operations on them). Neither do they form an algebra in the sense of Definition 1. The laws for `map` or `flatMap` are in no way related to “algebraic data types” of Definition 2. So here the word “algebra” is used in a way that is unrelated to the three previous definitions. To me, it does not seem helpful to say the word “algebra” or “algebraic” when talking about equational laws. These laws are “algebraic” in a trivial sense – i.e. they are written as equations. In mathematics, “algebraic” equations are different from “differential” or “integral” equations. In functional programming, all equational laws are of the same kind: some code on the left-hand side must be equal to some code on the right-hand side of the equation. So calling them

B.1 On the current misuse of the term “algebra”

“algebraic” does not help and does not clarify anything. I call them “equational laws” or just “laws”.

Definition 4. In the Church encoding of a free monad (nowadays known as the “final tagless” encoding), the term “algebra” refers to the *type constructor* parameter F . This definition has nothing to do with any of the previous definitions. Clearly, Definition 0 cannot apply to a type constructor. Definition 1 does not apply since F is not itself a function type of the form $G^A \Rightarrow A$. (A function of type $F^A \Rightarrow A$, which I call a “runner” for the type constructor F , is not usually called an “algebra” in discussions of the “final tagless” encoding.) Definition 2 seems to be the most closely related meaning, since F is *sometimes* a polynomial functor in practical usage (although in most cases F will be an unfunctor). However, it is not helpful to call the polynomial functor F an “algebraic data type” and, at the same time, an “algebra”. Definition 3 does not apply since the free monad construction does not assume that any laws hold about F , nor has any means of imposing such laws. The type constructor F is used to parameterize the effects described by the free monad, so it seems more reasonable to call it the “effect constructor”.

So, it seems that the current usage of the word “algebra” in functional programming is both inconsistent and unhelpful to practitioners. In this book, I reserve the word “algebra” to denote the branch of mathematics, as in “school-level algebra” or “graduate-level algebra”. Instead of “algebra” as in Definitions 1 to 4, I talk about “ F -algebras” with a specific functor F ; “polynomial types” or “polynomial functors” or “exponential-polynomial functors” etc.; “equational laws”; and an “effect constructor” F .

C Scala syntax and features

C.0.1 Function syntax

Functions have arguments, body, and type. The function type lists the type of all arguments and the type of the result value.

```
def f(x: Int, y: Int): Int ⇒ Int = { z ⇒ x + y + z }
```

Functions may be used with infix syntax as well. For this syntax to work, the function must be defined **as a Scala method**, that is, using `def` within the declaration of `x`'s class, or as an extension method. The infix syntax cannot work with functions defined using `val`. For clarity, I call Scala functions **infix methods** when defined and used in this way.

The syntax `List[Int]` means “a list of integer values.” In the type expression `List[Int]`, the “`Int`” is called the **type parameter** and `List` is called the **type constructor**.

A list can contain values of any type; for example, `List[List[List[Int]]]` means a list of lists of lists of integers. So, a type constructor can be seen as a function from types to types. A type constructor takes a type parameter as an argument, and produces a new type as a result.

C.0.2 Scala collections

The Scala standard library defines collections of several kinds, the main ones being sequences, sets, and dictionaries. These collections have many map/reduce-style methods defined on them.

Sequences are “subclasses” of the class `Seq`. The standard library will sometimes choose automatically a suitable subclass of `Seq`, such as `List`, `IndexedSeq`, `Vector`, `Range`, etc.; for example:

```
scala> 1 to 5
scala> (1 to 5).map(x ⇒ x*x)
scala> (1 to 5).toList
```

C Scala syntax and features

```
scala> 1 until 5  
scala> (1 until 5).toList
```

For our purposes, all these “sequence-like” types are equivalent.

Sets are values of class `Set`, and dictionaries are values of class `Map`.

```
scala> Set(1, 2, 3).filter(x ⇒ x % 2 == 0)
```

D Intuitionistic (constructive) propositional logic

The intuitionistic propositional logic describes how programs in functional programming languages may be able to compute values of different types.

The main difference between IPL and the classical (Boolean) logic is that IPL does not include the axiom of excluded middle (“*tertium non datur*”), which is

$$\forall A : (A \text{ or } (\text{not}A)) \text{ is true} \quad .$$

The reason this axiom is not included in IPL is that IPL propositions such as $\mathcal{CH}(A)$ correspond to the *possibility* of values of type A *to be computed* by a program. For the proposition $\mathcal{CH}(A)$ to be true in IPL, a program needs to actually compute a value of type A . It is not sufficient merely to show that the non-existence of such a value would be somehow contradictory. (In classical logic, that would be sufficient.)

Another significant difference between IPL and the Boolean logic is that propositions in IPL cannot be assigned a fixed set of “truth values”. This was proved by Gödel in 1935. It means that a proposition in IPL cannot be decided by writing out a truth table, even if we allow many truth values.

E Category theory

Examples of categories

1. Objects: types Int , String , ...; morphisms (arrows) are functions $\text{Int} \rightarrow \text{String}$ etc. – this is the “standard” category corresponding to a given programming language
2. Objects: types A , B , ...; morphisms are pairs of functions $(A \rightarrow B), (B \rightarrow A)$
3. * Objects: types List^A , List^B , ...; morphisms are functions of type $\text{List}^A \rightarrow \text{List}^B$
4. Objects: types A , B , ...; morphisms are functions of type $\text{List}^A \rightarrow \text{List}^B$
5. Objects: types A , B , ...; morphisms are functions of type $A \rightarrow \text{List}^B$
6. * Objects: types List^A , List^B , ...; morphisms are functions $A \rightarrow B$
7. Objects: types A , B , ...; morphisms are $\text{List}^{A \rightarrow B}$
8. Objects: types A , B , ...; morphisms are functions $B \rightarrow A$
9. * Objects: things A , B , ...; morphisms are pairs (A, B) of things – this is the same as a preorder

Examples marked with * are for illustration only, they are probably not very useful

F What is “applied functional type theory”

F.1 The import and scope of AFTT

Applied functional type theory (AFTT) is what I call the area of computer science that should serve the needs of functional programmers who are working as software engineers.

It is for these practitioners (I am one myself), rather than for academic researchers, that I set out to examine the incredible wealth of functional programming inventions over the last 30 years, – such as the “functional pearls” papers – and to determine the scope of theoretical material that has demonstrated its pragmatic usefulness and thus belongs to AFTT, as opposed to material that is purely academic and may be tentatively omitted.

What exactly is the extent of “theory” that a practicing functional programmer should know in order to be effective at writing code? In my view, this question is not yet resolved. Once it is resolved, AFTT will be that theory.

Traditional courses of theoretical computer science (algorithms, traditional data structures, complexity theory, formal languages, formal semantics, compiler techniques, databases, networking, operating systems, etc.) are largely not relevant to AFTT.

Here is an example: To an academic computer scientist, the “science behind Haskell” is the theory of lambda-calculus, the type-theoretic “System $F\omega$ ”, and formal semantics. These theories guided the design of the Haskell language and define rigorously what a Haskell program “means” in a mathematical sense. Academic computer science courses teach these theories.

However, a practicing Haskell or Scala programmer is not concerned with designing Haskell or Scala, or with proving any theoretical properties of those languages. A practicing programmer is mainly concerned with *using* a chosen programming language to *write code*.

Neither the theory of lambda-calculus, nor proofs of type-theoretical proper-

F What is “applied functional type theory”

ties of “System $F\omega$ ”, nor theories of formal semantics will actually help a programmer to write code. So all these theories are not within the scope of AFTT.

As an example of theoretical material that *is* within the scope of AFTT, consider the equational laws imposed on applicative functors (see Chapter 8).

It is essential for a practicing functional programmer to be able to recognize and use applicative functors. An applicative functor is a data structure can specify declaratively a set of operations that do not depend on each other’s effects. Programs can then easily combine these operations and effects, for example, in order to execute them in parallel, or to refactor the program for better maintainability.

To use this functionality, the programmer must begin by checking whether a given data structure that satisfies the laws of applicative functors. In a given application, a data structure may be dictated in part by the business logic that cannot be changed. The programmer first writes down the type of that data structure and the code implementing the required methods, and then checks that the laws hold. The data structure may need to be adjusted in order to fit the definition of an applicative functor or its laws.

This work is done using pen and paper, in a mathematical notation. Once the applicative laws are verified, the programmer proceeds to write code using that data structure.

Because of the proofs, it is assured that the data structure satisfies the known properties of applicative functors, no matter how the rest of the program is written. So, for example, it is assured that the relevant effects can be automatically parallelized, as is usual with applicative functors.

In this way, AFTT directly guides the programmer and helps to write correct code.

Applicative functors were discovered by practitioners who were using Haskell for writing code, in applications such as parser combinators, compilers, and domain-specific languages for parallel computations. However, applicative functors are not a feature of Haskell: they are the same in Scala, OCaml, or any other functional programming language. And yet, no standard computer science textbook defines applicative functors, motivates their laws, explores their structure on basic examples, or shows data structures that are *not* applicative functors and explains why. (For that matter, no book on category theory or type theory mentions applicative functors either.)

So far it appears that AFTT will be a mixture of certain areas of category theory, formal logic, and type theory. However, software engineers would not derive much benefit from following traditional academic courses in these subjects,

because their choice of material is too abstract and lacks specific results necessary for practical software engineering. In other words, the traditional academic courses answer questions that academic computer scientists have, not questions that software engineers have.

There exist several books intended as presentations of category theory “for computer scientists” or “for programmers”. However, these books fail to give examples vitally relevant to everyday programming, such as applicative or traversable functors. Instead, these books contain purely theoretical topics such as limits, adjunctions, or toposes, – concepts that have no applications in practical functional programming today.

Typical questions in academic books: “is this an introduction or an elimination rule” and “does this property hold in non-small categories, or are we limited to the category *Set*”. Typical questions a Scala programmer might have: “can we compute a value of type `Either[Z, R \Rightarrow A]` from a value of type `R \Rightarrow Either[Z, A]`” and “is the type constructor `F[A] = Option[(A, A, A)]` a monad”. The proper scope of AFTT includes answering the last two questions, but not the first two.

A software engineer hoping to understand the foundations of functional programming will not find the concepts of foldable, filterable, applicative, or traversable functors in any books on category theory, including books intended for programmers. And yet, these concepts are necessary to obtain a mathematically correct implementation of such foundationally important operations as fold, filter, zip, and traverse – operations that functional programmers use every day in their code.

To compensate for the lack of AFTT textbooks, programmers have written many online tutorials for each other, trying to explain the theoretical concepts necessary for practical work. There are the infamous “monad tutorials”, but also tutorials about applicative functors, traversable functors, free monads, and so on. These tutorials tend to be very hands-on and narrow in scope, limited to one or two specific questions and specific applications. Such tutorials usually do not present a sufficiently broad picture and do not illustrate deeper connections between these mathematical constructions.

For example, “free monads” became popular in the Scala community around 2015. Many talks about free monads were presented at Scala engineering conferences, each giving their own slightly different implementation but never formulating rigorously the required properties for a piece of code to be a valid implementation of the free monad.

Without knowledge of mathematical principles behind free monads, a pro-

F What is “applied functional type theory”

grammer cannot make sure that a given implementation is correct. However, books on category theory present free monads in a way that is unsuitable for programming applications: a free monad is just an adjoint functor to a forgetful functor into the category of sets.¹ This definition is too abstract and cannot be used to verify whether a given implementation of the free monad in Scala is correct.

Perhaps the best selection of AFTT tutorial material can be found in the [Haskell Wikibooks](#). However, those tutorials are incomplete and limited to explaining the use of Haskell. Many of them are suitable neither as a first introduction nor as a reference on AFTT.

Apart from referring to some notions from category theory, AFTT also uses some concepts from type theory and formal logic. However, existing textbooks on type theory and formal logic focus on domain theory and proof theory – which is a lot of information that practicing programmers will have difficulty assimilating and yet will have no chance of ever applying in their daily work. At the same time, these books never mention practical techniques used in many functional programming libraries today, such as quantified types, types parameterized by type constructors, or partial type-level functions (known as “typeclasses”).

Type theory and formal logic can, in principle, help the programmer with certain practical tasks, such as:

- deciding whether two data structures are equivalent as types, and implementing the isomorphism transformation; for example, the Scala type `(A, Either[B, C])` is equivalent to `Either[(A, B), (A, C)]`
- detecting whether a definition of a recursive type is “valid”, i.e. does not lead to a useless infinite recursion; an example of a “useless” recursive type definition in Scala is `case class Bad(x: Bad)`
- deriving an implementation of a function from its type signature and required laws; for example, deriving the `flatMap` function for the `Reader` monad from the type signature `def flatMap[Z, A, B](r: Z => A)(f: A => Z => B): Z => B` and verifying that the monad laws hold
- deciding whether a generic pure function with a given signature can be implemented; for example, `def f[A, B]: (A => B) => A` cannot be implemented but `def g[A, B]: A => (B => A)` can be implemented

¹“What’s your problem?” as the joke goes.

I mention these practical tasks as examples because they are perhaps the only real-world-coding applications of domain theory and the Curry-Howard correspondence theory, besides programming language design. However, existing books on type theory and logic do not give practical recipes for resolving these questions.

On the other hand, books such as “[Scala with Cats](#)” and “[Functional programming, simplified](#)” are focused on explaining the practical aspects of programming and do not adequately treat the algebraic laws that the mathematical structures require (such as the laws for applicative or monadic functors).

The only existing Scala-based AFTT textbook aiming at the proper scope is the [Bjarnason-Chiusano book](#), which balances practical considerations with theoretical developments such as algebraic laws. This book is written at about the same level but goes deeper into the mathematical foundations and at the same time gives a wider range of examples.

This book is an attempt to delineate the proper scope of AFTT and to develop a rigorous yet clear and approachable presentation of the chosen material. The goal is to teach programmers how to reason mathematically about types and code, in a way that is directly relevant to practical programming.

In this book, I demonstrate code examples in Scala because this is what I am most familiar with. However, most of this material will work equally well in Haskell, OCaml, and other FP languages. This is so because AFTT is not specific to Scala or Haskell. A serious user of any other functional programming language will have to face the same questions and struggle with the same practical issues.

F.2 Approach to presentation

The presentation is self-contained. I introduce and explain all the required notations, concepts, and Scala language features. The emphasis is on clarity and understandability of all examples, mathematical notions, derivations, and code. I introduce a fair amount of non-standard terminology and notations because this allows me to achieve a clearer and more logically coherent presentation of the material, especially for readers not familiar with the theoretical literature.

The main content of AFTT is to expose mathematical principles that guide the practice of functional programming – that is, help people to write code. Therefore, all mathematical developments in this book are thoroughly motivated by practical programming issues.

F What is “applied functional type theory”

Each mathematical construction is accompanied by code examples and unit tests to illustrate its usage and check correctness. For example, the equational laws for every standard typeclass (functor, applicative, monad, etc.) are first motivated heuristically before deriving a set of mathematical equations and formulating the laws in terms of axioms of a category. In this book, all theory is developed in order to write some new code as a result, and in order to answer a question arising from programming practice. Theoretical material from category theory, type theory, and formal logic is kept to the absolute minimum necessary. Mathematical generalizations are not pursued beyond either practical usefulness or immediate pedagogical usefulness. This limits the scope of required mathematical knowledge to bare rudiments of category theory, type theory, and formal logic. For instance, I never mention “introduction/elimination rules”, “strong normalization”, “complete partial order domain”, “adjoint functor”, “pullback”, “pushout”, “limit”, “co-limit”, or the word “algebra”, because using these concepts will not help a functional programmer to write code. Instead, I focus on practically useful material – including some little-known but useful constructions such as the “partial type function”, “filterable” and “applicative contrafunctor” typeclasses, and “free pointed” monad.

Each new concept or technique is fully explained via “worked examples” and drilled via provided exercises. Answers to exercises are not provided, but it is verified that the exercises are doable and free of errors. More difficult exercises are marked by an asterisk.

F.3 Intended audience

The material is presented here at medium to advanced level. It requires a certain amount of mathematical experience and is not suitable for people unfamiliar with school-level algebra, or for people who are generally allergic to mathematics, or for people unwilling to learn new and difficult concepts through prolonged mental concentration.

The first two chapters are introductory and may be suitable for beginners in programming. Starting from the middle of chapter 3, the material becomes unsuitable for beginners.

G GNU Free Documentation License

Version 1.2, November 2002

Copyright (c) 2000,2001,2002 Free Software Foundation, Inc. 59 Temple Place, Suite 330, Boston, MA 02111-1307, USA
Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

Preamble

The purpose of this License is to make a manual, textbook, or other functional and useful document free in the sense of freedom: to assure everyone the effective freedom to copy and redistribute it, with or without modifying it, either commercially or noncommercially. Secondly, this License preserves for the author and publisher a way to get credit for their work, while not being considered responsible for modifications made by others.

This License is a kind of “copyleft”, which means that derivative works of the document must themselves be free in the same sense. It complements the GNU General Public License, which is a copyleft license designed for free software.

We have designed this License in order to use it for manuals for free software, because free software needs free documentation: a free program should come with manuals providing the same freedoms that the software does. But this License is not limited to software manuals; it can be used for any textual work, regardless of subject matter or whether it is published as a printed book. We recommend this License principally for works whose purpose is instruction or reference.

G.0.0 Applicability and definitions

This License applies to any manual or other work, in any medium, that contains a notice placed by the copyright holder saying it can be distributed under the terms of this License. Such a notice grants a world-wide, royalty-free license, unlimited in duration, to use that work under the conditions stated herein. The “Document”, below, refers to any such manual or work. Any member of the public is a licensee, and is addressed as “you”. You accept the license if you copy, modify or distribute the work in a way requiring permission under copyright law.

A “Modified Version” of the Document means any work containing the Document or a portion of it, either copied verbatim, or with modifications and/or translated into another language.

A “Secondary Section” is a named appendix or a front-matter section of the Document that deals exclusively with the relationship of the publishers or authors of the Document to the Document’s overall subject (or to related matters) and contains nothing that could fall directly within that overall subject. (Thus, if the Document is in part a textbook of mathematics, a Secondary Section may not explain any mathematics.) The relationship could be a matter of historical connection with the subject or with related matters, or of legal, commercial, philosophical, ethical or political position regarding them.

The “Invariant Sections” are certain Secondary Sections whose titles are designated, as being those of Invariant Sections, in the notice that says that the Document is released under this License. If a section does not fit the above definition of Secondary then it is not allowed to be designated as Invariant. The Document may contain zero Invariant Sections. If the Document does not identify any Invariant Sections then there are none.

The “Cover Texts” are certain short passages of text that are listed, as Front-Cover Texts or Back-Cover Texts, in the notice that says that the Document is released under this License. A Front-Cover Text may be at most 5 words, and a Back-Cover Text may be at most 25 words.

A “Transparent” copy of the Document means a machine-readable copy, represented in a format whose specification is available to the general public, that is suitable for revising the document straightforwardly with generic text editors or (for images composed of pixels) generic paint programs or (for drawings) some widely available drawing editor, and that is suitable for input to text formatters or for automatic translation to a variety of formats suitable for input to text formatters. A copy made in an otherwise Transparent file format whose markup, or absence of markup, has been arranged to thwart or discourage subsequent modification by readers is not Transparent. An image format is not Transparent if used for any substantial amount of text. A copy that is not “Transparent” is called “Opaque”.

Examples of suitable formats for Transparent copies include plain ASCII without markup, Texinfo input format, LaTeX input format, SGML or XML using a publicly available DTD, and standard-conforming simple HTML, PostScript or PDF designed for human modification. Examples of transparent image formats include PNG, XCF and JPG. Opaque formats include proprietary formats that can be read and edited only by proprietary word processors, SGML or XML for which the DTD and/or processing tools are not generally available, and the machine-generated HTML, PostScript or PDF produced by some word processors for output purposes only.

The “Title Page” means, for a printed book, the title page itself, plus such following pages as are needed to hold, legibly, the material this License requires to appear in the title page. For works in formats which do not have any title page as such, “Title Page” means the text near the most prominent appearance of the work’s title, preceding the beginning of the body of the text.

A section “Entitled XYZ” means a named subunit of the Document whose title either is precisely XYZ or contains XYZ in parentheses following text that translates XYZ in another language. (Here XYZ stands for a specific section name mentioned below, such as “Acknowledgements”, “Dedications”, “Endorsements”, or “History”.) To “Preserve the Title” of such a section when you modify the Document means that it remains a section “Entitled XYZ” according to this definition.

G GNU Free Documentation License

The Document may include Warranty Disclaimers next to the notice which states that this License applies to the Document. These Warranty Disclaimers are considered by reference in this License, but only as regards disclaiming warranties: any other implication that these Warranty Disclaimers may have is void and has no effect on the meaning of this License.

G.0.1 Verbatim copying

You may copy and distribute the Document in any medium, either commercially or noncommercially, provided that this License, the copyright notices, and the license notice saying this License applies to the Document are reproduced in all copies, and that you add no other conditions whatsoever to those of this License. You may not use technical measures to obstruct or control the reading or further copying of the copies you make or distribute. However, you may accept compensation in exchange for copies. If you distribute a large enough number of copies you must also follow the conditions in section [G.0.2](#).

You may also lend copies, under the same conditions stated above, and you may publicly display copies.

G.0.2 Copying in quantity

If you publish printed copies (or copies in media that commonly have printed covers) of the Document, numbering more than 100, and the Document's license notice requires Cover Texts, you must enclose the copies in covers that carry, clearly and legibly, all these Cover Texts: Front-Cover Texts on the front cover, and Back-Cover Texts on the back cover. Both covers must also clearly and legibly identify you as the publisher of these copies. The front cover must present the full title with all words of the title equally prominent and visible. You may add other material on the covers in addition. Copying with changes limited to the covers, as long as they preserve the title of the Document and satisfy these conditions, can be treated as verbatim copying in other respects.

If the required texts for either cover are too voluminous to fit legibly, you should put the first ones listed (as many as fit reasonably) on the actual cover, and continue the rest onto adjacent pages.

If you publish or distribute Opaque copies of the Document numbering more than 100, you must either include a machine-readable Transparent copy along with each Opaque copy, or state in or with each Opaque copy a computer-network location from which the general network-using public has access to download using public-standard network protocols a complete Transparent copy of the Document, free of added material. If you use the latter option, you must take reasonably prudent steps, when you begin distribution of Opaque copies in quantity, to ensure that this Transparent copy will remain thus accessible at the stated location until at least one year after the last time you distribute an Opaque copy (directly or through your agents or retailers) of that edition to the public.

If it is requested, but not required, that you contact the authors of the Document well before redistributing any large number of copies, to give them a chance to provide you with an updated version of the Document.

G.0.3 Modifications

You may copy and distribute a Modified Version of the Document under the conditions of sections [G.0.1](#) and [G.0.2](#) above, provided that you release the Modified Version under precisely this License, with the Modified Version filling the role of the Document, thus licensing distribution and modification of the Modified Version to whoever possesses a copy of it. In addition, you must do these things in the Modified Version:

- A. Use in the Title Page (and on the covers, if any) a title distinct from that of the Document, and from those of previous versions (which should, if there were any, be listed in the History section of the Document). You may use the same title as a previous version if the original publisher of that version gives permission.
- B. List on the Title Page, as authors, one or more persons or entities responsible for authorship of the modifications in the Modified Version, together with at least five of the principal authors of the Document (all of its principal authors, if it has fewer than five), unless they release you from this requirement.
- C. State on the Title page the name of the publisher of the Modified Version, as the publisher.
- D. Preserve all the copyright notices of the Document.
- E. Add an appropriate copyright notice for your modifications adjacent to the other copyright notices.
- F. Include, immediately after the copyright notices, a license notice giving the public permission to use the Modified Version under the terms of this License, in the form shown in the Addendum below.
- G. Preserve in that license notice the full lists of Invariant Sections and required Cover Texts given in the Document's license notice.
- H. Include an unaltered copy of this License.
- I. Preserve the section Entitled "History", Preserve its Title, and add to it an item stating at least the title, year, new authors, and publisher of the Modified Version as given on the Title Page. If there is no section Entitled "History" in the Document, create one stating the title, year, authors, and publisher of the Document as given on its Title Page, then add an item describing the Modified Version as stated in the previous sentence.
- J. Preserve the network location, if any, given in the Document for public access to a Transparent copy of the Document, and likewise the network locations given in the Document for previous versions it was based on. These may be placed in the "History" section. You may omit a network location for a work that was published at least four years before the Document itself, or if the original publisher of the version it refers to gives permission.
- K. For any section Entitled "Acknowledgements" or "Dedications", Preserve the Title of the section, and preserve in the section all the substance and tone of each of the contributor acknowledgements and/or dedications given therein.
- L. Preserve all the Invariant Sections of the Document, unaltered in their text and in their titles. Section numbers or the equivalent are not considered part of the section titles.
- M. Delete any section Entitled "Endorsements". Such a section may not be included in the Modified Version.
- N. Do not retile any existing section to be Entitled "Endorsements" or to conflict in title with any Invariant Section.
- O. Preserve any Warranty Disclaimers.

If the Modified Version includes new front-matter sections or appendices that qualify as Secondary Sections and contain no material copied from the Document, you may at your option designate some or all of these sections as invariant. To do this, add their titles to the list of Invariant Sections in the Modified Version's license notice. These titles must be distinct from any other section titles.

You may add a section Entitled "Endorsements", provided it contains nothing but endorsements of your Modified Version by various parties—for example, statements of peer review or that the text has been approved by an organization as the authoritative definition of a standard.

You may add a passage of up to five words as a Front-Cover Text, and a passage of up to 25 words as a Back-Cover Text, to the end of the list of Cover Texts in the Modified Version. Only one passage of Front-Cover Text and one of Back-Cover Text may be added by (or through arrangements made by) any one entity.

If the Document already includes a cover text for the same cover, previously added by you or by arrangement made by the same entity you are acting on behalf of, you may not add another; but you may replace the old one, on explicit permission from the previous publisher that added the old one.

The author(s) and publisher(s) of the Document do not by this License give permission to use their names for publicity for or to assert or imply endorsement of any Modified Version.

Combining documents

You may combine the Document with other documents released under this License, under the terms defined in section 4 above for modified versions, provided that you include in the combination all of the Invariant Sections of all of the original documents, unmodified, and list them all as Invariant Sections of your combined work in its license notice, and that you preserve all their Warranty Disclaimers.

The combined work need only contain one copy of this License, and multiple identical Invariant Sections may be replaced with a single copy. If there are multiple Invariant Sections with the same name but different contents, make the title of each such section unique by adding at the end of it, in parentheses, the name of the original author or publisher of that section if known, or else a unique number. Make the same adjustment to the section titles in the list of Invariant Sections in the license notice of the combined work.

In the combination, you must combine any sections Entitled “History” in the various original documents, forming one section Entitled “History”; likewise combine any sections Entitled “Acknowledgements”, and any sections Entitled “Dedications”. You must delete all sections Entitled “Endorsements.”

Collections of documents

You may make a collection consisting of the Document and other documents released under this License, and replace the individual copies of this License in the various documents with a single copy that is included in the collection, provided that you follow the rules of this License for verbatim copying of each of the documents in all other respects.

You may extract a single document from such a collection, and distribute it individually under this License, provided you insert a copy of this License into the extracted document, and follow this License in all other respects regarding verbatim copying of that document.

Aggregation with independent works

A compilation of the Document or its derivatives with other separate and independent documents or works, in or on a volume of a storage or distribution medium, is called an “aggregate” if the copyright resulting from the compilation is not used to limit the legal rights of the compilation’s users beyond what the individual works permit. When the Document is included in an aggregate, this License does not apply to the other works in the aggregate which are not themselves derivative works of the Document.

If the Cover Text requirement of section G.0.2 is applicable to these copies of the Document, then if the Document is less than one half of the entire aggregate, the Document’s Cover Texts may be placed on covers that bracket the Document within the aggregate, or the electronic equivalent of covers if the Document is in electronic form. Otherwise they must appear on printed covers that bracket the whole aggregate.

Translation

Translation is considered a kind of modification, so you may distribute translations of the Document under the terms of section G.0.3. Replacing Invariant Sections with translations requires special permission from their copyright holders, but you may include translations of some or all Invariant Sections in addition to the original versions of these Invariant Sections. You may include a translation of this License, and all the license notices in the Document, and any Warranty Disclaimers, provided that you also include the original English version of this License and the original versions of those notices and disclaimers. In case of a disagreement between the translation and the original version of this License or a notice or disclaimer, the original version will prevail.

If a section in the Document is Entitled “Acknowledgements”, “Dedications”, or “History”, the requirement (section G.0.3) to Preserve its Title (section G.0.0) will typically require changing the actual title.

Termination

You may not copy, modify, sublicense, or distribute the Document except as expressly provided for under this License. Any other attempt to copy, modify, sublicense or distribute the Document is void, and will automatically terminate your rights under this License. However, parties who have received copies, or rights, from you under this License will not have their licenses terminated so long as such parties remain in full compliance.

G GNU Free Documentation License

Future revisions of this license

The Free Software Foundation may publish new, revised versions of the GNU Free Documentation License from time to time. Such new versions will be similar in spirit to the present version, but may differ in detail to address new problems or concerns. See <http://www.gnu.org/copyleft/>.

Each version of the License is given a distinguishing version number. If the Document specifies that a particular numbered version of this License “or any later version” applies to it, you have the option of following the terms and conditions either of that specified version or of any later version that has been published (not as a draft) by the Free Software Foundation. If the Document does not specify a version number of this License, you may choose any version ever published (not as a draft) by the Free Software Foundation.

ADDENDUM: How to use this License for your documents

To use this License in a document you have written, include a copy of the License in the document and put the following copyright and license notices just after the title page:

Copyright (c) <year> <your name>. Permission is granted to copy, distribute and/or modify this document under the terms of the GNU Free Documentation License, Version 1.2 or any later version published by the Free Software Foundation; with no Invariant Sections, no Front-Cover Texts, and no Back-Cover Texts. A copy of the license is included in the section entitled “GNU Free Documentation License”.

If you have Invariant Sections, Front-Cover Texts and Back-Cover Texts, replace the “with...Texts.” line with this:

with the Invariant Sections being <list their titles>, with the Front-Cover Texts being <list>, and with the Back-Cover Texts being <list>.

If you have Invariant Sections without Cover Texts, or some other combination of the three, merge those two alternatives to suit the situation.

If your document contains nontrivial examples of program code, we recommend releasing these examples in parallel under your choice of free software license, such as the GNU General Public License, to permit their use in free software.

Copyright

Copyright (c) 2000, 2001, 2002 Free Software Foundation, Inc. 59 Temple Place, Suite 330, Boston, MA 02111-1307, USA

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

H A humorous disclaimer

The following text is quoted in part from an anonymous source ("Project Guten Tag") dating back to 1997. The original text is no longer available on the Internet.

WARRANTO LIMITENSIS; DISCLAMATANTUS DAMAGENSIS
Solut exceptus "Rectum Replacator Refundiens" describitus ecci,

1. Projectus (etque nunquam partum quis hic etext remittibus cum PROJECT GUTEN TAG-tm identificator) disclamabat omni liabilitus tuus damagensis, pecuniensisque, includibantus pecunia legalitus, et
2. REMEDIA NEGLIGENTITIA NON HABET TUUS, WARRANTUS DESTRUCTIBUS CONTRACTUS NULLIBUS NI LIABILITUS SUMUS, INCLUTATIBUS NON LIMITATUS DESTRUCTIO DIRECTIBUS, CONSEQUENTIUS, PUNITIO, O INCIDENTUS, NON SUNT SI NOS NOTIFICAT VOBIS.

Sit discubriatus defectus en etextum sic entram diariam noventam recibidio, pecuniam tuum refundatorium receptorus posset, sic scribatis vendor. Sit veniabat medium physicalis, vobis idem reternat et replacator possit copius. Sit venitabat electronicabilis, sic viri datus chansus secundibus.

HIC ETEXT VENID "COMO-ASI". NIHIL WARRANTI NUNQUAM CLASSUM, EXPRESSITO NI IMPLICATO, LE MACCHEN COMO SI ETEXTIO BENE SIT O IL MEDIO BENE SIT, INCLUTAT ET NON LIMITAT WARRANTI MERCATENSIS, APPROPRIATENSIS PURPOSEM.

Statuen varias non permitatent disclamabaris ni warranti implicatorem ni exclusioni limitatio damagaren consequentialis, ecco lo qua disclamatori exclusatorique non vobis applicant, et potat optia alia legali.

Index

- F*-algebra, [141](#)
- aggregation, [12](#)
- algebra, [141](#)
- algebraic, [141](#)
- algebraic data type, [139](#)
- applied functional type theory, [151](#)
- assembly language, [26](#)
- backward composition, [136](#)
- bound variable, [7](#)
- case expression, [32](#)
- co-product, [139](#)
- codensity monad, [104](#)
- computation expressions (F##), [140](#)
- constant combinator, [52](#)
- constant function, [51](#)
- contrafunctor choice, [93](#)
- curried arguments, [46](#)
- curried function, [45](#)
- Curry-Howard correspondence, [128](#), [129](#)
- destructuring, [31](#)
- dictionary, [26](#)
- disjunctive type, [139](#)
- do-notation (Haskell), [140](#)
- exponential-polynomial type, [141](#)
- expression, [1](#)
- expression block, [5](#)
- factorial, [1](#)
- flipped Kleisli, [94](#)
- for-comprehension, [140](#)
- forward composition, [136](#)
- free variable, [7](#)
- function value, [3](#), [48](#), [50](#)
- functional programming, [19](#)
- functional type theory, [122](#)
- functor block, [140](#)
- GADT (generalized algebraic data type), [140](#)
- higher-order function, [49](#)
- identity function, [51](#)
- infix method, [145](#)
- infix syntax, [6](#), [145](#)
- isomorphic, [46](#)
- Kleisli arrow, [140](#)
- Kleisli morphism, [140](#)
- Kleisli function, [140](#)
- local scope, [52](#)
- local value, [5](#)

Index

- map/reduce style, 12
- mathematical induction, 17, 23
- method, 50, 140, 145
- monad transformers
 - stacking, 71
- monads
 - stack of, 71
- name shadowing, 53
- nameless function, 3
- nameless function, 139
- order of a function, 49
- paradigm, 19
- partial application, 46
- partial function, 43
- pattern matching, 31, 32
 - infallible, 44
- pattern variables, 32
- polynomial functor, 139
- predicate, 5
- product type, 139
- pure compatibility laws, 88
- rigid functor, 105
- Scala method, 145
- search monad, 93
- selector monad, 103
- short type notation, 141
- tagged union, 139
- total function, 43
- transformation, 12
- tuple, 29
 - nested, 30
 - parts, 29
- type, 21
- type constructor, 145
- type error, 30
- type parameter, 145
- type parameters, 34
- uncurried function, 46
- unfunctor, 139
- variable, 20