

The Science of Functional Programming

A Tutorial, with Examples in Scala

by Sergei Winitzki, Ph.D.

draft version 0.3, May 30, 2019

Published by **lulu.com** 2019

Copyright © 2018-2019 by Sergei Winitzki.

Published and printed by lulu.com

ISBN XXXXXX

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).

The PDF file `sofp.pdf` contains the full source code of the book as an “attachment”. To extract the source file `soft-src.tar.bz2` to the current directory, run the command ‘`pdftk sofp.pdf unpack_files output .`’

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 as well as exercises. 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

Preface	1
1 Mathematical formulas as code. I. Nameless functions	3
1.1 Translating mathematics into code	3
1.1.1 First examples	3
1.1.2 Nameless functions	4
1.1.3 Nameless functions and bound variables	7
1.2 Aggregating data from sequences	10
1.3 Filtering and truncating sequences	13
1.4 Solved examples	14
1.4.1 Aggregation	14
1.4.2 Transformation	17
1.5 Summary	18
1.6 Exercises	20
1.6.1 Aggregation	20
1.6.2 Transformation	20
1.7 Discussion	21
1.7.1 Functional programming as a paradigm	21
1.7.2 Functional programming languages	22
1.7.3 The mathematical meaning of variables	22
1.7.4 Iteration without loops	24
1.7.5 Nameless functions in mathematical notation	25
1.7.6 Named and nameless expressions and their uses	27
1.7.7 Nameless functions: historical perspective	29
2 Mathematical formulas as code. II. Mathematical induction	31
2.1 Tuple types	31
2.1.1 Examples of using tuples	31
2.1.2 Pattern matching on tuples	33
2.1.3 Using tuples with collections	35
2.1.4 Using dictionaries (Scala's Maps) as sequences	36

Contents

2.1.5	Solved examples: Tuples and collections	41
2.1.6	Reasoning about types of sequences	47
2.1.7	Exercises: Tuples and collections	49
2.2	Converting a sequence into a single value	51
2.2.1	Inductive definitions of aggregation functions	52
2.2.2	Implementing functions by recursion	53
2.2.3	Tail recursion	55
2.2.4	Implementing a generic aggregation function (<code>foldLeft</code>)	61
2.2.5	Solved examples: using <code>foldLeft</code>	63
2.2.6	Exercises: Using <code>foldLeft</code>	69
2.3	Converting a single value into a sequence	71
2.4	Transforming a sequence into another sequence	74
2.5	Summary	75
2.5.1	Solved examples	76
2.5.2	Exercises	86
2.6	Discussion	90
2.6.1	Total and partial functions	90
2.6.2	Scope of pattern matching variables	91
2.6.3	Lazy values and sequences: Iterators and streams	92
3	The formal logic of types. I. Higher-order functions	99
3.1	Types of higher-order functions	99
3.1.1	Curried functions	99
3.1.2	Calculations with nameless functions	101
3.1.3	Short syntax for function applications	103
3.1.4	Higher-order functions	104
3.1.5	Worked examples: higher-order functions	105
3.2	Discussion	107
3.2.1	Scope of bound variables	107
3.3	Exercises	108
4	The formal logic of types. II. Disjunctive types	109
4.1	Discussion	109
4.1.1	Scala's case classes as "named tuple" types	109
5	The formal logic of types. III. The Curry-Howard correspondence	111
5.0.1	Discussion	111

6	Functors	113
6.1	Discussion	113
6.2	Practical use	113
6.3	Laws and structure	113
7	Type-level functions and typeclasses	115
7.1	Combining typeclasses	115
7.2	Inheritance	115
7.3	Functional dependencies	115
7.4	Discussion	115
8	Computations in functor blocks. I. Filterable functors	117
8.1	Practical use	117
8.1.1	Discussion	117
8.2	Laws and structure	117
8.2.1	Discussion	117
9	Computations in functor blocks. II. Semimonads and monads	119
9.1	Practical use	119
9.1.1	Discussion	119
9.2	Laws and structure	119
9.2.1	Discussion	119
10	Applicative functors, contrafunctors, and profunctors	121
10.1	Practical use	121
10.1.1	Discussion	121
10.2	Laws and structure	121
11	Traversable functors and profunctors	123
11.1	Discussion	123
12	“Free” type constructions	125
12.1	Discussion	125
13	Computations in functor blocks. III. Monad transformers	127
13.1	Practical use	127
13.2	Laws and structure	127
13.2.1	Laws of monad transformers	127
13.2.2	Examples of incorrect monad transformers	128

Contents

13.2.3	Examples of failure to define a generic monad transformer	129
13.2.4	Properties of monadic morphisms	131
13.2.5	Functor composition with transformed monads	133
13.2.6	Stacking two monads	133
13.2.7	Stacking any number of monads	137
13.3	Monad transformers via functor composition: General properties	138
13.3.1	Motivation for the swap function	139
13.3.2	Deriving the necessary laws for swap	141
13.3.3	Intuition behind the laws of swap	146
13.3.4	Deriving swap from flatten	147
13.3.5	Laws of monad transformer liftings: Proofs	151
13.3.6	Laws of monad transformer runners: Proofs	152
13.3.7	Summary of results	157
13.4	Composed-inside transformers: Linear monads	158
13.4.1	Definitions of swap and flatten	159
13.4.2	Laws of swap	160
13.4.3	Composition of transformers for linear monads	167
13.5	Composed-outside transformers: Rigid monads	167
13.5.1	Rigid monad construction 1: choice	168
13.5.2	Rigid monad construction 2: composition	184
13.5.3	Rigid monad construction 3: product	186
13.5.4	Rigid monad construction 4: selector	186
13.5.5	Rigid functors	186
13.6	Recursive monad transformers	189
13.6.1	Transformer for the free monad <code>FreeT</code>	189
13.6.2	Transformer for the list monad <code>ListT</code>	189
13.7	Monad transformers for monad constructions	189
13.7.1	Product of monad transformers	189
13.7.2	Free pointed monad transformer	189
13.8	Irregular and incomplete monad transformers	189
13.8.1	The state monad transformer <code>StateT</code>	189
13.8.2	The continuation monad transformer <code>ContT</code>	189
13.8.3	The codensity monad transformer <code>CodT</code>	189
13.9	Summary and discussion	190
14	Recursive types	191
14.1	Fixpoints and type recursion schemes	191
14.2	Row polymorphism and OO programming	191

14.3	Column polymorphism	191
14.4	Discussion	191
15	Co-inductive typeclasses. Comonads	193
15.1	Practical use	193
15.2	Laws and structure	193
15.3	Co-free constructions	193
15.4	Co-free comonads	193
15.5	Comonad transformers	193
15.6	Discussion	193
16	Irregular typeclasses*	195
16.1	Distributive functors	195
16.2	Monoidal monads	195
16.3	Lenses and prisms	195
16.4	Discussion	195
17	Summary and discussion	197
18	Essay: Software engineers and software artisans	199
18.1	Engineering disciplines	199
18.2	Artisanship: Trades and crafts	200
18.3	Programmers today are artisans, not engineers	200
18.3.1	No requirement of formal study	200
18.3.2	No mathematical formalism to guide software development	202
18.3.3	Programmers avoid academic terminology	203
18.4	Towards software engineering	204
18.5	Does software need engineers, or are artisans good enough?	206
19	Essay: Towards functional data engineering with Scala	209
19.1	Data is math	209
19.2	Functional programming is math	210
19.3	The power of abstraction	211
19.4	Scala is Java on math	212
19.5	Conclusion	213
20	“Applied functional type theory”: A proposal	215
20.1	AFTT is not found in computer science curricula	215
20.2	AFTT is not category theory, type theory, or formal logic	217

Contents

A	Notations	221
A.1	Summary table	221
A.2	Detailed explanations	223
B	Glossary of terms	227
B.1	On the current misuse of the term “algebra”	229
C	Scala syntax and features	233
C.0.1	Function syntax	233
C.0.2	Functions of several arguments vs. tuples	233
C.0.3	Scala collections	234
D	Intuitionistic propositional logic (IPL)	235
D.1	Example: The logic of types is not Boolean	235
D.2	Using truth values in Boolean logic and in IPL	237
E	Category theory	239
F	A humorous disclaimer	241
G	GNU Free Documentation License	243
G.0.0	Applicability and definitions	243
G.0.1	Verbatim copying	244
G.0.2	Copying in quantity	244
G.0.3	Modifications	244

Preface

The goal of this book is to teach programmers how to reason mathematically about types and code, in a way that is directly relevant to software practice.

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 unwilling to learn new and difficult concepts through prolonged mental concentration and effort.

The first two chapters are introductory and may be suitable for beginners in programming. Starting from the middle of Chapter 5, the material becomes unsuitable for beginners.

The presentation in this book is self-contained. I define 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 use some *non-standard* notations (Appendix A) and terminology (Appendix B) to achieve a clearer presentation of the material, especially for readers not familiar with today's research literature in the theory of programming languages.

The main intent of this book is to explain the mathematical principles that guide the practice of functional programming – that is, help people to write code. Therefore, all mathematical developments in this book are motivated and justified by practical programming issues, and are accompanied by code examples that illustrate their usage. For example, the equational laws for standard type-classes (functor, applicative, monad, etc.) are first motivated heuristically before deriving a set of mathematical equations and formulating the laws in more abstract terms. Each new concept or technique is explained via solved 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 examples and exercises are marked by an asterisk (*).

A software engineer needs to know only a little of mathematical theory – the parts that answer questions arising in the practice of functional programming. So I keep the theoretical material to the minimum; *ars longa, vita brevis est*. I do not pursue mathematical generalizations beyond practical relevance or im-

Contents

mediate pedagogical usefulness. This limits the scope of required mathematical knowledge to bare rudiments of category theory, type theory, and formal logic. For instance, I do not use “introduction/elimination rules”, “strong normalization”, “complete partial order domains”, “adjoint functors”, “limits”, “co-limits”, “pullbacks”, “pushouts”, “topoi”, or even the word “algebra”, because learning these concepts will not help a functional programmer write code. Instead, I focus on practically useful material – including some rarely mentioned constructions, e.g. the “filterable functor” and “applicative contrafunctor” typeclasses.

Some formatting conventions used in this book:

- Text in boldface indicates a new concept or term that is being defined. Text in italics is a logical emphasis. Example:

An **aggregation** is a function from a list of values to a *single* value.

- Sample Scala code is written inline using a small monospaced font, such as this: `val a = "xyz"`. Longer code examples are written in separate code blocks, which may also show the output from the Scala interpreter:

```
val s = (1 to 10).toList

scala> s.product
res0: Int = 3628800
```

- Derivations of laws are written in a two-column notation where the right column contains the code and the left column indicates the property or law used to derive the expression at right. A green underline shows the part of the *previous* expression that we rewrite using the indicated law. Example:

$$\begin{array}{ll} & \text{pu}_M^{\uparrow \text{Id}} \circ \text{pu}_M \circ \text{ftn}_M \\ \text{raising into the identity functor:} & = \text{pu}_M \circ \underline{\text{pu}_M \circ \text{ftn}_M} \\ \text{left identity law for } M : & = \text{pu}_M \quad . \end{array}$$

A green underline is sometimes also used at the *last* step of the derivation, to indicate the part of the expression that resulted from the most recent rewriting.

1 Mathematical formulas as code. I. Nameless 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`. To define a name for this value, we use the “`val`” syntax:

```
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 that (1) the code can be evaluated (e.g. using the Scala interpreter) and yields a value, and (2) the code can be inserted as a part of a larger expression. For example, we could write

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

1 Mathematical formulas as code. I. Nameless functions

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 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 a function does not really need a name, – for instance, if the function is used only once. I denote “nameless” mathematical functions like this:

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

1.1 Translating mathematics into code

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 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 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 (`Int`) argument and returns an integer result value. What is the value of the function `fac` *itself*? As we have just seen, the Scala interpreter prints `<function1>` as the “value” of `fac`. An alternative Scala interpreter² called `ammonite` 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.

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 doing 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.

²See <https://ammonite.io/>

1 Mathematical formulas as code. I. Nameless functions

I usually imagine that a “function value” represents a block of compiled machine code, – code that will actually run and evaluate the function’s body when the function is applied to its argument.

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

```
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
```

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 **locally scoped** definition.

Nameless functions are most useful when they are themselves arguments of

1.1 Translating mathematics into code

other functions, as we will see next.

Example: prime numbers Let us define a function that takes an integer argument n and determines whether n is a prime number.

A simple mathematical formula for this function can be written as

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

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, $n \neq 0 \bmod k$. Formula (1.1) is translated into 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,

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

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

1 Mathematical formulas as code. I. Nameless functions

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

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

does not seem to involve any nameless functions. Instead, the mathematical formula defines the symbol k that “goes over the range $[2, n-1]$,” as one might say. The symbol k is then used for writing the predicate $n \neq 0 \bmod k$.

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

1.1 Translating mathematics into code

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 (it is 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 z and write instead of Eq. (1.1) an equivalent definition

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

or in Scala code,

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

In the nameless function `k => n % k != 0`, the argument `k` may be renamed to `z` 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 `z`. 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 => n % k != 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 use bound variables in various constructions such as $\forall k : p(k)$, $\exists k : p(k)$, $\sum_{k=a}^b f(k)$, $\int_0^1 k^2 dk$, $\lim_{n \rightarrow \infty} f(n)$, and $\text{argmax}_k f(k)$. When

1 Mathematical formulas as code. I. Nameless functions

translating 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

1.2 Aggregating data from sequences

Consider the task of finding how many even numbers there are in a given list L of integers. For example, the list $[5, 6, 7, 8, 9]$ contains *two* even numbers: 6 and 8.

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

$$\begin{aligned} \text{count_even}(L) &= \sum_{k \in L} \text{is_even}(k) \quad , \\ \text{is_even}(k) &= \begin{cases} 1 & \text{if } k = 0 \bmod 2 \\ 0 & \text{otherwise} \end{cases} . \end{aligned}$$

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

1.2 Aggregating data from sequences

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 0 => 1 // First, check if it is zero.  
  case _ => 0 // The underscore matches everything else.  
}
```

For such a simple computation, we could also write shorter code using 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  
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, the argument of `.map` is the nameless function $x \Rightarrow x^2 + 100x$. This function is repeatedly applied by `.map` to transform each of the values from a given list, creating a new list as a result.

1 Mathematical formulas as code. I. Nameless functions

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 “dotless” 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
```

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`.

1.3 Filtering and truncating sequences

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

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, there is no mathematical notation for operations such as “removing elements from a list”, so we will focus on the Scala syntax for these functions.

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

1 Mathematical formulas as code. I. Nameless functions

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

The `.takeWhile` method truncates a given list, returning a *new list* with the initial portion of values from the original list for which predicate keeps being `true`:

```
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). A **transformation** is a function from a list of values to another list of values; examples of transformation functions are `.filter` and `.map`. An **aggregation** is 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 functions of transformation and aggregation is known as programming in the **map/reduce style**.

1.4 Solved examples

1.4.1 Aggregation

Example 1.4.1.1 Improve the code for `is_prime` by limiting the search to $k^2 \leq n$:

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

Solution: Use `.takeWhile` to truncate the initial list when $k^2 \leq n$ becomes false:

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

Example 1.4.1.2 Compute $\prod_{k \in [1,10]} |\sin(k+2)|$.

Solution:

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

Example 1.4.1.3 Compute $\sum_{k \in [1,10]; \cos k > 0} \sqrt{\cos k}$.**Solution:**

```
(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)
```

Example 1.4.1.4 Compute the average of a non-empty list of type `List[Double]`,

$$\text{average}(s) = \frac{1}{n} \sum_{i=0}^{n-1} s_i \quad .$$

Solution: We need to divide the sum by the length of the list:

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

scala> average(List(1.0, 2.0, 3.0))
res0: Double = 2.0
```

Example 1.4.1.5 Given n , compute the Wallis product truncated up to $\frac{2n}{2n+1}$:

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

Solution: We will define the helper function `wallis_frac(i)` that computes the i^{th} fraction. The method `.toDouble` converts integers to `Double` numbers.

1 Mathematical formulas as code. I. Nameless functions

```
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.  
res0: Double = 3.9267453954401036E-5  
  
scala> math.cos(wallis(100000)) // Should be even closer to 0.  
res1: 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 .

Example 1.4.1.6 Another known series related to π is

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

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.

Solution:

```
def euler_series(n: Int): Double =  
  (1 to n).map(k => 1.0/k/k).sum  
  
scala> euler_series(100000)  
res0: Double = 1.6449240668982423  
  
scala> val pi = 4*math.atan(1)  
pi: Double = 3.141592653589793  
  
scala> pi*pi/6  
res1: Double = 1.6449340668482264
```

Example 1.4.1.7 Check numerically the infinite product formula

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

Solution: We will compute this product up to $k = n$ for $x = 0.1$ 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(n = 100000, x = 0.1) // Arguments may be named, for clarity.
res0: Double = 0.9836317414461351

scala> math.sin(pi*0.1)/pi/0.1
res1: Double = 0.9836316430834658
```

Example 1.4.1.8 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.

Solution:

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

Here is an example test for this function:

```
scala> p(List(2, 3, 4, 5), x => 60 / x)
res0: Int = 30
```

1.4.2 Transformation

Example 1.4.2.1 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 $\text{List}[\text{List}[\text{Int}]]$.

Solution: To “select the inner lists” 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, $t \Rightarrow t.\text{size} \geq 3$:

```
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) ))
res0: List[List[Int]] = List(List(1, 2, 3), List(1, 2, 3, 4))
```

The predicate in the argument of `.filter` is a nameless function $t \Rightarrow t.\text{size} \geq 3$ whose argument t is of type $\text{List}[\text{Int}]$. 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.

1 Mathematical formulas as code. I. Nameless functions

Example 1.4.2.2 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 $j^2 > 2k$.

Solution:

```
scala> (1 to 10).toList.filter(k => (1 to k).
    | filter(j => j*j > 2*k).size >= 3)
res0: 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 compares the size of that list with 3 and so imposes the requirement that there should be at least 3 values of j . We can see how the Scala code closely follows the mathematical formulation of the problem.

1.5 Summary

The following table translates mathematical formulas into code.

Mathematical notation	Scala code
$x \Rightarrow \sqrt{x^2 + 1}$	<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>
$\sum_{k \in S \text{ such that } p(k) \text{ holds}} f(k)$	<code>s.filter(p).map(f).sum</code>

What problems can one 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:

$$r_n = \frac{1}{n} \sum_{k=0}^{n-1} s_k \quad .$$

- Example 3: Perform binary search over a sorted list of integers.

These computations involve *mathematical induction*, which we have not yet learned to 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.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 write the solution of Example 1 as a formula by using mathematical induction, but we have not yet seen how to implement that in Scala code.

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

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

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.

1.6 Exercises

1.6.1 Aggregation

Exercise 1.6.1.1 **Machin's formula** converges to π faster than Example 1.4.1.5:

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

$$\arctan \frac{1}{n} = \frac{1}{n} - \frac{1}{3} \frac{1}{n^3} + \frac{1}{5} \frac{1}{n^5} - \dots = \sum_{k=1}^{\infty} \frac{(-1)^k}{2k+1} n^{-2k-1} \quad .$$

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 .

Exercise 1.6.1.2 Using the function `is_prime`, check numerically the **Euler product formula** for the Riemann zeta function $\zeta(4)$; **it is known that** $\zeta(4) = \frac{\pi^4}{90}$:

$$\prod_{k \geq 2; k \text{ is prime}} \frac{1}{1 - p^{-4}} = \frac{\pi^4}{90} \quad .$$

1.6.2 Transformation

Exercise 1.6.2.1 Define a function `add_20` of type `List[List[Int]] => List[List[Int]]` that adds 20 to every element of every inner list. A sample test:

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

Exercise 1.6.2.2 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 $n \in [1, \dots, 1000]$.

Exercise 1.6.2.3 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 1.6.2.2 using that function.

Exercise 1.6.2.4 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. Test with:

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

Exercise 1.6.2.5 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.

1.7 Discussion

1.7.1 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 espe-

cially 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.

1.7.2 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, and PureScript, have 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, or PureScript will not play a significant role.

1.7.3 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 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 that 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 while 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

1 Mathematical formulas as code. I. Nameless functions

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 $z = \frac{x-y}{x+y}$ and now compute $\cos z + \cos 2z + \cos 3z$. Again, the variable z remains immutable, and its type remains fixed.

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 an integer. 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 use wrong types:

```
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.
              ^
```

1.7.4 Iteration without loops

Another distinctive feature of the FP paradigm is the absence of explicit loops.

Iterative computations are ubiquitous in mathematics; as an example, consider the formula for the standard deviation estimated from a sample,

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

To compute these expressions, we need to iterate over values of i and j . And yet, no mathematics textbook uses “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

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

Instead of loops, mathematicians write *expressions* such as $\sum_{i=1}^n s_i$, where symbols such as $\sum_{i=1}^n$ or $\prod_{i=1}^n$ denote iterative computations. Such computations are defined using mathematical induction. The functional programming 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. The machine code most likely *will* contain loops; but the functional programmer does not need to see that code or to reason about it.

1.7.5 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

1 Mathematical formulas as code. I. Nameless functions

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

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

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:

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

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:

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

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

```
def summation(a: Int, b: Int, g: Int => Double) = (a to b).map(g).sum

scala> summation(1, 10, x => math.sqrt(x))
res0: Double = 22.4682781862041
```

Numerical integration requires longer code, since the formulas are more complicated. For instance, **Simpson's rule** can be written as

$$\text{integration}(a, b, g) = \frac{\delta}{3} (g(a) + g(b) + 4s_1 + 2s_2) \quad ,$$

$$\text{where } n = 2 \left\lfloor \frac{b-a}{\varepsilon} \right\rfloor, \quad \delta_x = \frac{b-a}{n} \quad ,$$

$$s_1 = \sum_{i=1,3,\dots,n-1} g(a + i\delta_x) \quad ,$$

$$s_2 = \sum_{i=2,4,\dots,n-2} g(a + i\delta_x) \quad .$$

A straightforward translation of this formula into Scala is

```
def integration(a: Double, b: Double, g: Double => Double, eps: Double) = {
  // First, we define some helper values and functions that replace
  // the definitions "where n = ..." in the mathematical formula.
  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
  // Now we write the expression for 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 = 4.0000000000000003

scala> integration(0, 7, x => x*x*x*x*x*x, 0.001) // Exact answer is 117649
res1: Double = 117649.000000000023
```

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.

1.7.6 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 Mathematical formulas as code. I. Nameless 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 $(x+123)*y/(2+x)$ could 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).

So, programmers become more productive when their programming language supports nameless expressions.

This is also common practice in mathematics; names are assigned when needed, but most 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 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 us to build larger programs from simpler parts in a uniform way.

1.7.7 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, because programming in the map/reduce style (which invites frequent use of nameless functions) turned out to be immensely useful. Table 1.1 shows the year when nameless functions were introduced in each language.

What I call a “nameless function” is also 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 Mathematical formulas as code. I. Nameless 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
Scheme	1975	(lambda (k) (+ k k))
OCaml	1985	fun (k:int) -> k + k
Haskell	1990	\ k -> (k::Int) + k
Oz	1991	fun {\$ K} K + K
R	1993	function(k) k + k
Python 1.0	1994	lambda k: k + k
JavaScript	1995	function(k) { return k + k; }
Mercury	1995	func(K) = K + K
Ruby	1995	lambda { k k + k }
Lua 3.1	1998	function(k) return k + k end
Scala	2003	(k:Int) => k + k
F#	2005	fun (k:int) -> k + k
C# 3.0	2007	delegate(int k) { return 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 formulas as code.

II. Mathematical induction

We will now study more flexible ways of working with data collections in the functional programming paradigm. 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 focus of this chapter.

But first we need to become fluent in using tuple types with Scala collections.

2.1 Tuple types

2.1.1 Examples of using tuples

Many standard library methods in Scala require working with **tuple** types. A simple example of a tuple is a *pair* of values, – such as, 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.

A **triple** is defined in Scala like this:

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

Pairs and triples are examples of tuples. A **tuple** can contain any 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. Also, the number of parts in a tuple is fixed. 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;
```

2 Mathematical formulas as code. II. Mathematical induction

```
found    : Int(2)
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 accessor 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
      ^
```

Type errors are detected at compile time, before any computations begin.

Tuples can be **nested**: any part of a tuple can be itself 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
res0: Boolean = true

scala> c._2
res1: (String, Int) = (abc,3)
```


To define functions whose arguments are tuples, we could 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 accessor methods (`._1` and `._2`) to compute the 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 the function *f* to an argument *p* that has a wrong type, e.g. the type `(Int, Int)` instead of `(Boolean, Int)`.

2.1.2 Pattern matching on tuples

Instead of using accessor methods when working with tuples, it is often convenient to use **pattern matching**. Pattern matching occurs in two situations:

- destructuring definition: `val pattern = ...`
- case expression: `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, then assigns the names *x*, *y*, *z* to these parts. The types 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 far *undefined*. These names are called **pattern variables**. The destructuring definition checks whether the structure of the value of *g* “matches” the three pattern variables. (If *g* does not contain a tuple with exactly three parts, the definition will fail.) This computation is called **pattern matching**.

Pattern matching is often used when working with tuples:

2 Mathematical formulas as code. II. Mathematical induction

```
scala> (1, 2, 3) match { case (a, b, c) => a + b + c }  
res0: Int = 6
```

The **case expression** (`case (a, b, c) => ...`) performs pattern matching on the tuple argument `p`. The pattern matching will “destructure” (i.e. decompose) the tuple and try to match it to the given pattern `(a, b, c)`. In this pattern, `a, b, c` are as yet undefined new variables, – that is, they are pattern variables. If the pattern matching succeeds, the pattern variables `a, b, c` are assigned their values, and the function body can proceed to perform its computation. In this example, the pattern variables `a, b, c` 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 t1(p: (Int, (String, Int))): String = p match {  
  case (x, (str, y)) => str + (x + y).toString  
}  
  
scala> t1((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. If we rewrite the same code using the tuple accessor methods instead of pattern matching, the code will look like this:

```
def t2(p: (Int, (String, Int))): String = p._2._1 + (p._1 + p._2._2).toString
```

This code is shorter but harder to read: For example, it is not immediately clear what `p._2._1` refers to. It is also harder to change this code: Suppose we want to change the type of the tuple `p` to `((Int, String), Int)`. Then the new code is

```
def t3(p: ((Int, String), Int)): String = p._1._2 + (p._1._1 + p._2).toString
```

It takes time to verify, by going through every accessor method, that the function `t3` computes the same expression as `t2`. In contrast, the code is changed easily when using the pattern matching expression instead of the accessor methods:

```
def t4(p: ((Int, String), Int)): String = p match {  
  case ((x, str), y) => str + (x + y).toString  
}
```

2.1 Tuple types

The only change in the function body, compared to `t1`, is in the pattern matcher, so it is visually clear that `t4` computes the same expression as `t1`.

Sometimes we do not need some of the tuple parts in a pattern match. The following syntax is used to make this intention clear:

```
scala> val (x, _, _, z) = ("abc", 123, false, true)
x: String = abc
z: Boolean = true
```

The underscore symbol `_` denotes the parts of the pattern that we want to ignore. The underscore will always match any value regardless of type.

A feature of Scala is a short syntax for functions such as `{case (x, y) => y}` that extract elements from tuples. The shorter syntax is `(t => t._2)` or even shorter, `(_. _2)`, as illustrated here:

```
scala> val p: ((Int, Int)) => Int = { case (x, y) => y }
p: ((Int, Int)) => Int = <function1>

scala> p((1, 2))
res0: Int = 2

scala> val q: ((Int, Int)) => Int = (_. _2)
q: ((Int, Int)) => Int = <function1>

scala> q((1, 2))
res1: Int = 2

scala> Seq( (1,10), (2,20), (3,30) ).map(_. _2)
res2: Seq[Int] = List(10, 20, 30)
```

2.1.3 Using tuples with collections

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

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

We can create a list of tuples,

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

We could define a tuple of lists of tuples of functions, or any other combination.

2 Mathematical formulas as code. II. Mathematical induction

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. The names “fruit”, “count” are chosen to help us remember the meaning of the parts of tuples.

We can easily transform a list of tuples into a list of values of a 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))
```

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 Maps) as sequences

In the Scala standard library, tuples are frequently used as types of intermediate values. For instance, tuples are used when iterating over dictionaries. 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**. Type parameters represent unknown types that will be chosen later, when working with values having specific types.

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

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

2.1 Tuple types

This is equivalent to first creating a sequence of key/value *pairs* and then converting that sequence into a dictionary.

Pairs are used often, so the Scala library defines a special infix syntax for pairs via 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, it is easier to read the code for creating a dictionary:

```
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 collection types such as `Seq`, `Vector`, `Stream`, and `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 in the Scala library. All these data structures are gathered under the common “sequence” type called `Seq`. The methods defined in the Scala standard library sometimes return different implementations of the `Seq` type for reasons of performance.

The standard library has several useful methods that use tuple types, such as `.map` and `.filter` (with dictionaries), `.toMap`, `.zip`, and `.zipWithIndex`. The methods `.flatten`, `.flatMap`, `.groupBy`, and `.sliding` also work with most collection types, including dictionaries and sets. It is important to become familiar with these methods, because it will help writing code that uses sequences, sets, and dictionaries. Let us now look at these methods one by one.

The `.map` and `.toMap` methods Chapter 1 showed how the `.map` method works on sequences: the expression `xs.map(f)` applies a given function `f` to each element of the sequence `xs`, gathering the results in a new sequence. In this sense, we can say that the `.map` method “iterates over” sequences. The `.map` method works similarly on dictionaries, except that iterating over a dictionary of type `Map[K, V]` when applying `.map` looks like iterating over a sequence of *pairs*, `Seq[(K,V)]`. If `d: Map[K,V]` is a dictionary, the argument `f` of `d.map(f)` must be a function operating on tuples of type `(K,V)`. Typically, such functions are written using `case` expressions:

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

scala> m1.map { case (fruit, count) => count * 2 }
```

2 Mathematical formulas as code. II. Mathematical induction

```
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 a dictionary with the `.toMap` method:

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

The `.filter` method works on dictionaries by iterating on key/value pairs. The filtering predicate must be a function of type `((K, V)) => Boolean`. For example:

```
scala> m1.filter { case (fruit, count) => count > 0 }.toMap
res2: Map[String,Int] = Map(apples -> 6, pears -> 4)
```

The `.zip` and `.zipWithIndex` methods 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, the equivalent “dotless” infix syntax (`s zip t`) is shown just to illustrate the flexibility of syntax conventions in Scala.

The `.zip` method works equally well on dictionaries: in that case, dictionaries are automatically converted to sequences of tuples before applying `.zip`.

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)
```

The “flattening” operation computes the concatenation of all inner sequences. In Scala, sequences are concatenated using the operation `++`, e.g.:

```
scala> List(1, 2, 3) ++ List(4, 5, 6) ++ List(0)
res7: List[Int] = List(1, 2, 3, 4, 5, 6, 0)
```

So the `.flatten` method inserts the operation `++` between all the inner sequences.

Keep in mind that `.flatten` removes *only one* level of nesting, which is at the “outside” of the data structure. If applied to a `List[List[List[Int]]]`, the `.flatten` method returns a `List[List[Int]]`:

```
scala> List(List(List(1), List(2)), List(List(2), List(3))).flatten
res8: 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, equivalent to first applying `.map` and then `.flatten`:

```
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)

scala> List(1,2,3,4).flatMap(n => (1 to n).toList)
res11: List[Int] = List(1, 1, 2, 1, 2, 3, 1, 2, 3, 4)
```

The `.flatMap` operation transforms a sequence by mapping each element to a potentially different number of new elements.

At first sight, it may be unclear why `.flatMap` is useful. (Should we perhaps combine `.filter` and `.flatten` into a `.flatMapFilter`, or combine `.zip` and `.flatten` into a `.flatMapZip`?) However, we will see later in this book that the use of `.flatMap`, which is related to “monads”, is one of the most versatile and powerful design patterns in functional programming. In this chapter, several examples and exercises will illustrate the use of `.flatMap` for working on sequences.

The `.groupBy` method rearranges a sequence into a dictionary where some elements of the original sequence are grouped together into subsequences. For example, given a sequence of words, we can group all words that start with the letter “y” into one subsequence, and all other words into another subsequence. This is accomplished by the following code,

```
scala> Seq("wombat", "xanthan", "yoghurt", "zebra").
```

2 Mathematical formulas as code. II. Mathematical induction

```
groupBy(s => if (s startsWith "y") 1 else 2)
res12: Map[Int,Seq[String]] = Map(1 -> List(yoghurt), 2 -> List(wombat, xanthan,
zebra))
```

The argument of the `.groupBy` method is a *function* that computes a “key” out of each sequence element. The key can have an arbitrarily chosen type. (In the current example, that type is `Int`.) The result of `.groupBy` is a dictionary that maps each key to the sub-sequence of values that have that key. (In the current example, the type of the dictionary is therefore `Map[Int, Seq[String]]`.) The order of elements in the sub-sequences remains the same as in the original sequence.

As another example of using `.groupBy`, the following code will group together all numbers that have the same remainder after division by 3:

```
scala> List(1,2,3,4,5).groupBy(k => k % 3)
res13: Map[Int,List[Int]] = Map(2 -> List(2, 5), 1 -> List(1, 4), 0 -> List(3))
```

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

```
scala> (1 to 10).sliding(4).toList
res14: List[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))
```

Usually, this method is used together with an aggregation operation on the inner sequences. For example, the following code computes a sliding-window average with window width 50 over an array of 100 numbers:

```
scala> (1 to 100).map(x => math.cos(x)).sliding(50).
map(_._sum / 50).take(5).toList
res15: List[Double] = List(-0.005153079196990285, -0.0011160413780774369,
0.003947079736951305, 0.005381273944717851, 0.0018679497047270743)
```

The `.sortBy` method sorts a sequence according to a sorting key. The argument of `.sortBy` is a *function* that computes the sorting key from a sequence element. In this way, we can sort elements in an arbitrary way:

```
scala> Seq(1, 2, 3).sortBy(x => -x)
res0: Seq[Int] = List(3, 2, 1)

scala> Seq("z", "xxx", "yy").sortBy(word => word)
res1: Seq[String] = List("xxx", "yy", "z")
```



```
scala> Seq("z", "xxx", "yy").sortBy(word => word.length)
res2: Seq[String] = List("z", "yy", "xxx")
```

Sorting by the elements themselves, as we have done here with `.sortBy(word => word)`, is only possible if the element's type has a well-defined ordering. For strings, this is the alphabetic ordering, and for integers, the standard arithmetic ordering. For such types, a convenience method `.sorted` is defined, and works equivalently to `.sortBy(x => x)`:

```
scala> Seq("z", "xxx", "yy").sorted
res3: Seq[String] = List("xxx", "yy", "z")
```

2.1.5 Solved 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 corresponding to each element of the original sequence. This transformation is exactly what the `.map` method does. So the code is

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

Example 2.1.5.2 Count how many times $\cos x_i > \sin x_i$ occurs in a sequence x_i .

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 the elements of the sequence:

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

We could also reuse the solution of Exercise 2.1.5.1 that computed the cosine and the sine values. The code would then become

```
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`, and assume `as` and `bs` are of type `Seq[Double]`.

2 Mathematical formulas as code. II. Mathematical induction

Solution: We can use `.zip` on `as` and `bs`, which gives a sequence of pairs,

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

We then compute the differences $a_i - b_i$ by applying `.map` to this sequence:

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

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

Hint: use `.zip` and `.tail`.

Solution: Given `ps:Seq[Double]`, we can compute `ps.tail`. The result is a sequence that is 1 element shorter than `ps`, for example:

```
scala> val ps = Seq(1,2,3,4)
ps: Seq[Int] = List(1, 2, 3, 4)

scala> ps.tail
res0: Seq[Int] = List(2, 3, 4)
```

Taking a `.zip` of the two sequences `ps` and `ps.tail`, we get a sequence of pairs:

```
scala> ps.zip(ps.tail)
res1: Seq[(Int, Int)] = List((1,2), (2,3), (3,4))
```

Note that `ps.tail` is 1 element shorter than `ps`, and the resulting sequence of pairs is also 1 element shorter than `ps`. In other words, it is not necessary to truncate `ps` before computing `ps.zip(ps.tail)`. Now apply the `.count` method:

```
ps.zip(ps.tail).count { case (a, b) => a > b }
```

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

Hint: use `.sliding`.

Solution: Applying the `.sliding` method to a list gives a list of nested lists:

```
scala> val bs = List(1, 2, 3, 4, 5)
bs: List[Int] = List(1, 2, 3, 4, 5)

scala> bs.sliding(3).toList
res0: List[List[Int]] = List(List(1, 2, 3), List(2, 3, 4), List(3, 4, 5))
```

For each b_i , we need to obtain a list of $2k + 1$ nearby elements $(b_{i-k}, \dots, b_{i+k})$. So we need to use `.sliding(2*k+1)` to obtain a window of the required size. Now we can compute the maximum of each of the nested lists by using the `.map` method on

the outer list, with the `.max` method applied to the nested lists. So the argument of the `.map` method must be the function `nested => nested.max`. The final code is

```
bs.sliding(2 * k + 1).map(nested => nested.max)
```

In Scala, this code can be written more concisely using the syntax

```
bs.sliding(2 * k + 1).map(_.max)
```

because `_.max` means the nameless function `x => x.max`.

Example 2.1.5.6 Create a 10×10 multiplication table as a dictionary of type `Map[(Int, Int), Int]`. For example, a 3×3 multiplication table would be given by this dictionary,

```
Map( (1, 1) -> 1, (1, 2) -> 2, (1, 3) -> 3, (2, 1) -> 2,
      (2, 2) -> 4, (2, 3) -> 6, (3, 1) -> 3, (3, 2) -> 6, (3, 3) -> 9 )
```

Hint: use `.flatMap` and `.toMap`.

Solution: We are required to make a dictionary that maps pairs of integers (x, y) to $x * y$. Begin by creating the list of *keys* for that dictionary, which must be a list of pairs (x, y) of the form `List((1,1), (1,2), ..., (2,1), (2,2), ...)`. We need to start with a sequence of values of x , and for each x from that sequence, iterate over another sequence to provide values for y . Try this computation:

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

We would like to get `List((1,1), (1,2), (1,3))` etc., and so we use `.map` on the inner list with a nameless function `y => (x, y)` that converts a number into a tuple,

```
scala> List(1, 2, 3).map{ y => (x, y) }
res0: List[(Int, Int)] = List((1,1), (1,2), (1,3))
```

Here the curly braces `{ y => (x, y) }` are used only for clarity; we could equivalently use parentheses and write `(y => (x, y))`.

Using this `.map` operation, we obtain the code for a nested list of tuples:

```
scala> val s = List(1, 2, 3).map(x => List(1, 2, 3).map{ y => (x, y) })
s: List[List[(Int, Int)]] = List(List((1,1), (1,2), (1,3)), List((2,1), (2,2), (2,3)), List((3,1), (3,2), (3,3)))
```

This is almost what we need, except that the nested lists need to be concatenated

2 Mathematical formulas as code. II. Mathematical induction

into a single list. This is exactly what `.flatten` does:

```
scala> val s = List(1, 2, 3).map(x => List(1, 2, 3).map{ y => (x, y) }).flatten
s: List[(Int, Int)] = List((1,1), (1,2), (1,3), (2,1), (2,2), (2,3), (3,1),
  (3,2), (3,3))
```

It is shorter to write `.flatMap(...)` instead of `.map(...).flatten`:

```
scala> val s = List(1, 2, 3).flatMap(x => List(1, 2, 3).map{ y => (x, y) })
s: List[(Int, Int)] = List((1,1), (1,2), (1,3), (2,1), (2,2), (2,3), (3,1),
  (3,2), (3,3))
```

This is the list of keys for the required dictionary, which should map each *pair* of integers (x,y) to $x*y$. In Scala, a dictionary is usually created by applying `.toMap` to a sequence of pairs (key,value). So we need to create a sequence of nested tuples of the form $((x,y), \text{product})$. To achieve this, we use `.map` with a function that computes the product and creates a nested tuple:

```
scala> val s = List(1, 2, 3).flatMap(x => List(1, 2, 3).
  map{ y => (x, y) }).map{ case (x, y) => ((x, y), x * y) }
s: List[((Int, Int), Int)] = List(((1,1),1), ((1,2),2), ((1,3),3), ((2,1),2),
  ((2,2),4), ((2,3),6), ((3,1),3), ((3,2),6), ((3,3),9))
```

We can simplify this code if we notice that we are first mapping each y to a tuple (x, y) , and later map each tuple (x,y) to a nested tuple $((x,y), x*y)$. Instead, the entire computation can be done in the inner `.map` operation:

```
scala> val s = List(1, 2, 3).flatMap(x => List(1, 2, 3).
  map{ y => ((x,y), x*y) })
s: List[((Int, Int), Int)] = List(((1,1),1), ((1,2),2), ((1,3),3), ((2,1),2),
  ((2,2),4), ((2,3),6), ((3,1),3), ((3,2),6), ((3,3),9))
```

It remains to convert this list of tuples to a dictionary with `.toMap`. Also, for better readability, we can use Scala's pair syntax `key -> value`, which is completely equivalent to writing the tuple $(\text{key}, \text{value})$. The final code is

```
(1 to 10).flatMap(x => (1 to 10).map{ y => (x,y) -> x*y }).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 `.flatMap`, `.max`.

Solution: We will compute the required value if we take `.max` of a list containing all of the numbers. To do that, first map each element of the list `xs:Seq[Double]` into a sequence of three numbers:

```
scala> List(0.1, 0.5, 0.9)
res0: List[Double] = List(0.1, 0.5, 0.9)

scala> res0.map { x => Seq(x, math.cos(x), math.sin(x)) }
res1: List[Seq[Double]] = List(List(0.1, 0.9950041652780258,
  0.09983341664682815), List(0.5, 0.8775825618903728, 0.479425538604203),
  List(0.9, 0.6216099682706644, 0.7833269096274834))
```

This list is almost what we need, except we need to `.flatten` it:

```
scala> res1.flatten
res2: List[Double] = List(0.1, 0.9950041652780258, 0.09983341664682815, 0.5,
  0.8775825618903728, 0.479425538604203, 0.9, 0.6216099682706644,
  0.7833269096274834)
```

Now we just need to take the maximum of the resulting numbers:

```
scala> res2.max
res3: Double = 0.9950041652780258
```

The final code (starting from a given sequence `xs`) is

```
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 `.map` and `.toMap`.

Solution: Keep in mind that iterating over a dictionary looks like iterating over a list of (key, value) pairs, and use `.map` to reverse each pair:

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

Example 2.1.5.9 Write the solution of Example 2.1.5.8 as a function with type parameters `Name` and `Addr` instead of the fixed type `String`.

Solution: In Scala, the syntax for type parameters in a function definition is

```
def rev[Name, Addr](...) = ...
```

The type of the argument is `Map[Name, Addr]`, while the type of the result is `Map[Addr, Name]`. So we use the type parameters `Name` and `Addr` in the type signature of the function. The final code is

2 Mathematical formulas as code. II. Mathematical induction

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

The body of the function `rev` remains the same as in Example 2.1.5.8; only the type signature changed. This is because the procedure for reversing a dictionary works in the same way for dictionaries of any type. So the body of the function `rev` does not actually need to know the types of the keys and values in the dictionary. For this reason, it was easy for us to change the specific types (`String`) into type parameters in this function.

When the function `rev` is applied to a dictionary of a specific type, the Scala compiler will automatically set the type parameters `Name` and `Addr` that fit the required types of the dictionary's keys and values. For example, if we apply `rev` to a dictionary of type `Map[Boolean, Seq[String]]`, the type parameters will be set automatically as `Name = Boolean` and `Addr = Seq[String]`:

```
scala> val d = Map(true -> Seq("x", "y"), false -> Seq("z", "t"))  
d: Map[Boolean, Seq[String]] = Map(true -> List(x, y), false -> List(z, t))  
  
scala> rev(d)  
res0: Map[Seq[String], Boolean] = Map(List(x, y) -> true, List(z, t) -> false)
```

Type parameters can be also set explicitly when using the function `rev`. If the type parameters are chosen incorrectly, the program will not compile:

```
scala> rev[Boolean, Seq[String]](d)  
res1: Map[Seq[String], Boolean] = Map(List(x, y) -> true, List(z, t) -> false)  
  
scala> rev[Int, Double](d)  
<console>:14: error: type mismatch;  
found    : Map[Boolean,Seq[String]]  
required: Map[Int,Double]  
    rev[Int, Double](d)  
                        ^
```

Example 2.1.5.10* Given a sequence `words:Seq[String]` of words, compute a sequence of type `Seq[(Seq[String], Int)]`, where each inner sequence contains all the words having the same length, paired with the integer value showing that length. So, the input `Seq("the", "food", "is", "good")` should produce the output

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

The resulting sequence must be ordered by increasing length of words.

Solution: It is clear that we need to begin by grouping 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` if 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 list of words having 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 need. If we convert this dictionary to a sequence, we will get 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 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) }
```

This code can be written somewhat shorter if we use the syntax `_. _1` for selecting the first part of the tuple, and `.swap` for swapping the two elements of the pair:

```
words.groupBy(_.length).toSeq.sortBy(_. _1).map(_.swap)
```

However, the program may now be harder to read and to modify.

2.1.6 Reasoning about types of sequences

In Example 2.1.5.10 we have applied a chain of operations to a sequence. Let us add comments showing the type of the intermediate result after each operation:

2 Mathematical formulas as code. II. Mathematical induction

```
words // Seq[String]
  .groupBy(_.length) // Map[Int, Seq[String]]
  .toSeq // Seq[ (Int, Seq[String]) ]
  .sortBy { case (len, words) => len } // Seq[ (Int, Seq[String]) ]
  .map { case (len, words) => (words, len) } // Seq[ (Seq[String], Int) ]
```

In computations like this, the Scala compiler verifies at each step that the operations are applied to values of the correct type.

For instance, `.sortBy` is defined for sequences but not for dictionaries, so it would be a type error to apply `.sortBy` to a dictionary without first converting it to a sequence using `.toSeq`. The type of the intermediate result after `.toSeq` is `Seq[(Int, Seq[String])]`, and the `.sortBy` operation is applied to that sequence. So the sequence element matched by `{ case (len, words) => len }` is a tuple `(Int, Seq[String])`, which means that the pattern variables `len` and `words` must have types `Int` and `Seq[String]` respectively. It would be a type error to use the sorting key function `{ case (len, words) => words }`: the sorting key can be an integer `len`, but not a string sequence `words` (because sorting by string sequences is not defined).

If we visualize how the type of the sequence should change at every step, we can more quickly understand how to implement the required task. Begin by writing down the intermediate types that would be needed during the computation:

```
words: Seq[String] // Need to group by word length.
Map[Int, Seq[String]] // Need to sort by word length; can't sort a dictionary!
// Need to convert this dictionary to a sequence:
Seq[ (Int, Seq[String]) ] // Now sort this! Sorting does not change the type.
// It remains to swap the parts of all tuples in the sequence:
Seq[ (Seq[String], Int) ] // We are done.
```

Having written down these types, we are better assured that the computation can be done correctly. Writing the code becomes straightforward, since we are guided by the already known types of the intermediate results:

```
words.groupBy(_.length).toSeq.sortBy(_._1).map(_._swap)
```

This example illustrates the main benefits of reasoning about types: it gives direct guidance about how to organize the computation, together with a greater assurance in the correctness of the code.

2.1.7 Exercises: Tuples and collections

Exercise 2.1.7.1 Find all pairs i, j within $(0, 1, \dots, 9)$ such that $i + 4 * j > i * j$.

Hint: use `.flatMap` and `.filter`.

Exercise 2.1.7.2 Same task as in Exercise 2.1.7.1, but for i, j, k and the condition $i + 4 * j + 9 * k > i * j * k$.

Exercise 2.1.7.3 Given two sequences `p: Seq[String]` and `q: 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 2.1.7.4 Convert a `Seq[Int]` into a `Seq[(Int, Boolean)]` where the `Boolean` value is `true` when the element is followed by a larger value. For example, `Seq(1,3,2,4)` is to be converted into `Seq((1,true), (3,false), (2,true), (4,false))`. (The last element, 4, has no following element.)

Exercise 2.1.7.5 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 2.1.7.6 Write the solution of Exercise 2.1.7.5 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 2.1.7.7 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`. For example, if `p = Seq("a", "b", "c")` and `q = Seq(10, -1, 5)` then the result must be `Seq("b", "c", "a")`.

Exercise 2.1.7.8 Write the solution of Exercise 2.1.7.7 as a function with type parameter `s` instead of the fixed type `String`. The required type signature and a sample test:

```
def reorder[S](p: Seq[S], q: Seq[Int]): Seq[S] = ???

scala> reorder(Seq(6.0,2.0,8.0,4.0), Seq(20,10,40,30))
res0: Seq[Double] = List(2.0, 6.0, 4.0, 8.0)
```

Exercise 2.1.7.9 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

2 Mathematical formulas as code. II. Mathematical induction

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

the output must be

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

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

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

Exercise 2.1.7.10 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 `.map`, `.sortBy`, `.take`.

Exercise 2.1.7.11 (a) Given two sets `p:Set[Int]` and `q:Set[Int]`, compute a set of type `Set[(Int, Int)]` as 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`.

(b) Implement this computation as a function with type parameters `I, J` instead of `Int`. The required type signature and a sample test:

```
def cartesian[I,J](p: Set[I], q: Set[J]): Set[(I, J)] = ???

scala> cartesian(Set("a", "b"), Set(10, 20))
res0: Set[(String, Int)] = Set((a,10), (a,20), (b,10), (b,20))
```

Hint: use `.flatMap` and `.map` on sets.

Exercise 2.1.7.12* 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 that `Person` and `Amount` are type parameters. The required type signature and a sample test:

```
def payments[Person, Amount](data: Seq[Map[Person, Amount]]): Map[Person,
  Seq[Amount]] = ???

scala> payments(Seq(Map("Tarski" -> 10, "Church" -> 20), Map("Church" -> 100,
  "Gentzen" -> 40), Map("Tarski" -> 50)))
res0: Map[String, Seq[Int]] = Map(Gentzen -> List(40), Church -> List(20, 100),
  Tarski -> List(10, 50))
```

Hint: use `.flatMap`, `.groupBy`, `.mapValues` on dictionaries.

2.2 Converting a sequence into a single value

Until this point, we have been working with sequences using methods such as `.map` and `.zip`. These techniques are powerful but still insufficient for solving certain problems.

A simple computation that is impossible to do using `.map` is to compute the sum of a sequence of numbers. The standard library method `.sum` already does this; but we cannot implement `.sum` ourselves by using `.map`, `.zip`, or `.filter`. These operations always compute *new sequences*, while we need to compute a single value (the sum of all elements) from a sequence.

We have seen a few library methods such as `.count`, `.length`, and `.max` that compute a single value from a sequence; but we still cannot implement `.sum` using these methods. What we need is a more general way of converting a sequence to a single value, such that we could ourselves implement `.sum`, `.count`, `.max`, and other similar computations.

Another task not solvable with `.map`, `.sum`, etc., is to compute a floating-point number from a given sequence of decimal digits (including a “dot” character):

```
def digitsToDouble(ds: Seq[Char]): Double = ???

scala> digitsToDouble(Seq('2', '0', '4', '.', '5'))
res0: Double = 204.5
```

Why is it impossible to implement this function using `.map`, `.sum`, `.zip` and other methods we have seen so far? In fact, the same task for *integer* numbers (not for floating-point numbers) is solvable using `.length`, `.map`, `.sum`, and `.zip`:

```
def digitsToInt(ds: Seq[Int]): Int = {
  val n = ds.length
  // Compute a sequence of powers of 10, e.g. [1000, 100, 100, 1]
  val powers: Seq[Int] = (0 to n-1).map(k => math.pow(10, n-1-k).toInt)
  // Sum the powers of 10 with coefficients from 'ds'.
  (ds zip powers).map { case (d, p) => d * p }.sum
}

scala> digitsToInt(Seq(2,4,0,5))
res0: Int = 2405
```

2 Mathematical formulas as code. II. Mathematical induction

The computation can be written as the formula

$$r = \sum_{k=0}^{n-1} d_k * 10^{n-1-k} .$$

The sequence of powers of 10 can be computed separately and “zipped” with the sequence of digits d_k . However, for floating-point numbers, the sequence of powers of 10 depends on the position of the “dot” character. Methods such as `.map` and `.zip` cannot compute a sequence whose next elements are not known in advance but depend on previous elements via a custom function.

2.2.1 Inductive definitions of aggregation functions

Mathematical induction is a general way of expressing the dependence of next values on previously computed values. To define a function from sequence to a single value (e.g. an aggregation function $f : \text{Seq}[\text{Int}] \Rightarrow \text{Int}$) by using mathematical induction, we need to specify two computations:

- (The **base case** of the induction.) We need to specify what value the function f returns for an empty sequence, `Seq()`. If the function is only defined for non-empty sequences, we need to specify what the function f returns for a one-element sequence such as `Seq(x)`, with any x .
- (The **inductive step**.) Assuming that the function f is already computed for some sequence xs (the **inductive assumption**), how to compute the function f for a sequence with one more element x ? The sequence with one more element is written as `xs ++ Seq(x)`. So, we need to specify how to compute $f(xs ++ Seq(x))$ assuming that $f(xs)$ is already known.

Once these two computations are specified, the function f is defined (and can in principle be computed) for an arbitrary input sequence. This is how induction works in mathematics, and it works in the same way in functional programming. With this approach, the inductive definition of the method `.sum` looks like this:

- The sum of an empty sequence is 0. That is, `Seq().sum = 0`.
- If the result is already known for a sequence `xs`, and we have a sequence that has one more element x , the new result is equal to `xs.sum + x`. In code, this is `(xs ++ Seq(x)).sum = xs.sum + x`.

2.2 Converting a sequence into a single value

The inductive definition of the function `digitsToInt` is:

- For an empty sequence of digits, `Seq()`, the result is 0. This is a convenient base case, even if we never call `digitsToInt` on an empty sequence.
- If `digitsToInt(xs)` is already known for a sequence `xs` of digits, and we have a sequence `xs ++ Seq(x)` with one more digit `x`, then

```
digitsToInt(xs ++ Seq(x)) = digitsToInt(xs) * 10 + x
```

Let us write inductive definitions for methods such as `.length`, `.max`, and `.count`:

- The length of a sequence:
 - for an empty sequence, `Seq().length = 0`
 - if `xs.length` is known then `(xs ++ Seq(x)).length = xs.length + 1`
- Maximum element of a sequence (undefined for empty sequences):
 - for a one-element sequence, `Seq(x).max = x`
 - if `xs.max` is known then `(xs ++ Seq(x)).max = math.max(xs.max, x)`
- Count the sequence elements satisfying a predicate `p`:
 - for an empty sequence, `Seq().count(p) = 0`
 - if `xs.count(p)` is known then `(xs ++ Seq(x)).count(p) = xs.count(p) + c`, where `c = 1` when `p(x) == true` and `c = 0` otherwise

There are two main ways of translating mathematical induction into code. The first way is to write a recursive function. The second way is to use a standard library function, such as `foldLeft` or `reduce`. Most often it is better to use the standard library functions, but sometimes the code is more transparent when using explicit recursion. So let us consider each of these ways in turn.

2.2.2 Implementing functions by recursion

A **recursive function** is any function that calls itself somewhere within its own body. The call to itself is the **recursive call**.

When the body of a recursive function is evaluated, it may repeatedly call itself with different arguments until the result value can be computed *without* any recursive calls. The last recursive call corresponds to the base case of the induction. It is an error if the base case is never reached, as in this example:

2 Mathematical formulas as code. II. Mathematical induction

```
scala> def infiniteLoop(x: Int): Int = infiniteLoop(x+1)
infiniteLoop: (x: Int)Int

scala> infiniteLoop(2) // You will need to press Ctrl-C to stop this.
```

We translate mathematical induction into code by first writing a condition to decide whether we are in the base case or in the inductive step. As an example, consider how we would define `.sum` by recursion. The base case returns 0, and the inductive step returns a value computed from the recursive call:

```
def sum(s: Seq[Int]): Int = if (s == Seq()) 0 else {
  val x = s.head // To split s = Seq(x) ++ xs, compute x
  val xs = s.tail // and xs.
  sum(prev) + next // Call sum(...) recursively.
}
```

In this example, we use the `if/else` expression to separate the base case from the inductive step. In the inductive step, we split the given sequence `s` into a single-element sequence `Seq(x)`, the “head” of `s`, and the remainder (“tail”) sequence `xs`. So, we split `s` as `s = Seq(x) ++ xs` rather than as `s = xs ++ Seq(x)`.

For computing the sum of a numerical sequence, the order of summation does not matter. However, the order of operations *will* matter for many other computational tasks. We need to choose whether the inductive step should split the sequence as `s = Seq(x) ++ xs` or as `s = xs ++ Seq(x)`, according to the task at hand.

Consider the implementation of `digitsToInt` according to the inductive definition shown in the previous subsection:

```
def digitsToInt(s: Seq[Int]): Int = if (s == Seq()) 0 else {
  val x = s.last // To split s = xs ++ Seq(x), compute x
  val xs = s.take(s.length - 1) // and xs.
  digitsToInt(xs) * 10 + x // Call digitstoInt(...) recursively.
}
```

In this example, it is important to split the sequence `s = xs ++ Seq(x)` in this order, and not in the order `Seq(x) ++ xs`. The reason is that digits increase their numerical value from right to left, so we need to multiply the value of the *left* subsequence, `digitsToInt(xs)`, by 10, in order to compute the correct result.

These examples show how mathematical induction is converted into recursive code. This approach often works but has two technical problems. The first problem is that the code will fail due to the “stack overflow” when the input sequence `s` is long enough. In the next subsection, we will see how this problem is solved

2.2 Converting a sequence into a single value

(at least in some cases) using “tail recursion”. The second problem is that each inductively defined function repeats the code for checking the base case and the code for splitting the sequence `s` into the subsequence `xs` and the extra element `x`. This repeated common code can be put into a library function, and the Scala library provides such functions. We will look at using them in Section [2.2.4](#).

2.2.3 Tail recursion

The code of `lengthS` will fail for large enough sequences. To see why, consider an inductive definition of the `.length` method as a function `lengthS`:

```
def lengthS(s: Seq[Int]): Int =
  if (s == Seq()) 0
  else 1 + lengthS(s.tail)

scala> lengthS((1 to 1000).toList)
res0: Int = 1000

scala> val s = (1 to 100000).toList
s: List[Int] = List(1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17,
  18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, ...)

scala> lengthS(s)
java.lang.StackOverflowError
  at .lengthS(<console>:12)
  at .lengthS(<console>:12)
  at .lengthS(<console>:12)
  at .lengthS(<console>:12)
  ...
```

The problem is not due to insufficient main memory: *we are* able to compute and hold in memory the sequence `s`. The problem is with the code of the function `lengthS`. This function calls itself *inside* an expression `1 + lengthS(...)`. So we can visualize how the computer evaluates this code:

```
lengthS(Seq(1, 2, ..., 100000))
= 1 + lengthS(Seq(2, ..., 100000))
= 1 + (1 + lengthS(Seq(3, ..., 100000)))
= ...
```

The function body of `lengthS` will evaluate the inductive step, that is, the “*else*” part of the “*if/else*”, about 100000 times. Each time, the sub-expression with nested computations `1+(1+(...))` will get larger. This intermediate sub-expression

2 Mathematical formulas as code. II. Mathematical induction

needs to be held somewhere in memory, until at some point the function body goes into the base case and returns a value. When that happens, the entire intermediate sub-expression will contain about 100000 nested function calls still waiting to be evaluated. This sub-expression is held in a special area of memory called **stack memory**, where the not-yet-evaluated nested function calls are held in the order of their calls, as if on a “stack”. Due to the way computer memory is managed, the stack memory has a fixed size and cannot grow automatically. So, when the intermediate expression becomes large enough, it causes an overflow of the stack memory, and the program may crash.

A way to solve this problem is to use a trick called **tail recursion**. Using tail recursion means rewriting the code so that all recursive calls occur at the end positions (at the “tails”) of the function body. In other words, each recursive call must be *itself* the last computation in the function body. Recursive calls cannot be placed inside other computations.

As an example, we can rewrite the code of `lengthS` in this way:

```
def lengthT(s: Seq[Int], res: Int): Int =  
  if (s == Seq()) res  
  else lengthT(s.tail, 1 + res)
```

In this code, one of the branches of the `if/else` returns a fixed value without doing any recursive calls, while the other branch returns the result of recursive call to `lengthT(...)`. In the code of `lengthT`, recursive calls do not occur within sub-expressions such as `1 + lengthT(...)`, unlike in the code of `lengthS`.

It is not a problem that the recursive call to `lengthT` has some sub-expressions such as `1+res` as its arguments, because all these sub-expressions will be computed *before* `lengthT` is recursively called. The recursive call to `lengthT` is the *last* computation performed by this branch of the `if/else`. This shows that the code of `lengthT` is tail-recursive.

A tail-recursive function can have many `if/else` or `match/case` branches, with or without recursive calls; but all recursive calls must be always the last expressions returned.

The Scala compiler has a feature for checking automatically that a function's code is tail-recursive: the `@tailrec` annotation. If a function with a `@tailrec` annotation is not tail-recursive, or is not recursive at all, the program will not compile.

```
@tailrec def lengthT(s: Seq[Int], res: Int): Int =  
  if (s == Seq()) res  
  else lengthT(s.tail, 1 + res)
```


2.2 Converting a sequence into a single value

Let us trace the evaluation of this function on a short example:

```
lengthT(Seq(1,2,3), 0)
= lengthT(Seq(2, 3), 1 + 0) // = lengthT(Seq(2, 3), 1)
= lengthT(Seq(3), 1 + 1) // = lengthT(Seq(3), 2)
= lengthT(Seq(), 1 + 2) // = lengthT(Seq(), 3)
= 3
```

The sub-expressions such as `1 + 2` are computed each time *before* each recursive call to `lengthT`. Because of that, sub-expressions do not grow within the stack memory. This is the main benefit of tail recursion.

How did we rewrite the code of `lengths` to obtain the tail-recursive code of `lengthT`? An important difference between `lengths` and `lengthT` is the additional argument, `res`, called the **accumulator argument**. This argument is equal to an intermediate result of the computation. The next intermediate result (`1 + res`) is computed and passed on to the next recursive call via the accumulator argument. In the base case of the recursion, the function now returns the accumulated result, `res`, rather than `0`, because at that time the computation is finished.

Rewriting code by adding an accumulator argument to achieve tail recursion is called the **accumulator technique** or the “accumulator trick”.

One consequence of using the accumulator trick is that the function `lengthT` now always needs a value for the accumulator argument. However, our goal is to implement a function such as `length(s)` with just one argument, `s: Seq[Int]`. We can define `length(s) = lengthT(s, ???)` if we supply an initial accumulator value. The correct initial value for the accumulator is `0`, since in the base case (an empty sequence `s`) we need to return `0`.

So, a tail-recursive implementation of `lengthT` requires us to define *two* functions: the tail-recursive `lengthT` and an “adapter” function that will set the initial value of the accumulator argument. To emphasize that `lengthT` is a helper function, one could define it *inside* the adapter function:

```
def length[A](s: Seq[A]): Int = {
  @tailrec def lengthT(s: Seq[A], res: Int): Int = {
    if (s == Seq()) res
    else lengthT(s.tail, 1 + res)
  }
  lengthT(s, 0)
}
```

With this code, users will not be able to call `lengthT` directly, since it is only visible within the body of the `length` function.

2 Mathematical formulas as code. II. Mathematical induction

Another possibility in Scala is to use a **default value** for the `res` argument:

```
@tailrec def length[A](s: Seq[A], res: Int = 0): Int =  
  if (s == Seq()) res  
  else length(s.tail, 1 + res)
```

In Scala, giving a default value for a function argument is the same as defining *two* functions: one with that argument and one without. For example, the syntax

```
def f(x: Int, y: Boolean = false): Int = ... // Function body.
```

is equivalent to defining two functions (with the same name),

```
def f(x: Int, y: Boolean) = ... // Function body.  
def f(x: Int): Int = f(Int, false)
```

Using a default argument value, we can define the tail-recursive helper function and the adapter function at once, making the code shorter.

The accumulator trick works in a large number of cases, but it may be far from obvious how to introduce the accumulator argument, what its initial value must be, and how to define the induction step for the accumulator. In the example with the `lengthT` function, the accumulator trick works because of the following mathematical property of the expression being computed:

$$1 + (1 + (1 + (... + 1))) = (((1 + 1) + 1) + ...) + 1 \quad .$$

This property is called the **associativity law** of addition. Due to this law, the computation can be rearranged so that additions associate to the left. In code, it means that intermediate expressions are computed immediately before making the recursive calls; this avoids the growth of the intermediate expressions.

Usually, the accumulator trick works because some associativity law is present. In that case, we are able to rearrange the order of recursive calls so that these calls always occur outside all other sub-expressions, – that is, in tail positions. However, not all computations obey a suitable associativity law. Even if a code rearrangement exists, it may not be immediately obvious how to find it.

As an example, consider a tail-recursive re-implementation of the function `digitsToInt` from the previous subsection, where the recursive call is within a sub-expression `digitsToInt(xs) * 10 + x`. To transform the code into a tail-recursive form, we need to rearrange the main computation,

$$r = d_{n-1} + 10 * (d_{n-2} + 10 * (d_{n-3} + 10 * (...d_0))) \quad ,$$

2.2 Converting a sequence into a single value

so that the operations group to the left. We can do this by rewriting r as

$$r = ((d_0 * 10 + d_1) * 10 + \dots) * 10 + d_{n-1} \quad .$$

It follows that the digit sequence s must be split into the *leftmost* digit and the rest, $s = s.\text{head} ++ s.\text{tail}$. So, a tail-recursive implementation of the above formula is

```
@tailrec def fromDigits(s: Seq[Int], res: Int = 0): Int =
  // 'res' is the accumulator.
  if (s == Seq()) res
  else fromDigits(s.tail, 10 * res + s.head)
```

Despite a certain similarity between this code and the code of `digitsToInt` from the previous subsection, the implementation `fromDigits` cannot be directly derived from the inductive definition of `digitsToInt`. One needs a separate proof that `fromDigits(s, 0)` computes the same result as `digitsToInt(s)`. The proof follows from the following property.

Statement 2.2.3.1 For any $xs: \text{Seq}[\text{Int}]$ and $r: \text{Int}$, we have

```
fromDigits(xs, r) = digitsToInt(xs) + r * math.pow(10, s.length)
```

Proof We prove this by induction. Let us use a short notation for sequences, $[1, 2, 3]$ instead of `Seq(1, 2, 3)`, and temporarily write $d(s)$ instead of `digitsToInt(s)` and $f(s, r)$ instead of `fromDigitsT(s, r)`. Then an inductive definition of $f(s, r)$ is

$$f([], r) = r \quad , \quad f([x] ++ s, r) = f(s, 10 * r + x) \quad . \quad (2.1)$$

Denoting by $|s|$ the length of a sequence s , we reformulate Statement 2.2.3.1 as

$$f(s, r) = d(s) + r * 10^{|s|} \quad , \quad (2.2)$$

We prove Eq. (2.2) by induction. To prove the base case $s = []$, we have $f([], r) = r$ and $d([]) + r * 10^0 = r$ since $d([]) = 0$ and $|s| = 0$. The resulting equality $r = r$ proves the base case.

To prove the inductive step, we assume that Eq. (2.2) holds for a given sequence s ; then we need to prove that

$$f([x] ++ s, r) = d([x] ++ s) + r * 10^{|s|+1} \quad . \quad (2.3)$$

To prove this, we transform the left-hand side and the right-hand side separately,

2 Mathematical formulas as code. II. Mathematical induction

hoping that we will obtain the same expression. The left-hand side of Eq. (2.3):

$$\begin{aligned} & f([x]++s, r) \\ \text{use Eq. (2.1) : } & = f(s, 10 * r + x) \\ \text{use Eq. (2.2) : } & = d(s) + (10 * r + x) * 10^{|s|} . \end{aligned}$$

The right-hand side of Eq. (2.3) contains $d([x]++s)$, which we somehow need to simplify. Assuming that $d(s)$ correctly calculates a number from its digits, we can use the basic property of decimal notation, which is that a digit x in front of n other digits has the value $x * 10^n$. This property can be formulated as

$$d([x]++s) = x * 10^{|s|} + d(s) . \quad (2.4)$$

So, the right-hand side of Eq. (2.3) can be rewritten as

$$\begin{aligned} & d([x]++s) + r * 10^{|s|+1} \\ \text{use Eq. (2.4) : } & = x * 10^{|s|} + d(s) + r * 10^{|s|+1} \\ \text{factor out } 10^{|s|} : & = d(s) + (10 * r + x) * 10^{|s|} . \end{aligned}$$

Now we have transformed both sides of Eq. (2.3) to the same expression.

We have not yet proved that the function d satisfies the property in Eq. (2.4). The proof uses induction and begins by writing the code of d in a short notation,

$$d([]) = 0 , \quad d(s++[y]) = d(s) * 10 + y . \quad (2.5)$$

The base case is Eq. (2.4) with $s = []$. It is proved by

$$x = d([]++[x]) = d([x]++) = x * 10^0 + d([]) = x .$$

The induction step assumes Eq. (2.4) for a given x and a given sequence s , and needs to prove that for any y , the same property holds with $s++[y]$ instead of s :

$$d([x]++s++[y]) = x * 10^{|s|+1} + d(s++[y]) . \quad (2.6)$$

The left-hand side of Eq. (2.6) is transformed into its right-hand side like this:

$$\begin{aligned} & d([x]++s++[y]) \\ \text{use Eq. (2.5) : } & = d([x]++s) * 10 + y \\ \text{use Eq. (2.4) : } & = (x * 10^{|s|} + d(s)) * 10 + y \\ \text{expand parentheses : } & = x * 10^{|s|+1} + d(s) * 10 + y \\ \text{use Eq. (2.5) : } & = x * 10^{|s|+1} + d(s++[y]) . \end{aligned}$$

This demonstrates Eq. (2.6) and so concludes the proof.

2.2.4 Implementing a generic aggregation function (foldLeft)

An **aggregation** converts a sequence of values into a single value. In general, the type of the result value may be different from the type of elements in the sequence. To describe this general situation, we introduce type parameters, A and B , so that the input sequence is of type `Seq[A]` and the aggregated value is of type B . Then an inductive definition of any aggregation function $f: \text{Seq}[A] \Rightarrow B$ looks like this:

- (Base case.) For an empty sequence, $f(\text{Seq}()) = b_0$ where $b_0:B$ is a given value.
- (Induction step.) Assuming that $f(xs) = b$ is already computed, we define $f(xs ++ \text{Seq}(x)) = g(x, b)$ where g is a given function with type signature $g:(A, B) \Rightarrow B$.

Then the code implementing f is written using recursion:

```
def f[A, B](s: Seq[A]): B =
  if (s == Seq()) b0
  else g(s.last, f(s.take(s.length - 1)))
```

We can now refactor this code into a generic utility function, by making b_0 and g into parameters. A possible implementation is

```
def f[A, B](s: Seq[A], b: B, g: (A, B) => B): B =
  if (s == Seq()) b
  else g(s.last, f(s.take(s.length - 1), b, g))
```

However, this implementation is not tail-recursive. Applying f to a sequence of, say, three elements, `Seq(x, y, z)`, will create an intermediate expression $g(z, g(y, g(x, b)))$. This expression will grow with the length of s , which is not acceptable. To rearrange the computation into a tail-recursive form, we need to start the base case at the innermost call $g(x, b)$, then compute $g(y, g(x, b))$ and continue. In other words, we need to traverse the sequence starting from its *leftmost* element x , rather than starting from the right. So, instead of splitting the sequence s into $s.last ++ s.take(s.length - 1)$ as we did in the code of f , we need to split s into $s.head ++ s.tail$. Let us also exchange the order of the arguments of g , in order to be more consistent with the way this code is implemented in the Scala library. The resulting code is now tail-recursive:

```
@tailrec def leftFold[A, B](s: Seq[A], b: B, g: (B, A) => B): B =
```

2 Mathematical formulas as code. II. Mathematical induction

```
if (s == Seq()) b
else leftFold(s.tail, g(b, s.head), g)
```

This function is called a “left fold” because it aggregates (or “folds”) the sequence starting from the leftmost element.

In this way, we have defined a general method of computing any inductively defined aggregation function on a sequence. The function `leftFold` implements the logic of aggregation defined via mathematical induction. Using `leftFold`, we can write concise implementations of methods such as `.sum` or `.max`, and of many other similar aggregation methods. The method `leftFold` already contains all the code necessary to set up the base case and the induction step. The programmer just needs to specify the expressions for the initial value `b` and for the updater function `g`.

As a first example, let us use `leftFold` for implementing the `.sum` method:

```
def sum(s: Seq[Int]): Int = leftFold(s, 0, { (x, y) => x + y })
```

To understand in detail how `leftFold` works, let us trace the evaluation of this function when applied to `Seq(1, 2, 3)`:

```
sum(Seq(1, 2, 3)) = leftFold(Seq(1, 2, 3), 0, g)
// Here, g = { (x, y) => x + y }, so g(x, y) = x + y
= leftFold(Seq(2, 3), g(0, 1), g) // g(0, 1) = 1
= leftFold(Seq(3), g(1, 2), g) // now expand the code of leftFold
= leftFold(Seq(3), g(1, 2), g) // g(1, 2) = 3; expand the code
= leftFold(Seq(), g(3, 3), g) // g(3, 3) = 6; expand the code
= 6
```

The second argument of `leftFold` is the accumulator argument. The initial value of the accumulator is specified when first calling `leftFold`. At each iteration, the new accumulator value is computed by calling the updater function `g`, which uses the previous accumulator value and the value of the next sequence element.

To visualize the process of evaluation, it is convenient to write a table showing the sequence elements and the accumulator values as they are updated:

Current element x	Old accumulator value	New accumulator value
1	0	1
2	1	3
3	3	6

2.2 Converting a sequence into a single value

In general, the type of the accumulator value can be different from the type of the sequence elements. An example is an implementation of `count`:

```
def count[A](s: Seq[A], p: A => Boolean): Int =  
  leftFold(s, 0, { (x, y) => x + (if (p(y)) 1 else 0) })
```

The accumulator is of type `Int`, while the sequence elements can have an arbitrary type, parameterized by `A`. The aggregation function `leftFold` works in the same way for all types of accumulators and all types of sequence elements.

In Scala's standard library, the `.foldLeft` method has a different type signature and, in particular, requires its arguments to be in separate argument groups. For comparison, the implementation of `sum` using our `leftFold` function is

```
def sum(s: Seq[Int]): Int = leftFold(s, 0, { (x, y) => x + y })
```

With the Scala library's `.foldLeft` method, the code is written as

```
def sum(s: Seq[Int]): Int = s.foldLeft(0){ (x, y) => x + y }
```

This syntax makes it more convenient to write nameless functions as arguments, since the updater argument of `.foldLeft` is separated from other arguments by curly braces. We will use the standard `.foldLeft` method from now on; the `leftFold` function was implemented only as an illustration.

The method `.foldLeft` is available in the Scala standard library for all collections, including dictionaries and sets. It is safe to use `.foldLeft`, in the sense that no stack overflows will occur even for very large sequences.

The Scala library contains several other methods similar to `.foldLeft`, such as `.foldRight` and `.reduce`. (However, `.foldRight` is not tail-recursive!)

2.2.5 Solved examples: using `foldLeft`

It is important to gain experience using the `.foldLeft` method.

Example 2.2.5.1 Use `.foldLeft` for implementing the `max` function for integer sequences. Return the special value `Int.MinValue` for empty sequences.

Solution: Write an inductive formulation of the `max` function:

- (Base case.) For an empty sequence, return `Int.MinValue`.
- (Inductive step.) If `max` is already computed on a sequence `xs`, say `max(xs) = b`, the value of `max` on a sequence `xs ++ Seq(x)` is `math.max(b, x)`.

2 Mathematical formulas as code. II. Mathematical induction

Now we can write the code:

```
def max(s: Seq[Int]): Int =  
  s.foldLeft(Int.MinValue){ (b, x) => math.max(b, x) }
```

If we are sure that the function will never be called on empty sequences, we can implement `max` in a simpler way by using the `.reduce` method:

```
def max(s: Seq[Int]): Int = s.reduce { (x, y) => math.max(x, y) }
```

Example 2.2.5.2 Implement the `count` method on sequences of type `Seq[A]`.

Solution: Using the inductive definition of the function `count` as shown in Section 2.2.1, we can write the code as

```
def count[A](s: Seq[A], p: A => Boolean): Int =  
  s.foldLeft(0){ (b, x) => b + (if (p(x)) 1 else 0) }
```

Example 2.2.5.3 Implement the function `digitsToInt` using `.foldLeft`.

Solution: The inductive definition of `digitsToInt` is directly translated into code:

```
def digitsToInt(d: Seq[Int]): Int = d.foldLeft(0){ (n, x) => n * 10 + x }
```

Example 2.2.5.4 For a given non-empty sequence `xs: Seq[Double]`, compute the minimum, the maximum, and the mean as a tuple $(x_{\min}, x_{\max}, x_{\text{mean}})$. The sequence should be traversed only once, i.e. we may call `xs.foldLeft` only once.

Solution: Without the requirement of using a single traversal, we would write

```
(xs.min, xs.max, xs.sum / xs.length)
```

However, this code traverses `xs` at least three times, since each of the aggregations `xs.min`, `xs.max`, and `xs.sum` iterates over `xs`. We need to combine the four inductive definitions of `min`, `max`, `sum`, and `length` into a single inductive definition of some function. What is the type of that function's return value? We need to accumulate intermediate values of *all four* numbers (`min`, `max`, `sum`, and `length`) in a tuple. So the required type of the accumulator is `(Double, Double, Double, Double)`. To avoid repeating a very long type expression, we can define a type alias for it, say, `D4`:

```
scala> type D4 = (Double, Double, Double, Double)  
defined type alias D4
```


2.2 Converting a sequence into a single value

The updater function must update each of the four numbers according to the definitions of their inductive steps:

```
def update(p: D4, x: Double): D4 = p match {  
  case (min, max, sum, length) =>  
    (math.min(x, min), math.max(x, max), x + sum, length + 1)  
}
```

Now we can write the code of the required function:

```
def f(xs: Seq[Double]): (Double, Double, Double) = {  
  val init: D4 = (Double.PositiveInfinity, Double.NegativeInfinity, 0, 0)  
  val (min, max, sum, length) = xs.foldLeft(init)(update)  
  (min, max, sum/length)  
}  
  
scala> f(Seq(1.0, 1.5, 2.0, 2.5, 3.0))  
res0: (Double, Double, Double) = (1.0,3.0,2.0)
```

Example 2.2.5.5* Implement the function `digitsToDouble` using `.foldLeft`. The argument is of type `Seq[Char]`. For example,

```
digitsToDouble(Seq('3', '4', '.', '2', '5')) = 34.25
```

Assume that all characters are either digits or the dot character (so, negative numbers are not supported).

Solution: The evaluation of a `.foldLeft` on a sequence of digits will visit the sequence from left to right. The updating function should work the same as in `digitsToInt` until a dot character is found. After that, we need to change the updating function. So, we need to remember whether a dot character has been seen. The only way for `.foldLeft` to “remember” any data is to hold that data in the accumulator value. We can choose the type of the accumulator according to our needs. So, for this task we can choose the accumulator to be a *tuple* that contains, for instance, the floating-point result constructed so far and a `Boolean` flag showing whether we have already seen the dot character.

To visualize what `digitsToDouble` must do, let us consider how the evaluation of `digitsToDouble(Seq('3', '4', '.', '2', '5'))` should go. We can write a table showing the intermediate result at each iteration. This will hopefully help us figure out what the accumulator and the updater function `g(...)` must be:

2 Mathematical formulas as code. II. Mathematical induction

Current digit c	Previous result n	New result $n' = g(n, c)$
'3'	0.0	3.0
'4'	3.0	34.0
','	34.0	34.0
'2'	34.0	34.2
'5'	34.2	34.25

Until the dot character is found, the updater function multiplies the previous result by 10 and adds the current digit. After the dot character, the updater function must add to the previous result the current digit divided by a factor that represents increasing powers of 10. In other words, the update computation $n' = g(n, c)$ must be defined by these formulas:

- Before the dot character: $g(n, c) = n * 10 + c$.
- After the dot character: $g(n, c) = n + \frac{c}{f}$, where f is 10, 100, 1000, etc., for each subsequent digit.

The updater function g has only two arguments: the current digit and the previous accumulator value. So, the changing factor f must be *part of* the accumulator value, and must be multiplied by 10 at each digit after the dot. If the factor f is not a part of the accumulator value, the function g will not have enough information for computing the next accumulator value correctly. So, the updater computation must be $n' = g(n, c, f)$, not $n' = g(n, c)$.

For this reason, we choose the accumulator type as a tuple (Double, Boolean, Double) where the first number is the result n computed so far, the Boolean flag indicates whether the dot was already seen, and the third number is f , that is, the power of 10 by which the current digit will be divided if the dot was already seen. Initially, the accumulator tuple will be equal to (0.0, false, 10.0). Then the updater function is implemented like this:

```
def update(acc: (Double, Boolean, Double), c: Char): (Double, Boolean, Double) =
  acc match { case (n, flag, factor) =>
    if (c == ',') (n, true, factor) // Set flag to 'true' if dot character.
    else {
      val digit = c - '0'
      if (flag) // This digit is after the dot. Update 'factor'.
        (n + digit/factor, flag, factor * 10)
      else // This digit is before the dot.
```

2.2 Converting a sequence into a single value

```
    (n * 10 + digit, flag, factor)
  }
}
```

Now we can implement `digitsToDouble` as follows,

```
def digitsToDouble(d: Seq[Char]): Double = {
  val initAccumulator = (0.0, false, 10.0)
  val (n, _, _) = d.foldLeft(initAccumulator)(update)
  n
}

scala> digitsToDouble(Seq('3', '4', '.', '2', '5'))
res0: Double = 34.25
```

The result of calling `d.foldLeft` is a tuple `(n, flag, factor)`, in which only the first part, `n`, is needed. In Scala's pattern matching expressions, the underscore symbol is used to denote the pattern variables whose values are not needed in the subsequent code. We could extract the first part using the accessor method `._1`, but the code is more readable if we show all parts of the tuple as `(n, _, _)`.

Example 2.2.5.6 Implement the `.map` method for sequences by using `.foldLeft`. The input sequence should be of type `Seq[A]`, and the output sequence of type `Seq[B]`, where `A` and `B` are type parameters. Here are the required type signature of the function and a sample test:

```
def map[A, B](xs: Seq[A])(f: A => B): Seq[B] = ???

scala> map(List(1, 2, 3)){ x => x * 10 }
res0: Seq[Int] = List(10, 20, 30)
```

Solution: The required code should build a new sequence by applying the function `f` to each element. How can we build a new sequence using `.foldLeft`? The evaluation of `.foldLeft` consists of iterating over the input sequence and accumulating some result value, which is updated at each iteration. Since the result of a `.foldLeft` is always equal to the last computed accumulator value, it follows that the new sequence should *be* the accumulator value. So, we need to update the accumulator by appending the value `f(x)`, where `x` is the current element of the input sequence. We can append elements to sequences using the `:+` operation:

```
def map[A, B](xs: Seq[A])(f: A => B): Seq[B] =
  xs.foldLeft(Seq[B]()) { (acc, x) => acc :+ f(x) }
```

2 Mathematical formulas as code. II. Mathematical induction

The operation `acc := f(x)` is equivalent to `acc ++ Seq(f(x))` but is shorter to write.

Example 2.2.5.7 Implement a function `toPairs` that converts a sequence of type `Seq[A]` to a sequence of pairs, `Seq[(A, A)]`, by putting together each pair of adjacent elements. If the initial sequence has an odd number of elements, a given default value of type `A` is used:

```
def toPairs[A](xs: Seq[A], default: A): Seq[(A, A)] = ???

scala> toPairs(Seq(1, 2, 3, 4, 5, 6), -1)
res0: Seq[(Int, Int)] = List((1,2), (3,4), (5,6))

scala> toPairs(Seq("a", "b", "c"), "<nothing>")
res1: Seq[(String, String)] = List((a,b), (c,<nothing>))
```

Solution: We need to use `.foldLeft` to accumulate a sequence of pairs. However, we iterate over elements of the input sequence one by one. So, a new pair can be added only once every two iterations. The accumulator needs to hold the information about the current iteration being even or odd. For odd-numbered iterations, the accumulator also needs to store the previous element that is still waiting for its pair. Therefore, we choose the type of the accumulator to be a tuple `(Seq[(A, A)], Seq(A))`. The first sequence is the intermediate result, and the second sequence is the “remainder”: it holds the previous element for odd-numbered iterations and is empty for even-numbered iterations. Initially, the accumulator should be empty. A trace of the accumulator updates is shown in this table:

Current element <i>x</i>	Previous accumulator	Next accumulator
"a"	(Seq(), Seq())	(Seq(), Seq("a"))
"b"	(Seq(), Seq("a"))	(Seq(("a", "b")), Seq())
"c"	(Seq(("a", "b")), Seq())	(Seq(("a", "b")), Seq("c"))

Now it becomes clear how to implement the updater function. The code calls `.foldLeft` and then performs some post-processing to make sure we create the last pair in case the last iteration is odd-numbered, i.e. when the “remainder” is not empty after `.foldLeft` is finished. In this implementation, we use pattern matching to decide whether a sequence is empty:

```
def toPairs[A](xs: Seq[A], default: A): Seq[(A, A)] = {
  type Acc = (Seq[(A, A)], Seq[A]) // Type alias, for brevity.
  def init: Acc = (Seq(), Seq())
```

2.2 Converting a sequence into a single value

```
def updater(acc: Acc, x: A): Acc = acc match {
  case (result, Seq()) => (result, Seq(x))
  case (result, Seq(prev)) => (result ++ Seq((prev, x)), Seq())
}
val (result, remainder) = xs.foldLeft(init)(updater)
// May need to append the last element to the result.
remainder match {
  case Seq() => result
  case Seq(x) => result ++ Seq((x, default))
}
}
```

This code shows examples of partial functions that are applied safely. One of these partial functions is used in the expression

```
remainder match {
  case Seq() => ...
  case Seq(a) => ...
}
```

This code works when `remainder` is empty or has length 1, but fails for longer sequences. So it is safe to apply the partial function as long as it is used on sequences of length at most 1, which is indeed the case for the code of `toPairs`.

2.2.6 Exercises: Using `foldLeft`

Exercise 2.2.6.1 Implement a function `fromPairs` that performs the inverse transformation to the `toPairs` function defined in Example 2.2.5.7. The required type signature and a sample test:

```
def fromPairs[A](xs: Seq[(A, A)]): Seq[A] = ???

scala> fromPairs(Seq((1, 2), (3, 4)))
res0: Seq[Int] = List(1, 2, 3, 4)
```

Hint: This can be done with `.foldLeft` or with `.flatMap`.

Exercise 2.2.6.2 Implement the `flatten` method for sequences by using `.foldLeft`. The required type signature and a sample test:

```
def flatten[A](xxs: Seq[Seq[A]]): Seq[A] = ???

scala> flatten(Seq(Seq(1, 2, 3), Seq(), Seq(4)))
res0: Seq[Int] = List(1, 2, 3, 4)
```

2 Mathematical formulas as code. II. Mathematical induction

Exercise 2.2.6.3 Use `.foldLeft` to implement the `zipWithIndex` method for sequences. The required type signature and a sample test:

```
def zipWithIndex[A](xs: Seq[A]): Seq[(A, Int)] = ???

scala> zipWithIndex(Seq("a", "b", "c", "d"))
res0: Seq[String] = List((a, 0), (b, 1), (c, 2), (d, 3))
```

Exercise 2.2.6.4 Use `.foldLeft` to implement a function `filterMap` that combines `.map` and `.filter` for sequences. The required type signature and a sample test:

```
def filterMap[A, B](xs: Seq[A])(pred: A => Boolean)(f: A => B): Seq[B] = ???

scala> filterMap(Seq(1, 2, 3, 4)) { x => x > 2 } { x => x * 10 }
res0: Seq[Int] = List(30, 40)
```

Exercise 2.2.6.5* Split a sequence into subsequences (“batches”) of length not larger than a given maximum length n . The required type signature and a sample test are:

```
def batching[A](xs: Seq[A], size: Int): Seq[Seq[A]] = ???

scala> batching(Seq("a", "b", "c", "d"), 2)
res0: Seq[Seq[String]] = List(List(a, b), List(c, d))

scala> batching(Seq(1, 2, 3, 4, 5, 6, 7), 3)
res1: Seq[Seq[Int]] = List(List(1, 2, 3), List(4, 5, 6), List(7))
```

Exercise 2.2.6.6* Split a sequence into batches by “weight” computed via a given function. The total weight of items in any batch should not be larger than a given maximum weight. The required type signature and a sample test:

```
def byWeight[A](xs: Seq[A], maxW: Double)(w: A => Double): Seq[Seq[A]] = ???

scala> byWeight((1 to 10).toList, 5.75){ x => math.sqrt(x) }
res0: Seq[Seq[Int]] = List(List(1, 2, 3), List(4, 5), List(6, 7), List(8),
    List(9), List(10))
```

Exercise 2.2.6.7* Use `.foldLeft` to implement a `groupBy` function. The required type signature and a sample test:

```
def groupBy[A, K](xs: Seq[A])(by: A => K): Map[K, Seq[A]] = ???
```

2.3 Converting a single value into a sequence

```
scala> groupBy(Seq(1, 2, 3, 4, 5)){ x => x % 2 }
res0: Map[Int, Seq[Int]] = Map(1 -> List(1, 3, 5), 0 -> List(2, 4))
```

Hints:

The accumulator should be of type `Map[K, Seq[A]]`.

To work with dictionaries, you will need to use the methods `.getOrElse` and `.updated`. The method `.getOrElse` fetches a value from a dictionary by key, and returns the given default value if the dictionary does not contain that key:

```
scala> Map("a" -> 1, "b" -> 2).getOrElse("a", 300)
res0: Int = 1

scala> Map("a" -> 1, "b" -> 2).getOrElse("c", 300)
res1: Int = 300
```

The method `.updated` produces a new dictionary that contains a new value for the given key, whether or not that key already exists in the dictionary:

```
scala> Map("a" -> 1, "b" -> 2).updated("c", 300) // Key is new.
res0: Map[String,Int] = Map(a -> 1, b -> 2, c -> 300)

scala> Map("a" -> 1, "b" -> 2).updated("a", 400) // Key already exists.
res1: Map[String,Int] = Map(a -> 400, b -> 2)
```

2.3 Converting a single value into a sequence

An aggregation converts or “folds” a sequence into a single value; the opposite operation (“unfolding”) converts a single value into a sequence. An example of this task is to compute the sequence of decimal digits for a given integer:

```
def digitsOf(x: Int): Seq[Int] = ???

scala> digitsOf(2405)
res0: Seq[Int] = List(2, 4, 0, 5)
```

We cannot implement this function using `.map`, `.zip`, or `.foldLeft`, because these methods may be applied only if we *already have* a sequence. A new sequence can be created, e.g., via the expression `(1 to n)`, but we do not know in advance how long the required sequence must be. The length of the required sequence is determined by a condition that we cannot easily evaluate in advance.

2 Mathematical formulas as code. II. Mathematical induction

A general “unfolding” operation requires us to build a sequence whose length is not determined in advance. This kind of sequence is called a **stream**. A stream is a sequence whose elements are computed only when necessary (unlike sequences such as a `List` or an `Array`, whose elements are all computed in advance and stored). The unfolding operation will keep computing the next element; this creates a stream. We can then apply `.takeWhile` to the stream, in order to stop it when a certain condition holds. Finally, if required, the truncated stream may be converted to a list or to another type of sequence. In this way, we can generate a sequence of initially unknown length according to the given requirements.

A general stream-producing function `Stream.iterate` is available in the Scala library. This function has two arguments, the initial value and a function that computes the next value from the previous one:

```
scala> Stream.iterate(2){ x => x + 10 }  
res0: Stream[Int] = Stream(2, ?)
```

The stream is ready to start computing the next elements of the sequence (so far, only the first element, 2, has been computed). In order to see the next elements, we need to stop the stream at a finite size and then convert the result to a list:

```
scala> Stream.iterate(2){ x => x + 10 }.take(6).toList  
res1: List[Int] = List(2, 12, 22, 32, 42, 52)
```

If we try to evaluate `.toList` on a stream without first limiting its size via `.take` or `.takeWhile`, the program will keep producing more and more elements of the stream, until it runs out of memory and crashes.

Streams are similar to sequences, and methods such as `.map`, `.filter`, `.flatMap` are also defined for streams. For instance, we can use the method `.drop` that skips a given number of initial elements:

```
scala> Seq(10, 20, 30, 40, 50).drop(3)  
res2: Seq[Int] = List(40, 50)  
  
scala> Stream.iterate(2){ x => x + 10 }.drop(3)  
res3: Stream[Int] = Stream(32, ?)
```

This example shows that in order to evaluate `.drop(3)`, the stream had to compute its elements up to 32 (but the subsequent elements are still not computed).

To figure out the code for `digitsOf`, we first write this function as a mathematical formula. To compute the sequence of digits for, say, $n = 2405$, we need to divide n repeatedly by 10, getting a sequence n_k of intermediate numbers ($n_0 = 2405$,

2.3 Converting a single value into a sequence

$n_1 = 240, \dots$) and the corresponding sequence of last digits, $n_k \bmod 10$ (in this example: 5, 0, ...). The sequence n_k is defined using mathematical induction:

- Base case: $n_0 = n$, where n is the given initial integer.
- Inductive step: $n_{k+1} = \left\lfloor \frac{n_k}{10} \right\rfloor$ for $k = 1, 2, \dots$

Here $\left\lfloor \frac{n_k}{10} \right\rfloor$ is the mathematical notation for the integer division by 10.

Let us trace the evaluation of the sequence n_k for $n = 2405$:

$k =$	0	1	2	3	4	5	6
$n_k =$	2405	240	24	2	0	0	0
$n_k \bmod 10 =$	5	0	4	2	0	0	0

The numbers n_k will remain all zeros after $k = 4$. It is clear that the useful part of the sequence is before it becomes all zeros. In this example, the sequence n_k needs to be stopped at $k = 4$. The sequence of digits then becomes [5, 0, 4, 2], and we need to reverse it to obtain [2, 4, 0, 5]. For reversing a sequence, the Scala library has the standard method `.reverse`. So the code is

```
def digitsOf(n: Int): Seq[Int] =  
  if (n == 0) Seq(0) else { // n == 0 is a special case.  
    Stream.iterate(n) { nk => nk / 10 }  
      .takeWhile { nk => nk != 0 }  
      .map { nk => nk % 10 }  
      .toList.reverse  
  }
```

By writing nameless functions such as `{ nk => nk % 10 }` in a shorter syntax such as `(_ % 10)`, we can shorten the code of `digitsOf`:

```
def digitsOf(n: Int): Seq[Int] =  
  if (n == 0) Seq(0) else { // n == 0 is a special case.  
    Stream.iterate(n)(_ / 10)  
      .takeWhile(_ != 0)  
      .map(_ % 10)  
      .toList.reverse  
  }
```

The type signature of the method `Stream.iterate` can be written as

```
def iterate[A](init: A)(next: A => A): Stream[A]
```

This shows how `Stream.iterate` constructs a sequence defined by mathematical induction. The base case is the first value, `init`, and the inductive step is a function, `next`, that computes the next element from the previous one. It is a flexible way of generating sequences whose length is not determined in advance.

2.4 Transforming a sequence into another sequence

We have seen methods such as `.map` and `.zip` that transform sequences into sequences. However, these methods cannot express a general transformation where the elements of the new sequence are defined by induction, depending on previous elements. An example of a task of this kind is to compute the partial sums of a given sequence x_i :

$$b_k = \sum_{i=0}^{k-1} x_i \quad .$$

This formula defines $b_0 = 0$, $b_1 = x_0$, $b_2 = x_0 + x_1$, $b_3 = x_0 + x_1 + x_2$, etc. A definition via mathematical induction may be written like this,

- (Base case.) $b_0 = 0$.
- (Induction step.) Given b_k , we define $b_{k+1} = b_k + x_k$ for $k = 0, 1, 2, \dots$

The Scala library method `.scanLeft` implements a general sequence-to-sequence transformation defined in this way. The code implementing the partial sums is

```
def partialSums(xs: Seq[Int]): Seq[Int] = xs.scanLeft(0){ (x, y) => x + y }  
  
scala> partialSums(Seq(1, 2, 3, 4))  
res0: Seq[Int] = List(0, 1, 3, 6, 10)
```

The first argument of `.scanLeft` is the base case, and the second argument is an updater function describing the induction step. In general, the type of elements of the second sequence is different from that of the first sequence. The updater function takes an element of the first sequence and a previous element of the second sequence, and returns the next element of the second sequence. Note that the result of `.scanLeft` is one element longer than the original sequence, because the base case provides an initial value.

Until now, we have seen that `.foldLeft` is sufficient to re-implement almost every method that work on sequences, such as `.map`, `.filter`, or `.flatten`. The method `.scanLeft` can be also implemented via `.foldLeft`. In the implementation, the accumulator contains the previous element of the second sequence together with a growing fragment of that sequence, which is updated as we iterate over the first sequence. The code (shown here only as illustration) is

```
def scanLeft[A, B](xs: Seq[A])(b0: B)(next: (B, A) => B): Seq[B] = {
  val init: (B, Seq[B]) = (b0, Seq(b0))
  val (_, result) = xs.foldLeft(init){ case ((b, seq), x) =>
    val newB = next(b, x)
    (newB, seq ++ Seq(newB))
  }
  result
}
```

To define the (nameless) updater function for `.foldLeft`, we used the Scala feature that makes it easier to define functions with several arguments containing tuples. In our case, the updater function in `.foldLeft` has two arguments: the first is a tuple `(B, Seq[B])`, the second is a value of type `A`. The pattern matching expression `{ case ((b, seq), x) => ... }` appears to match against a nested tuple. In reality, this expression matches the two arguments of the updater function and, at the same time, deconstructs the tuple argument as `(b, seq)`.

2.5 Summary

We have done a broad overview of translating mathematical induction into Scala code. What problems can we solve now?

- Compute mathematical expressions involving arbitrary recursion.
- Use the accumulator trick to enforce tail recursion.
- Use arbitrary inductive (i.e. recursive) formulas to:
 - convert sequences to single values (“aggregation”);
 - create new sequences from single values;
 - transform existing sequences into new sequences.

This table summarizes the Scala methods implementing these tasks:

2 Mathematical formulas as code. II. Mathematical induction

Definition via mathematical induction	Scala code example
$f([]) = b; f(s++[x]) = g(f(s), x)$	<code>f(xs) = xs.foldLeft(b)(g)</code>
$x_0 = b; x_{k+1} = g(x_k)$	<code>xs = Stream.iterate(b)(g)</code>
$y_0 = b; y_{k+1} = g(y_k, x_k)$	<code>ys = xs.scanLeft(b)(g)</code>

Using these methods, any iterative calculation is implemented by translating mathematical induction directly into code. In the functional programming paradigm, the programmer does not need to write any loops or check any array indices. Instead, the programmer reasons about sequences as mathematical values: “Starting from this value, we get that sequence, then transform it into this other sequence,” etc. This is a powerful way of working with sequences, dictionaries, and sets. Many kinds of programming errors (such as an incorrect array index) are avoided from the outset, and the code is shorter and easier to read than conventional code written using loops.

What tasks are not possible with these tools? We cannot implement a non-tail-recursive function without stack overflow (i.e. without unlimited growth of intermediate expressions). The accumulator trick does not always work! In some cases, it is impossible to implement tail recursion in a given recursive computation. An example of such a computation is the “merge-sort” algorithm where the function body must contain two recursive calls within a single expression. (It is impossible to rewrite *two* recursive calls as one!)

What if our recursive code cannot be transformed into tail-recursive code via the accumulator trick, but the depth of the recursion is so large that stack overflows are possible? We must then use more advanced tricks (for instance, the “continuation-passing” or “trampolines”) that convert non-tail-recursive code into iterative code without stack overflows. These tricks are beyond the scope of this chapter.

2.5.1 Solved examples

Example 2.5.1.1 Compute the smallest n such that $f(f(f(\dots f(1)\dots)) \geq 1000$, where the function f is applied n times. Write this as a function taking f , 1, and 1000 as arguments. Test with $f(x) = 2x + 1$.

Solution: We define a stream of values $[1, f(1), f(f(1)), \dots]$ and use `.takeWhile` to stop the stream when the given condition holds. The `.length` method then gives the length of the resulting sequence:

```
scala> Stream.iterate(1)(x => 2*x+1).takeWhile(x => x < 1000).toList
res0: List[Int] = List(1, 3, 7, 15, 31, 63, 127, 255, 511)

scala> Stream.iterate(1)(x => 2*x+1).takeWhile(x => x < 1000).length
res1: Int = 9
```

Example 2.5.1.2 (a) For a given stream of integers, compute the stream of the largest values seen so far.

(b) Compute the stream of k largest values seen so far (k is a given integer parameter).

Solution: We cannot use `.max` or sort the entire stream, since the length of the stream is not known in advance. So we need to use `.scanLeft`, which will build the output stream one element at a time.

(a) Maintain the largest value seen so far in the accumulator of the `.scanLeft`:

```
def maxSoFar(xs: Stream[Int]): Stream[Int] =
  xs.scanLeft(xs.head){ case (max, x) => math.max(max, x) }
    .drop(1)
```

We use `.drop(1)` to remove the initial value, `xs.head`, because it is not useful for our result but is necessary for the definition of `.scanLeft`.

To test this function, let us define a stream whose values go up and down:

```
scala> val s = Stream.iterate(0)(x => 1 - 2*x)
s: Stream[Int] = Stream(0, ?)

scala> s.take(10).toList
res0: List[Int] = List(0, 1, -1, 3, -5, 11, -21, 43, -85, 171)

scala> maxSoFar(s).take(10).toList
res1: List[Int] = List(0, 1, 1, 3, 3, 11, 11, 43, 43, 171)
```

(b) We again use `.scanLeft`, where now the accumulator needs to keep the largest k values seen so far. There are two ways of maintaining this accumulator: First, to have a sequence of k values that we sort and truncate each time. Second, to use a specialized data structure such as a priority queue that automatically keeps values sorted and its length bounded. For the purposes of this tutorial, let us avoid using specialized data structures:

```
def maxKSoFar(xs: Stream[Int], k: Int): Stream[Seq[Int]] = {
  // The initial value of the accumulator is an empty Seq() of type Seq[Int].
  xs.scanLeft(Seq[Int]()) { case (seq, x) =>
```

2 Mathematical formulas as code. II. Mathematical induction

```
// Sort in the descending order, and take the first k values.
(seq :+ x).sorted.reverse.take(k)
}.drop(1) // Skip the useless first value.
}

scala> maxKSoFar(s, 3).take(10).toList
res2: List[Seq[Int]] = List(List(0), List(1, 0), List(1, 0, -1), List(3, 1, 0),
    List(3, 1, 0), List(11, 3, 1), List(11, 3, 1), List(43, 11, 3), List(43,
    11, 3), List(171, 43, 11))
```

Example 2.5.1.3 Find the last element of a non-empty sequence. (Use `.reduce`.)

Solution: This function is available in the Scala library as the standard method `.last` on sequences. Here we need to re-implement it using `.reduce`. Begin by writing an inductive definition:

- (Base case.) $\text{last}(\text{Seq}(x)) = x$.
- (Inductive step.) $\text{last}(\text{Seq}(x) ++ xs) = \text{last}(xs)$ assuming xs is non-empty.

The `.reduce` method implements an inductive aggregation similarly to `.foldLeft`, except that for `.reduce` the base case is fixed – it always returns x for a 1-element sequence $\text{Seq}(x)$. This is exactly what we need here, so the inductive definition is directly translated into code, with the updater function $g(x, y) = y$:

```
def last[A](xs: Seq[A]): A = xs.reduce { case (x, y) => y }
```

Example 2.5.1.4 (a) Using tail recursion, implement the binary search algorithm in a given sorted sequence `xs: Seq[Int]` as a function returning the index of the requested number `n` (assume that `xs` contains the number `n`):

```
def binSearch(xs: Seq[Int], goal: Int): Int = ???

scala> binSearch(Seq(1, 3, 5, 7), 5)
res0: Int = 2
```

(b) Re-implement `binSearch` using `Stream.iterate` instead of explicit recursion.

Solution: (a) The well-known binary search algorithm splits the array into two halves and may continue the search recursively in one of the halves. We need to write the solution as a tail-recursive function with an additional accumulator argument. So we expect that the code should look like this,

```
def binSearch(xs: Seq[Int], goal: Int, acc: _ = ???): Int = {
  if (???are we done???) acc
```

```

else {
  // Determine which half of the sequence contains 'goal'.
  // Then update the accumulator accordingly.
  val newAcc = ???
  binSearch(xs, goal, newAcc) // Tail-recursive call.
}
}

```

It remains to determine the type and the initial value of the accumulator, as well as the code for updating it.

The information required for the recursive call is the remaining segment of the sequence where the target number is present. This segment is defined by two indices i, j representing the left and the right bounds of the sub-sequence, such that the target element is x_n with $x_i \leq x_n < x_{j-1}$. It follows that the accumulator should be a pair of two integers (i, j) . The initial value of the accumulator is the pair $(0, N)$ where N is the length of the entire sequence. The search is finished when $i + 1 = j$. We can now write the corresponding code, where for convenience we introduce *two* accumulator values:

```

@tailrec def binSearch(xs: Seq[Int], goal: Int)(left: Int = 0,
                                             right: Int = xs.length): Int = {
  // Check whether 'goal' is at one of the boundaries.
  if (right - left <= 1 || xs(left) == goal) left
  else {
    val middle = (left + right) / 2
    // Determine which half of the array contains 'target'.
    // Update the accumulator accordingly.
    val (newLeft, newRight) =
      if (goal < xs(middle)) (left, middle)
      else (middle, right)
    binSearch(xs, goal)(newLeft, newRight) // Tail-recursive call.
  }
}

scala> binSearch(0 to 10, 3)() // Default accumulator values.
res0: Int = 3

```

Here we used a feature of Scala that allows us to use `xs.length` as a default value for the argument `right` of `binSearch`. This is possible only because `right` is in a different **argument group** from `xs`. In Scala, values in an argument group may depend on arguments given in a *previous* argument group. However, the code

```
def binSearch(xs: Seq[Int], goal: Int, left: Int = 0, right: Int = xs.length)
```

2 Mathematical formulas as code. II. Mathematical induction

will generate an error: the arguments in the same argument group cannot depend on each other. (The error will say `not found: value xx.`)

(b) We can visualize the binary search as a procedure that generates a sequence of progressively tighter bounds for the location of `goal`. The initial bounds are $(0, \text{xs.length})$, and the final bounds are $(k, k+1)$ for some k . We can generate the sequence of bounds using `Stream.iterate` and stop the sequence when the bounds become sufficiently tight. To make the use of `.takeWhile` more convenient, we add an extra sequence element where the bounds (k, k) are equal. The code becomes

```
def binSearch(xs: Seq[Int], goal: Int): Int = {
  type Acc = (Int, Int)
  val init: Acc = (0, xs.length)
  val updater: Acc => Acc = { case (left, right) =>
    if (right - left <= 1) (left, left) // Extra element.
    else if (xs(left) == goal) (left, left + 1)
    else {
      val middle = (left + right) / 2
      // Determine which half of the array contains 'target'.
      // Update the accumulator accordingly.
      if (goal < xs(middle)) (left, middle)
      else (middle, right)
    }
  }
  Stream.iterate(init)(updater)
    .takeWhile{ case (left, right) => right > left }
    .last._1 // Take the 'left' boundary from the last element.
}
```

This code is clearer because recursion is delegated to `Stream.iterate`, and we only need to write the “business logic” (i.e. the base case and the inductive step) of our function.

Example 2.5.1.5 For a given positive $n:\text{Int}$, compute the sequence $[s_0, s_1, s_2, \dots]$ defined by $s_0 = SD(n)$ and $s_k = SD(s_{k-1})$ for $k > 0$, where $SD(x)$ is the sum of the decimal digits of the integer x , e.g. $SD(123) = 6$. Stop the sequence s_i when the numbers begin repeating. For example, $SD(99) = 18$, $SD(18) = 9$, $SD(9) = 9$. So, for $n = 99$, the sequence s_i must be computed as $[99, 18, 9]$.

Hint: use `Stream.iterate`; compute the decimal digits in the reverse order since the sum will be the same.

Solution: We need to implement a function `sdSeq` having the type signature

```
def sdSeq(n: Int): Seq[Int]
```


First we need to implement $SD(x)$. The sum of digits is obtained by almost the same code as in Section 2.3:

```
def SD(n: Int): Int = if (n==0) 0 else
  Stream.iterate(n)(_ / 10).takeWhile(_ != 0).map(_ % 10).sum
```

Now we can try evaluating SD on some numbers to see its behavior:

```
scala> (1 to 15).toList.map(SD)
res0: List[Int] = List(1, 2, 3, 4, 5, 6, 7, 8, 9, 1, 2, 3, 4, 5, 6)
```

It is clear that $SD(n) < n$ as long as $n \geq 10$. So the sequence elements s_i will not repeat until they become smaller than 10, and then they will always repeat. This seems to be an easy way of stopping the sequence. Let us try that:

```
scala> Stream.iterate(99)(SD).takeWhile(x => x >= 10).toList
res1: List[Int] = List(99, 18)
```

We are missing the last element of the sequence, $SD(18) = 9$, because `.takeWhile` stops the sequence too early. In order to obtain the correct sequence, we need to compute one more element. To fix this, we can generate a stream of *pairs*:

```
scala> Stream.iterate((0, 99)){ case (prev, x) => (x, SD(x)) }.
  takeWhile{ case (prev, x) => prev >= 10 || x >= 10 }.toList
res2: List[(Int, Int)] = List((0,99), (99,18), (18,9))
```

This looks right; it remains to remove the first parts of the tuples:

```
def sdSeq(n: Int): Seq[Int] =
  Stream.iterate((0, n)){ case (prev, x) => (x, SD(x)) } // Stream[(Int, Int)]
    .takeWhile{ case (prev, x) => prev >= 10 || x >= 10 } // Stream[(Int, Int)]
    .map(_._2) // Stream[Int]
    .toList // List[Int]

scala> sdSeq(99)
res3: Seq[Int] = List(99, 18, 9)
```

Example 2.5.1.6 For a given stream $[s_0, s_1, s_2, \dots]$ of type `Stream[T]`, compute the “half-speed” stream $h = [s_0, s_0, s_1, s_1, s_2, s_2, \dots]$. (The half-speed sequence h can be defined by the formula $s_k = h_{2k} = h_{2k+1}$.)

Solution: We use `.map` to replace each element s_i by a sequence containing two copies of s_i . Let us try this on a sample sequence:

```
scala> Seq(1,2,3).map( x => Seq(x, x))
```

2 Mathematical formulas as code. II. Mathematical induction

```
res0: Seq[Seq[Int]] = List(List(1, 1), List(2, 2), List(3, 3))
```

The result is almost what we need, except we need to `.flatten` the nested list:

```
scala> Seq(1,2,3).map( x => Seq(x, x)).flatten
res1: Seq[Seq[Int]] = List(1, 1, 2, 2, 3, 3)
```

The composition of `.map` and `.flatten` is `.flatMap`, so the final code is

```
def halfSpeed[T](str: Stream[T]): Stream[T] = str.flatMap(x => Seq(x, x))

scala> halfSpeed(Seq(1,2,3).toStream)
res2: Stream[Int] = Stream(1, ?)

scala> halfSpeed(Seq(1,2,3).toStream).toList
res3: List[Int] = List(1, 1, 2, 2, 3, 3)
```

Example 2.5.1.7 Stop a given stream $[s_0, s_1, s_2, \dots]$ at a place k where the sequence repeats itself; that is, an element s_k equals some earlier element s_i with $i < k$.

Solution: The trick is to create a half-speed sequence h_i out of s_i and then find an index $k > 0$ such that $h_k = s_k$. (The condition $k > 0$ is needed because we will always have $h_0 = s_0$.) If we find such an index k , it would mean that either $s_k = s_{k/2}$ or $s_k = s_{(k-1)/2}$; in either case, we will have found an element s_k that equals an earlier element.

As an example, take $s = [1, 3, 5, 7, 9, 3, 5, 7, 9, \dots]$ and compute the half-speed sequence $h = [1, 1, 3, 3, 5, 5, 7, 7, 9, 9, 3, 3, \dots]$. Looking for an index $k > 0$ such that $h_k = s_k$, we find that $s_7 = h_7 = 7$. This is indeed an element of s_i that repeats an earlier element (although s_7 is not the first such repetition).

There are in principle two ways of finding an index $k > 0$ such that $h_k = s_k$: First, to iterate over a list of indices $k = 1, 2, \dots$ and evaluate the condition $h_k = s_k$ as a function of k . Second, to build a sequence of pairs (h_i, s_i) and use `.takeWhile` to stop at the required index. In the present case, we cannot use the first way because we do not have a fixed set of indices to iterate over. Also, the condition $h_k = s_k$ cannot be directly evaluated as a function of k because s and h are streams that compute elements on demand, not lists whose elements are computed in advance and ready to be used.

So the code must iterate over a stream of pairs (h_i, s_i) :

```
def stopRepeats[T](str: Stream[T]): Stream[T] = {
  val halfSpeed = str.flatMap(x => Seq(x, x))
  val result = halfSpeed.zip(str) // Stream[(T, T)]
```

```

    .drop(1) // Enforce the condition k > 0.
    .takeWhile { case (h, s) => h != s } // Stream[(T, T)]
    .map(_._2) // Stream[T]
    str.head += result // Prepend the first element that was dropped.
  }

scala> stopRepeats(Seq(1, 3, 5, 7, 9, 3, 5, 7, 9).toStream).toList
res0: List[Int] = List(1, 3, 5, 7, 9, 3, 5)

```

Example 2.5.1.8 Reverse each word in a string, but keep the order of words:

```

def revWords(s: String): String = ???

scala> revWords("A quick brown fox")
res0: String = A kciuq nworb xof

```

Solution: The standard method `.split` converts a string into an array of words:

```

scala> "pa re ci vo mu".split(" ")
res0: Array[String] = Array(pa, re, ci, vo, mu)

```

Each word is reversed with `.reverse`; the resulting array is concatenated into a string with `.mkString`. So the code is

```

def revWords(s: String): String = s.split(" ").map(_.reverse).mkString(" ")

```

Example 2.5.1.9 Remove adjacent repeated characters from a string:

```

def noDups(s: String): String = ???

scala> noDups("abbcddeeeefddgggggh")
res0: String = abcdefdgh

```

Solution: A string is automatically converted into a sequence of characters when we use methods such as `.map` or `.zip` on it. So, we can use `s.zip(s.tail)` to get a sequence of pairs (s_k, s_{k+1}) where s_k is the k -th character of the string s . Now we can use `.filter` to remove the elements s_k for which $s_{k+1} = s_k$:

```

scala> val s = "abbcd"
s: String = abbcd

scala> s.zip(s.tail).filter { case (sk, skPlus1) => sk != skPlus1 }
res0: IndexedSeq[(Char, Char)] = Vector((a,b), (b,c), (c,d))

```

2 Mathematical formulas as code. II. Mathematical induction

It remains to convert this sequence of pairs into the string "abcd". One way of doing this is to project the sequence of pairs onto the second parts of the pairs,

```
scala> res0.map(_._2).mkString
res1: String = bcd
```

We just need to add the first character, 'a'. The resulting code is

```
def noDups(s: String): String = if (s == "") "" else {
  val pairs = s.zip(s.tail).filter { case (x, y) => x != y }
  pairs.head._1 +: pairs.map(_._2).mkString
}
```

The method `+:` prepends an element to a sequence, so `x +: xs` is equivalent to `Seq(x) ++ xs`.

Example 2.5.1.10 (a) Count the occurrences of each distinct word in a string:

```
def countWords(s: String): Map[String, Int] = ???

scala> countWords("a quick a quick a fox")
res0: Map[String, Int] = Map("a" -> 3, "quick" -> 2, "fox" -> 1)
```

(b) Count the occurrences of each distinct element in a sequence of type `Seq[A]`.

Solution: **(a)** We split the string into an array of words via `s.split(" ")`, and apply a `.foldLeft` to that array, since the computation is a kind of aggregation over the array of words. The accumulator of the aggregation will be the dictionary of word counts for all the words seen so far:

```
def countWords(s: String): Map[String, Int] = {
  val init: Map[String, Int] = Map()
  s.split(" ").foldLeft(init) { (dict, word) =>
    val newCount = dict.getOrElse(word, 0) + 1
    dict.updated(word, newCount)
  }
}
```

(b) The main code of `countWords` does not depend on the fact that words are of type `String`. It will work in the same way for any other type of keys for the dictionary. So we keep the same code and define the type signature of the function to contain a type parameter `A` instead of `String`:

```
def countValues[A](xs: Seq[A]): Map[A, Int] =
  xs.foldLeft(Map[A, Int]()) { (dict, word) =>
    val newCount = dict.getOrElse(word, 0) + 1
```

```

    dict.updated(word, newCount)
  }

scala> countValues(Seq(100, 100, 200, 100, 200, 200, 100))
res0: Map[Int,Int] = Map(100 -> 4, 200 -> 3)

```

Example 2.5.1.11 For a given sequence of type `Seq[A]`, find the longest subsequence that does not contain any adjacent duplicate values.

```

def longestNoDups[A](xs: Seq[A]): Seq[A] = ???

scala> longestNoDups(Seq(1, 2, 2, 5, 4, 4, 4, 8, 2, 3, 3))
res0: Seq[Int] = List(4, 8, 2, 3)

```

Solution: This is a dynamic programming problem. Many such problems are solved with a single `.foldLeft`. The accumulator represents the current “state” of the dynamic programming solution, and the “state” is updated with each new element of the input sequence.

To obtain the solution, we first need to determine the type of the accumulator value, or the “state”. The task is to find the longest subsequence without adjacent duplicates. So the accumulator should represent the longest subsequence found so far, as well as any required extra information about other subsequences that might grow as we iterate over the elements of `xs`. What is this extra information in our case?

Imagine that we wanted to build set of *all* subsequences without adjacent duplicates. In the example where the input sequence is `[1, 2, 2, 5, 4, 4, 4, 8, 2, 3, 3]`, this set of subsequences should be `{[1, 2], [2, 5, 4], [4, 8, 2, 3]}`. We can build this set incrementally in the accumulator value of a `.foldLeft`. To visualize how this set would be built, consider the partial result after seeing the first 8 elements of the input sequence, `[1, 2, 2, 5, 4, 4, 4, 8]`. The partial set of non-repeating subsequences is `{[1, 2], [2, 5, 4], [4, 8]}`. As we add another element, 2, we update the partial set to `{[1, 2], [2, 5, 4], [4, 8, 2]}`.

It is now clear that the subsequence `[1, 2]` has no chance of being the longest subsequence, since `[2, 5, 4]` is already longer. However, we do not yet know whether `[2, 5, 4]` or `[4, 8, 2]` is the winner, because the subsequence `[4, 8, 2]` could still grow and become the longest one (and it does become `[4, 8, 2, 3]` later). At this point, we need to keep both of these two subsequences in the accumulator, but we may already discard `[1, 2]`.

We have deduced that the accumulator needs to keep only *two* sequences: the first sequence is already terminated and will not grow, the second sequence ends

2 Mathematical formulas as code. II. Mathematical induction

with the current element and may yet grow. The initial value of the accumulator is empty. The first subsequence is discarded when it becomes shorter than the second. The code can be written now:

```
def longestNoDups[A](xs: Seq[A]): Seq[A] = {
  val init: (Seq[A], Seq[A]) = (Seq(), Seq())
  val (first, last) = xs.foldLeft(init) { case ((first, current), x) =>
    // If 'current' is empty, 'x' cannot be repeated.
    val xWasRepeated = current != Seq() && current.last == x
    val firstIsLongerThanCurrent = first.length > current.length
    // Compute the new pair '(first, current)'.
    // Keep 'first' only if it is longer; otherwise replace it by 'current'.
    val newFirst = if (firstIsLongerThanCurrent) first else current
    // Append 'x' to 'current' if 'x' is not repeated.
    val newCurrent = if (xWasRepeated) Seq(x) else current :+ x
    (newFirst, newCurrent)
  }
  // Return the longer of the two subsequences; prefer 'first'.
  if (first.length >= last.length) first else last
}
```

2.5.2 Exercises

Exercise 2.5.2.1 Compute the sum of squared digits of a given integer; e.g., `dsq(123) = 14` (see Example 2.5.1.5). Generalize the solution to take an arbitrary function `f : Int => Int` as a parameter, instead of the squaring operation. The type signature and a sample test:

```
def digitsMapSum(x: Int)(f: Int => Int): Int = ???

scala> digitsMapSum(123){ x => x * x }
res0: Int = 14

scala> digitsMapSum(123){ x => x * x * x }
res1: Int = 36
```

Exercise 2.5.2.2 Compute the **Collatz sequence** c_i as a stream defined by

$$c_0 = n \quad ; \quad c_{k+1} = \begin{cases} \frac{c_k}{2} & \text{if } c_k \text{ is even,} \\ 3c_k + 1 & \text{if } c_k \text{ is odd.} \end{cases}$$

Stop the stream when it reaches 1 (as one would expect it will).

Exercise 2.5.2.3 For a given integer n , compute the sum of cubed digits, then the sum of cubed digits of the result, etc.; stop the resulting sequence when it repeats itself, and so determine whether it ever reaches 1. (Use Exercise 2.5.2.1.)

```
def cubes(n: Int): Stream[Int] = ???

scala> cubes(123).take(10).toList
res0: List[Int] = List(123, 36, 243, 99, 1458, 702, 351, 153, 153, 153)

scala> cubes(2).take(10).toList
res1: List[Int] = List(2, 8, 512, 134, 92, 737, 713, 371, 371, 371)

scala> cubes(4).take(10).toList
res2: List[Int] = List(4, 64, 280, 520, 133, 55, 250, 133, 55, 250)

def cubesReach1(n: Int): Boolean = ???

scala> cubesReach1(10)
res3: Boolean = true

scala> cubesReach1(4)
res4: Boolean = false
```

Exercise 2.5.2.4 For a, b, c of type `Set[Int]`, compute the set of all sets of the form `Set(x, y, z)` where x is from a , y from b , and z from c . The required type signature and a sample test:

```
def prod3(a: Set[Int], b: Set[Int], c: Set[Int]): Set[Set[Int]] = ???

scala> prod3(Set(1,2), Set(3), Set(4,5))
res0: Set[Set[Int]] = Set(Set(1,3,4), Set(1,3,5), Set(2,3,4), Set(2,3,5))
```

Hint: use `.flatMap`.

Exercise 2.5.2.5* Same task as in Exercise 2.5.2.4 for a set of sets, i.e. given a `Set[Set[Int]]` instead of just three sets a, b, c . The required type signature and a sample test:

```
def prodSet(si: Set[Set[Int]]): Set[Set[Int]] = ???

scala> prodSet(Set(Set(1,2), Set(3), Set(4,5), Set(6)))
res0: Set[Set[Int]] = Set(Set(1,3,4,6), Set(1,3,5,6), Set(2,3,4,6), Set(2,3,5,6))
```

Hint: use `.foldLeft` and `.flatMap`.

2 Mathematical formulas as code. II. Mathematical induction

Exercise 2.5.2.6* In a sorted array `xs:Array[Int]` where no values are repeated, find all pairs of values whose sum equals a given number n . Use tail recursion. A possible type signature and a sample test:

```
def pairs(goal: Int, xs: Array[Int]): Set[(Int, Int)] = ???

scala> pairs(10, Array(1, 2, 3, 4, 5, 6, 7, 8))()
res0: Set[(Int, Int)] = Set((2,8), (3,7), (4,6), (5,5))
```

Exercise 2.5.2.7 Reverse a sentence's word order, but keep the words unchanged:

```
def revSentence(s: String): String = ???

scala> revSentence("A quick brown fox")
res0: String = "fox brown quick A"
```

Exercise 2.5.2.8 Reverse an integer's digits (see Example 2.5.1.5) as shown:

```
def revDigits(n: Int): Int = ???

scala> revDigits(12345)
res0: Int = 54321
```

A **palindrome number** is an integer n such that `revDigits(n) == n`. Write a function `Int => Boolean` that checks whether a given positive integer is a palindrome.

Exercise 2.5.2.9 Starting from a given integer n , compute `revDigits(n) + n`; the function `revDigits` was defined in Exercise 2.5.2.8. Check whether the result is a palindrome integer. If it is not, repeat the same operation until a palindrome number is found, and return that number. The required type signature and a test:

```
def findPalindrome(n: Int): Int = ???

scala> findPalindrome(123)
res0: Int = 444

scala> findPalindrome(83951)
res1: Int = 869363968
```

Exercise 2.5.2.10 (a) For a given integer interval $[n_1, n_2]$, find the largest integer $k \in [n_1, n_2]$ such that the decimal representation of k does *not* contain any of the digits 3, 5, or 7. (b) For a given integer interval $[n_1, n_2]$, find the integer $k \in [n_1, n_2]$ with the largest sum of decimal digits. (c) A positive integer n is called a **perfect**

number if it is equal to the sum of its divisors (other integers k such that $k < n$ and n/k is an integer). For example, 6 is a perfect number because its divisors are 1, 2, and 3, and $1 + 2 + 3 = 6$, while 8 is not a perfect number because its divisors are 1, 2, and 4, and $1 + 2 + 4 = 7 \neq 8$. Write a function that determines whether a given number n is perfect. Determine all perfect numbers up to one million.

Exercise 2.5.2.11 Remove adjacent repeated elements from a sequence of type `Seq[A]` when they are repeated more than k times. Repetitions up to k times should remain unchanged. The required type signature and a sample test:

```
def removeDups[A](s: Seq[A], k: Int): Seq[A] = ???

scala> removeDups(Seq(1, 1, 1, 1, 5, 2, 2, 5, 5, 5, 5, 5, 1), 3)
res0: Seq[Int] = List(1, 1, 1, 5, 2, 2, 5, 5, 5, 1)
```

Exercise 2.5.2.12 (a) Remove repeated elements (whether adjacent or not) from a sequence of type `Seq[A]`. (This re-implements the standard method `.distinct`.)

(b) For a sequence of type `Seq[A]`, remove all elements that are repeated (whether adjacent or not) more than k times:

```
def removeK[A](k: Int, xs: Seq[A]): Seq[A] = ???

scala> removeK(2, Seq("a", "b", "a", "b", "b", "c", "b", "a"))
res0: Seq[String] = List(a, b, a, b, c)
```

Exercise 2.5.2.13* For a given sequence `xs:Seq[Double]`, find a subsequence that has the largest sum of values. The sequence `xs` is not sorted, and its values may be positive or negative. The required type signature and a sample test:

```
def maxsub(xs: Seq[Double]): Seq[Double] = ???

scala> maxsub(Seq(1.0, -1.5, 2.0, 3.0, -0.5, 2.0, 1.0, -10.0, 2.0))
res0: Seq[Double] = List(2.0, 3.0, -0.5, 2.0, 1.0)
```

Hint: use dynamic programming and `.foldLeft`.

Exercise 2.5.2.14* Find all common integers between two sorted sequences:

```
def commonInt(xs: Seq[Int], ys: Seq[Int]): Seq[Int] = ??? // Use tail recursion.

scala> commonInt(Seq(1, 3, 5, 7), Seq(2, 3, 4, 6, 7, 8))
res0: Seq[Int] = List(3, 7)
```

2.6 Discussion

2.6.1 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 simple example of a partial function in Scala is the `.max` method: it only works for non-empty sequences. Trying to evaluate it on an empty sequence generates an error called an “exception”:

```
scala> Seq(1).tail
res0: Seq[Int] = List()
scala> res0.max
java.lang.UnsupportedOperationException: empty.max
    at scala.collection.TraversableOnce$class.max(TraversableOnce.scala:229)
    at scala.collection.AbstractTraversable.max(Traversable.scala:104)
    ... 32 elided
```

This kind of error may crash the entire program at run time. Unlike the type errors we saw before, which occur at compilation time (i.e. before the program can start), **run-time errors** occur while the program is running, and only when some partial function happens to get an incorrect input. The incorrect input may occur at any point after the program started running, which may crash the entire program in the middle of a long computation.

So, it seems clear that we should write code that does not generate such errors. For instance, it is safe to apply `.max` to a sequence if we know that it is non-empty.

Sometimes, a function that uses pattern matching turns out to be a partial function because its pattern matching code fails on certain input data.

If a pattern matching expression 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. In most cases, programs can be written to avoid the possibility of match errors. An example of an unsafe pattern matching expression 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.Tuple2$mcII$sp)
  at .h(<console>:12)
  ... 32 elided
```

Here 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. So, `h` is a partial function.

Pattern matching failures never happen if we match a tuple of correct size with a pattern such as `(x, y, z)`, because pattern variables will always match whatever values the tuple has. So, pattern matching with a pattern such as `(x, y, z)` is **infallible** (never fails at run time) when applied to a tuple with 3 elements.

Another way in which pattern matching can be made infallible is by including a pattern that matches everything:

```
p match {
  case (x, 0) => ... // This only matches some tuples.
  case _ => ... // This matches everything.
}
```

If the first pattern `(x, 0)` fails to match the value of `p`, the second pattern will be tried (and will always succeed). When a `match` expression has several `case` patterns, the patterns are tried in the order they are written. So, a match expression can be made infallible by adding a “match-all” underscore pattern.

2.6.2 Scope of pattern matching variables

Pattern matching introduces **locally scoped** variables – that is, variables defined only on the right-hand side of the pattern match expression. As an example, consider this code:

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

scala> f( (2,4) )
res0: Int = 6
```

The argument of `f` is the variable `x` of a tuple type `(Int,Int)`, but there is also a pattern variable `x` in the case expression. The pattern variable `x` matches the first part of the tuple and has type `Int`. Because variables are locally scoped, the pattern variable `x` is only defined within the expression `x + y`. The argument `x: (Int,Int)` is a completely different variable whose value has a different type.

The code works correctly but is confusing to read because of the name clash

between the two quite different variables, both named `x`. Another negative consequence of the name clash is that the argument `x: (Int, Int)` is *invisible* within the case expression: if we write “`x`” in that expression, we will get the pattern variable `x: Int`. One says that the argument `x: (Int, Int)` has been **shadowed** by the pattern variable `x`.

The problem is easy to correct: we can give the pattern variable some other name. Since the pattern variable is locally scoped, it can be renamed within its scope without having to change any other code. A completely equivalent code is

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

scala> f( (2,4) )
res0: Int = 6
```

2.6.3 Lazy values and sequences: iterators and streams

We have used streams to create sequences whose length is not known in advance. An example is a stream containing a sequence of increasing positive integers:

```
scala> val p = Stream.iterate(1)(_ + 1)
p: Stream[Int] = Stream(1, ?)
```

At this point, we have not defined a stopping condition for this stream. In some sense, streams are “infinite” sequences, although in practice a stream is always finite because computers cannot run infinitely long. Also, computers cannot store infinitely many values in memory.

To be more precise, streams are “not fully computed” rather than “infinite”. The main difference between arrays and streams is that a stream’s elements are computed on demand and not initially available (except perhaps for the first element), while an array’s elements are all computed in advance and are available immediately. Generally, there are four possible ways a value could be available:

Availability	Explanation	Example Scala code
“eager”	computed in advance	<code>val x = f(123)</code>
“lazy”	computed upon first request	<code>lazy val y = f(123)</code>
“on-call”	computed each time it is requested	<code>def z = f(123)</code>
“never”	cannot be computed due to errors	<code>val (x, y) = "abc"</code>

A **lazy value** (declared as `lazy val` in Scala) is computed only when used in some other expression. Once computed, a lazy value stays in memory and will not be re-computed.

An “on-call” value is re-computed every time it is used. In Scala, this is the behavior of a `def` declaration.

Most collection types in Scala (such as `List`, `Array`, `Set`, and `Map`) are **eager**: all the elements inside these collections are already evaluated. A stream can be seen as a **lazy collection**. Elements of a stream are computed only when needed; after that, they stay in memory and will not be computed again:

```
scala> val str = Stream.iterate(1)(_ + 1)
str: Stream[Int] = Stream(1, ?)

scala> str.take(10).toList
res0: List[Int] = List(1, 2, 3, 4, 5, 6, 7, 8, 9, 10)

scala> str
res1: Stream[Int] = Stream(1, 2, 3, 4, 5, 6, 7, 8, 9, 10, ?)
```

In many cases, it is not necessary to keep previous values of a sequence in memory. For example, consider the computation

```
scala> (1L to 1000000000L).sum
res0: Long = 500000000500000000
```

We do not actually need to keep a billion numbers in memory if we only want to compute their sum. Indeed, the computation just shown does *not* keep all the numbers in memory. The same computation fails if we use a list or a stream:

```
scala> (1L to 1000000000L).toStream.sum
java.lang.OutOfMemoryError: GC overhead limit exceeded
```

The code `(1L to 1000000000L).sum` works because the operation `(1 to n)` produces a sequence whose elements are computed when needed but are not stored in memory. This can be seen as a sequence with the “on-call” availability of elements. Sequences of this sort are called **iterators**. Here are some examples:

```
scala> 1 to 5
res0: scala.collection.immutable.Range.Inclusive = Range(1, 2, 3, 4, 5)

scala> 1 until 5
res1: scala.collection.immutable.Range = Range(1, 2, 3, 4)
```

2 Mathematical formulas as code. II. Mathematical induction

The types `Range` and `Range.Inclusive` are defined in the Scala standard library and are iterators. They behave as collections and support the usual methods (`.map`, `.filter`, etc.), but they do not store previously computed values in memory.

The `.view` method Eager collections such as `List` or `Array` can be converted to iterators by using the `.view` method. This is necessary when intermediate collections consume too much memory when fully evaluated. For example, consider the computation of Example 2.1.5.7 where we used `.flatMap` to replace each element of an initial sequence by three new numbers before computing `.max` of the resulting collection. If instead of three new numbers we wanted to compute *three million* new numbers each time, the intermediate collection created by `.flatMap` would require too much memory, and the computation would crash:

```
scala> (1 to 10).flatMap(x => 1 to 3000000).max
java.lang.OutOfMemoryError: GC overhead limit exceeded
```

Even though the range expression `(1 to 10)` produces an iterator, a subsequent `.flatMap` operation creates an intermediate collection that is too large for our computer's memory. We can use `.view` to avoid this:

```
scala> (1 to 10).view.flatMap(x => 1 to 3000000).max
res0: Int = 3000000
```

The choice between using streams and using iterators is dictated by the memory considerations. Except for that, streams and iterators behave similarly to other sequences. We may write programs in the map/reduce style, applying the standard methods such as `.map`, `.filter`, etc., to streams and iterators. Mathematical reasoning about sequences is the same, whether they are eager, lazy, or on-call.

The broken `Iterator` class The Scala library contains a class called `Iterator`, which has methods such as `Iterator.iterate` and other methods similar to `Stream`. However, `Iterator` actually not an “iterator” in the sense I explained. It cannot be treated as a *value* in the mathematical sense:

```
scala> val iter = (1 until 10).toIterator
iter: Iterator[Int] = non-empty iterator

scala> iter.toList // Look at the elements of 'iter'.
res0: List[Int] = List(1, 2, 3, 4, 5, 6, 7, 8, 9)

scala> iter.toList // Look at these elements again...??
res1: List[Int] = List()
```

```
scala> iter
res2: Iterator[Int] = empty iterator
```

Evaluating the expression `iter.toList` two times produces a *different* result the second time! Also, we see that `iter` became “empty” after the first use.

This situation is impossible in mathematics: if x is some value, such as 100, and f is some function, such as $f(x) = \sqrt{x}$, then $f(x)$ will be the same, $f(100) = \sqrt{100} = 10$, no matter how many times we compute $f(x)$. For instance, we can compute $f(x) + f(x) = 20$ and obtain the correct result. The number $x = 100$ does not “become empty” after the first use; its value remains the same. This behavior is called the **value semantics** of numbers. One says that integers “are values” in the mathematical sense. Alternatively, one says that numbers are **immutable**, i.e. cannot be changed. (What would it mean to “modify” the number 10?)

In programming, a type has value semantics if any computation applied to it always gives the same result. Usually, this means that the type contains immutable data. We can see that Scala’s `Range` has value semantics and is immutable:

```
scala> val x = 1 until 10
x: scala.collection.immutable.Range = Range(1, 2, 3, 4, 5, 6, 7, 8, 9)

scala> x.toList
res0: List[Int] = List(1, 2, 3, 4, 5, 6, 7, 8, 9)

scala> x.toList
res1: List[Int] = List(1, 2, 3, 4, 5, 6, 7, 8, 9)
```

Collections such as `List`, `Map`, or `Stream` are immutable. Some elements of a `Stream` may not be evaluated yet, but this does not affect its value semantics:

```
scala> val str = (1 until 10).toStream
str: scala.collection.immutable.Stream[Int] = Stream(1, ?)

scala> str.toList
res0: List[Int] = List(1, 2, 3, 4, 5, 6, 7, 8, 9)

scala> str.toList
res1: List[Int] = List(1, 2, 3, 4, 5, 6, 7, 8, 9)
```

Iterators produced by applying `.view` also have value semantics:

```
scala> val v = (1 until 10).view
v: scala.collection.SeqView[Int,IndexedSeq[Int]] = SeqView(...)
```

2 Mathematical formulas as code. II. Mathematical induction

```
scala> v.toList
res0: List[Int] = List(1, 2, 3, 4, 5, 6, 7, 8, 9)

scala> v.toList
res1: List[Int] = List(1, 2, 3, 4, 5, 6, 7, 8, 9)
```

Due to the lack of value semantics, programs written using `Iterator` cannot use the tools of mathematical reasoning. This makes it easy to write wrong code that looks correct!

To illustrate the problem, let us re-implement Example 2.5.1.7 by keeping the same code but using `Iterator` instead of `Stream`:

```
def stopRepeatsBad[T](iter: Iterator[T]): Iterator[T] = {
  val halfSpeed = iter.flatMap(x => Seq(x, x))
  halfSpeed.zip(iter) // Do not prepend the first element. It won't help.
    .drop(1)
    .takeWhile { case (h, s) => h != s }
    .map(_._2)
}

scala> stopRepeatsBad(Seq(1, 3, 5, 7, 9, 3, 5, 7, 9).toIterator).toList
res0: List[Int] = List(5, 9, 3, 7, 9)
```

The result `[5,9,3,7,9]` is incorrect, but not in an obvious way: the sequence *was* stopped at a repetition, as we expected, but some of the elements of the given sequence are missing (while other elements are present). It is difficult to debug a program when it produces numbers that are *partially* correct!

The error in this code occurs in the expression `halfSpeed.zip(iter)` due to the fact that `halfSpeed` was itself defined via `iter`. The result is that `iter` is *used twice* in this code, which leads to errors because `iter` is not immutable and does not behave as a value. Creating an `Iterator` and using it twice in the same expression can even fail with an exception:

```
scala> val s = (1 until 10).toIterator
s: Iterator[Int] = non-empty iterator

scala> val t = s.zip(s).toList
java.util.NoSuchElementException: next on empty iterator
```

It is surprising and counter-intuitive that a variable cannot be used twice! We expect code such as `s.zip(s)` to work correctly even though the variable `s` is used twice. When we read the expression `s.zip(s)`, we imagine a given sequence `s` being “zipped” with itself. So we reason that `s.zip(s)` should produce a sequence

of pairs. But Scala's `Iterator` is not immutable, which breaks the usual ways of mathematical reasoning about code.

An `Iterator` can be converted to a `Stream` using the `.toStream` method. This restores the value semantics, since streams are values:

```
scala> val iter = (1 until 10).toIterator
iter: Iterator[Int] = non-empty iterator

scala> val str = iter.toStream
str: Stream[Int] = Stream(1, ?)

scala> str.toList
res0: List[Int] = List(1, 2, 3, 4, 5, 6, 7, 8, 9)

scala> str.toList
res1: List[Int] = List(1, 2, 3, 4, 5, 6, 7, 8, 9)

scala> str.zip(str).toList
res2: List[(Int, Int)] = List((1,1), (2,2), (3,3), (4,4), (5,5), (6,6), (7,7),
    (8,8), (9,9))
```

Instead of `Iterator`, we can use `Stream` and `.view` when lazy or on-call collections are required.

3 The formal logic of types. I.

Higher-order functions

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 => x - y`, which is a function of type `Int => Int`. So the type of `f0` is written as `Int => (Int => Int)`.

To use `f0`, we must 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, `Int => Int => Int` means the same as `Int => (Int => Int)`, and `x => y => x - y` means the same as `x => (y => x - y)`. In other words, the function symbol `=>` associates to the right. Thus, the type signature of `f0` may be equivalently written

3 The formal logic of types. I. Higher-order functions

as `Int => Int => Int`.

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

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

Let us compare the function `f1` with a function that takes its two arguments at once. Such a function will have a different type signature, for instance

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

has type signature `(Int, Int) => Int`.

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, – that is, applied only to its first argument, 20. The result of `f1(20)` is 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 => B => C => P => Z`, where `A`, `B`, `C`, ..., `Z` are some types. I prefer to think about this kind of type signature as having arguments of types `A`, `B`, ..., `P`, called the **curried arguments** of the function, and the “final” return value of type `Z`. The “final” return value of the function is returned after supplying all arguments.

A function with this type signature is, in a sense, equivalent to an **uncurried** function with type signature `(A,B,C,...,P) => 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. One says that a curried function is **isomorphic** or **equivalent** to an uncurried function.

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 sometimes harder to read.

3.1 Types of higher-order functions

In the syntax 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 tradition 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, compared with the syntax where commas are used to separate the arguments. (Consider, for example, the Haskell expression `f2 (g 10) (h 20) + 30`.) To improve readability of code, programmers may prefer to first define short names for complicated expressions and then use these names as curried arguments.

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 give 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 that calls `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`).

A feature of Scala is the “dotless” method syntax: for example, `xs map f` is equivalent to `xs.map(f)`. The “dotless” syntax works only for infix methods, such as `.map`, defined on specific types such as `Seq`. Do not confuse Scala’s “dotless” method syntax with the short function application syntax such as `fmap f xs`, used in Haskell and some other languages.

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

3 The formal logic of types. I. Higher-order functions

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 $\text{Int} \Rightarrow \text{Int}$.

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 a colon in the superscript, for example: $x^{\text{Int}} \Rightarrow x + 2$.

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

3.1 Types of higher-order functions

these examples.

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

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 $((f\ x)\ y)\ z$.

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. I. Higher-order functions

- 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 work with 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.

Examples:

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

The function `f1` has type signature `Int \Rightarrow 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 `Int \Rightarrow Int \Rightarrow 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 `(Int \Rightarrow Int) \Rightarrow 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
```


3.1 Types of higher-order functions

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 (defined using `def`), another as a function value (defined 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 => 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`.

3 The formal logic of types. I. Higher-order functions

- 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  $\Rightarrow$  Boolean) = x  $\Rightarrow$  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  $\Rightarrow$  Int) = k  $\Rightarrow$  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  $\Rightarrow$  Int) = x  $\Rightarrow$  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  $\Rightarrow$  Int) = _  $\Rightarrow$  123
```

- 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  $\Rightarrow$  Int = _  $\Rightarrow$  x
val ff_v: (Int  $\Rightarrow$  Int  $\Rightarrow$  Int) = x  $\Rightarrow$  _  $\Rightarrow$  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  $\Rightarrow$  Double, g: Double  $\Rightarrow$  String): (Int  $\Rightarrow$  String) =
  x  $\Rightarrow$  g(f(x))
scala> val f: Int  $\Rightarrow$  Double = x  $\Rightarrow$  5.67 + x
scala> val g: Double  $\Rightarrow$  String = x  $\Rightarrow$  f"x=%3.2f"
scala> val h = comp(f, g)
scala> h(10)
```

The function `comp` has two arguments, of types `Int \Rightarrow Double` and `Double \Rightarrow String`. The result value of `comp` is of type `Int \Rightarrow 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 `Int => Int => Int` and returns an uncurried equivalent function of type `(Int, Int) => 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 $x \Rightarrow \frac{1}{1+x}$.

3 The formal logic of types. I. Higher-order functions

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,

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

Let us decipher this confusing 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 clear that this function actually ignores its argument x .

3.3 Exercises

1. Define a function of type `Int => List[List[Int]] => List[List[Int]]` similar to Exercise 1.6.2.4 except that the hard-coded number 100 must be a *curried* first argument. Implement Exercise 1.6.2.4 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 `(Int, Int) => Int` and returns a curried equivalent function of type `Int => Int => Int`.

4 The formal logic of types. II.

Disjunctive types

4.1 Discussion

4.1.1 Scala's case classes as “named tuple” types

It is often convenient to use names for the different parts of a tuple. Scala's “case classes” allow programmers to put names on each tuple part, as well as on the tuple type itself. The syntax is

```
case class Person(firstName: String, lastName: String, age: Int)
```

This data type carries the same information as a tuple (`String, String, Int`). Without using a case class to hold this information, we could define a type alias for this tuple,

```
type PersonTuple = (String, String, Int)
```

and write code like this,

```
scala> val p: PersonTuple = ("Albert", "Einstein", 140)
p: PersonTuple = (Albert,Einstein,140)

scala> val einsteinName = p._2
einsteinName: String = Einstein
```

However, this type alias only creates an alternative name for an existing tuple of data. The declaration of a `case class Person` gives the programmer several methods that make working with this tuple more convenient.

In one sense, case classes behave as named tuples with named parts. Creating a value of a case class and accessing its parts looks like this:

```
scala> val einstein = Person(firstName = "Albert", lastName = "Einstein", age = 140)
```

4 The formal logic of types. II. Disjunctive types

```
einstein: Person = Person(Albert,Einstein,140)

scala> einstein.firstName
res0: String = Albert

scala> einstein.age
res1: Int = 140
```

Part names are optional when creating a value:

```
scala> val poincare = Person("Henri", "Poincare", 165)
poincare: Person = Person(Henri,Poincare,165)
```

It is a type error to use wrong types with a case class:

```
scala> val p = Person(140, "Einstein", "Albert")
<console>:13: error: type mismatch;
 found   : Int(140)
 required: String
    val p = Person(140, "Einstein", "Albert")
                      ^
<console>:13: error: type mismatch;
 found   : String("Albert")
 required: Int
    val p = Person(140, "Einstein", "Albert")
                      ^
```

However, parts can be specified in any order when using part names:

```
scala> val p = Person(age = 140, lastName = "Einstein", firstName = "Albert")
p: Person = Person(Albert,Einstein,140)
```

Another feature of case classes is that each case class is a *different type*, even if the data it contains is the same. A function whose argument is of type `Person` requires an argument of that type. It will be a type error to call that function with an argument of type `(String, String, Int)`. One cannot make a mistake using one case class instead of another.

5 The formal logic of types. III. The Curry-Howard correspondence

5.0.1 Discussion

6 Functors

6.1 Discussion

6.2 Practical use

6.3 Laws and structure

7 Type-level functions and typeclasses

7.1 Combining typeclasses

7.2 Inheritance

7.3 Functional dependencies

7.4 Discussion

8 Computations in functor blocks.

I. Filterable functors

8.1 Practical use

8.1.1 Discussion

8.2 Laws and structure

8.2.1 Discussion

9 Computations in functor blocks. II. Semimonads and monads

9.1 Practical use

9.1.1 Discussion

9.2 Laws and structure

9.2.1 Discussion

10 Applicative functors, contrafunctors, and profunctors

10.1 Practical use

10.1.1 Discussion

10.2 Laws and structure

11 Traversable functors and profunctors

11.1 Discussion

12 “Free” type constructions

12.1 Discussion

13 Computations in functor blocks.

III. Monad transformers

13.1 Practical use

13.2 Laws and structure

13.2.1 Laws of monad transformers

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 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 ftr_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 laws:** 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 lifts monadic morphisms from $M \rightsquigarrow N$ to $T_L^M \rightsquigarrow T_L^N$ and must satisfy the corresponding **lifting laws**:

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

It follows from the identity law $T_L^{\text{Id}} \cong L$ that the base monad L can be lifted into T_L^M : Setting $\phi = \text{pu}_M : \text{Id} \rightsquigarrow M$, we obtain

$$\text{mrun}(\text{pu}_M) : T_L^{\text{Id}} \rightsquigarrow T_L^M = L \rightsquigarrow T_L^M .$$

13 Computations in functor blocks. III. Monad transformers

This function is called the **base lifting**, $\text{mr}(\text{pu}_M) \triangleq \text{blift} : L^A \Rightarrow T_L^{M,A}$. The base lifting automatically satisfies the non-degeneracy law,

$$\text{blift} \circ \text{mr}(\phi^{M \sim \text{Id}}) = \text{id} \quad ,$$

for any monadic morphism $\phi : M \rightsquigarrow \text{Id}$, because the left-hand side equals $\text{mr}(\text{pu}_M \circ \phi)$, and the composition law for monadic morphisms gives $\text{pu}_M \circ \phi = \text{pu}_{\text{Id}} = \text{id}$.

5. **Base runner laws:** For any monadic morphism $\theta : L \rightsquigarrow \text{Id}$ and for any monad M , the base runner $\text{br}(\theta) : T_L^M \rightsquigarrow M$ is a monadic morphism. The base runner must also satisfy the **non-degeneracy law**,

$$\text{lift} \circ \text{br}(\theta) = \text{id} \quad .$$

Since it is not possible to transform the base monad L into an arbitrary other monad, there are no lifting laws for br , unlike mr . So the non-degeneracy law is not an automatic consequence of other laws.

13.2.2 Examples of incorrect monad transformers

The laws of monad transformers guarantee that the transformed monad is able to represent, without loss of information, the operations of the base monad as well as the operations of the foreign monad. If some of these laws are omitted, we may obtain a type constructor that has methods with the required type signatures but does not work correctly.

The simplest example of an incorrect monad transformer is obtained if we define the transformed monad to be the unit monad, $T_L^{M,A} \triangleq 1$, for any monads L and M . It is clear that this monad transformer is “fake”: it cannot possibly keep the information about the monads L and M , because the methods of the unit monad discard all information and return 1. However, the transformer $T_L^{M,A}$ has all the required methods pu_T , ftn_T , lift_T , mr_T , and br_T with correct type signatures (they are constant functions returning 1). All these functions are automatically monadic morphisms, since the morphism from any monad to the unit monad is a monadic morphism. We find that many of the monad transformer laws hold! However, the identity law $T_L^{\text{Id}} \cong L$ and the non-degeneracy law $\text{lift} \circ \text{br}(\theta) = \text{id}$ are violated since $T_L^{\text{Id}} = 1 \not\cong L$ and $\text{lift} \circ \text{br}(\theta) = _ \Rightarrow 1 \neq \text{id}$. For this reason, the unit monad is not a lawful monad transformer.

This simple example demonstrates the importance of the monad transformer laws. A malicious programmer could give us a “fake” implementation of a monad transformer that appears to have all the methods with the correct type signatures but, instead of a bigger monad T_L^M , constructs a unit monad dressed up as a type constructor $T_L^{M,A}$. The only way for us to detect the fraud is to find that the identity law and the non-degeneracy law are violated.

Other examples of “fake” transformers violating some of the laws are $T_L^M = L$ (no lifting law) and $T_L^M = M$ (no identity law).

In these cases, it is intuitively clear that the “fake” monad transformer definitions are incorrect because the information about either L or M is missing in T_L^M . A potentially working definition of T_L^M must be a type constructor whose definition somehow combines *both* type constructors L and M .

13.2.3 Examples of failure to define a generic monad transformer

It appears to be impossible to define T_L^M as a generic construction that works in the same way for all monads L and M . We will now consider a few ways of combining the type constructors L and M in a way that is independent of their structure. In all these cases, we will find that some of the monad transformer laws are violated.

General ways of combining two type constructors L^\bullet and M^\bullet are functor composition L^{M^\bullet} or M^{L^\bullet} , disjunction $L^\bullet + M^\bullet$, and product $L^\bullet \times M^\bullet$.

Functor composition A general way of combining two type constructors L^\bullet and M^\bullet is the functor composition L^{M^\bullet} or M^{L^\bullet} . However, the functor composition works only for certain monads and only in a certain order; so it cannot work as a generic monad transformer. A simple counter-example is $L^A \triangleq 1 + A$ and $M^A \triangleq A \times A$ where M^{L^A} is a monad but L^{M^A} is not (see Section ???). Another counter-example is the `State` monad, $\text{State}_S^A \triangleq S \Rightarrow S \times A$, for which we have already shown that $1 + \text{State}_S^A$ is not a monad and $\text{State}_S^{Z \Rightarrow A}$ is not a monad (see Section ???). In other words, the `State` monad does not compose with arbitrary monads M in either order.

Functor disjunction The functor disjunction $L^\bullet + M^\bullet$ is in general not a monad when L and M are arbitrary monads. An immediate counter-example is found by using two `Reader` monads, $L^A \triangleq R \Rightarrow A$ and $M^A \triangleq S \Rightarrow A$. The disjunction $(R \Rightarrow A) + (S \Rightarrow A)$ is a functor that is not a monad (and not even applicative, see Section ???).

Functor product The functor product $L^\bullet \times M^\bullet$ is a monad for arbitrary monads L and M . However, there is no naturally defined lift : $M^\bullet \rightsquigarrow L^\bullet \times M^\bullet$ because we cannot create values of type L^A out of values of type M^A for arbitrary monads L and M .

Using the free monad The functor composition L^{M^\bullet} and the disjunction $L^\bullet + M^\bullet$ may not always be monads, but they are always functors. So we can make monads out of them, by using the free monad construction. We get $\text{Free}^{L^{M^\bullet}}$, the free monad over L^{M^\bullet} , and $\text{Free}^{L^\bullet + M^\bullet}$, the free monad over $L^\bullet + M^\bullet$. Many laws of the monad transformer are satisfied by these constructions. However, the identity laws fail because

$$\text{Free}^{L^{\text{Id}^\bullet}} \cong \text{Free}^{L^\bullet} \not\cong L \quad , \quad \text{Free}^{L^\bullet + \text{Id}^\bullet} \not\cong L \quad ,$$

and the lifting laws are also violated because lift : $M^A \Rightarrow \text{Free}^{L^\bullet + M^\bullet, A}$ is not a monad morphism because it maps pu_M into a non-pure value of the free monad. Nevertheless, these constructions are not useless. Once we interpret the free monad into a concrete (non-free) monad, we could arrange to hide the violations of these laws, so that the monad laws hold for the resulting (non-free) monad.

“Monoidal convolution” The construction called “**monoidal convolution**” defines a new functor $L \star M$ via

$$(L \star M)^A \triangleq \exists P \exists Q. (P \times Q \Rightarrow A) \times L^P \times M^Q \quad . \quad (13.1)$$

This formula can be seen as a combination of the co-Yoneda identities

$$L^A \cong \exists P. L^P \times (P \Rightarrow A) \quad , \quad M^A \cong \exists Q. M^Q \times (Q \Rightarrow A) \quad .$$

The functor product $L \times M$ is equivalent to

$$\begin{aligned} L^A \times M^A & \\ \text{co-Yoneda identities for } L^A \text{ and } M^A : & \cong \exists P. L^P \times \underline{(P \Rightarrow A)} \times \exists Q. M^Q \times \underline{(Q \Rightarrow A)} \\ \text{equivalence in Eq. (13.3) :} & \cong \exists P. \exists Q. L^P \times M^Q \times (P + Q \Rightarrow A) \end{aligned} \quad (13.2)$$

where we used the type equivalence

$$(P \Rightarrow A) \times (Q \Rightarrow A) \cong P + Q \Rightarrow A \quad . \quad (13.3)$$

If we (arbitrarily) replace $P + Q \Rightarrow A$ by $P \times Q \Rightarrow A$ in Eq. (13.2), we will obtain Eq. (13.1).

The monoidal convolution $L \star M$ always produces a functor since Eq. (13.1) is covariant in A . An example where the monoidal convolution fails to produce a monad transformer is $L^A \triangleq 1 + A$ and $M^A \triangleq R \Rightarrow A$. We compute the functor $L \star M$ and establish that it is not a monad:

$$\begin{aligned}
 (L \star M)^A & \\
 \text{definitions of } L, M, \star : &= \exists P \exists Q. (P \times Q \Rightarrow A) \times (\mathbb{1} + P) \times (R \Rightarrow Q) \\
 \text{curry the arguments, move quantifier :} &= \exists P. (\mathbb{1} + P) \times \exists Q. (Q \Rightarrow P \Rightarrow A) \times (R \Rightarrow Q) \\
 \text{co-Yoneda identity with } \exists Q : &= \exists P. (\mathbb{1} + P) \times (R \Rightarrow P \Rightarrow A) \\
 \text{swap curried arguments :} &= \exists P. (\mathbb{1} + P) \times (P \Rightarrow R \Rightarrow A) \\
 \text{co-Yoneda identity with } \exists P : &= \mathbb{1} + (R \Rightarrow A) \quad .
 \end{aligned}$$

This functor is not a monad (see Section ???).

Codensity tricks***

- codensity monad over L^{M^\bullet} : $F^A \triangleq \forall B. (A \Rightarrow L^{M^B}) \Rightarrow L^{M^B}$ – no lift
- Codensity- L transformer: $\text{Cod}_L^{M,A} \triangleq \forall B. (A \Rightarrow L^B) \Rightarrow L^{M^B}$ – no lift
 - applies the continuation transformer to $M^A \cong \forall B. (A \Rightarrow B) \Rightarrow M^B$
- Codensity composition: $F^A \triangleq \forall B. (M^A \Rightarrow L^B) \Rightarrow L^B$ – not a monad
 - Counterexample: $M^A \triangleq R \Rightarrow A$ and $L^A \triangleq S \Rightarrow A$

13.2.4 Properties of monadic morphisms

Statement 13.2.4.1 For any monad M , the method $\text{pu}_M : A \Rightarrow M^A$ is also a monadic morphism $\text{pu}_M : \text{Id} \rightsquigarrow M$ between the identity monad and M .

Proof The identity law requires $\text{pu}_{\text{Id}} \circ \text{pu}_M = \text{pu}_M$. This holds because $\text{pu}_{\text{Id}} = \text{id}$. The composition law requires $\text{ftn}_{\text{Id}} \circ \text{pu}_M = \text{pu}_M^{\uparrow \text{Id}} \circ \text{pu}_M \circ \text{ftn}_M$. Since $\text{ftn}_{\text{Id}} = \text{id}$ and $f^{\uparrow \text{Id}} = f$ for any function f , we simplify both sides of the composition law

13 Computations in functor blocks. III. Monad transformers

to the same expression pu_M :

$$\begin{aligned} \text{ftn}_{\text{Id}} \circ \text{pu}_M &= \text{pu}_M \quad , \\ \text{pu}_M^{\uparrow \text{Id}} \circ \text{pu}_M \circ \text{ftn}_M & \\ \text{raising into the identity functor:} &= \text{pu}_M \circ \text{pu}_M \circ \text{ftn}_M \\ \text{left identity law for } M : &= \text{pu}_M \quad . \end{aligned}$$

Exercise 13.2.4.2 Suppose M is a given monad, Z is a fixed type, and a fixed value $m : M^Z$ is given. (a) Consider the function f defined as

$$\begin{aligned} f : (Z \Rightarrow A) &\Rightarrow M^A \quad , \\ f(q^{Z \Rightarrow A}) &\triangleq q^{\uparrow M} m \quad . \end{aligned}$$

Prove that f is *not* a monadic morphism from the reader monad $R^A \triangleq Z \Rightarrow A$ to the monad M , despite having the correct type signature.

(b) Under the same assumptions, consider the function ϕ defined as

$$\begin{aligned} \phi : (Z \Rightarrow M^A) &\Rightarrow M^A \quad , \\ \phi(q^{Z \Rightarrow M^A}) &\triangleq \text{flm}_M(q)(m) \quad . \end{aligned}$$

Show that ϕ is *not* a monadic morphism from the monad $Q^A \triangleq Z \Rightarrow M^A$ to M .

Statement 13.2.4.3 If $\phi : M^\bullet \rightsquigarrow N^\bullet$, equivalently written as $\phi : M^A \Rightarrow N^A$, is a monadic morphism between monads M and N , then ϕ is also a natural transformation between functors M and N .

Statement 13.2.4.4 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 13.2.4.5 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 13.2.4.6 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 13.2.4.7 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 13.2.4.8 For any monads M and N , the component-swapping function $\sigma : M^\bullet \times N^\bullet \leadsto N^\bullet \times M^\bullet$ is a monadic morphism.

Proof The code for σ can be written as $\sigma = \Delta_! (\nabla_2 \boxtimes \nabla_1)$. When σ is written in this way, it follows that σ is a composition of Δ , which is a monadic morphism by Statement 13.2.4.5, and $\nabla_2 \boxtimes \nabla_1$, which is a monadic morphism by Statement 13.2.4.6. A composition of monadic morphisms is a monadic morphism by Statement 13.2.4.4, so σ is a monadic morphism.

13.2.5 Functor composition with 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 the transformed 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 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 into functors as necessary. The result will be a number of nested layers of T . Second, we use `ftnT` as many times as necessary to flatten all nested layers of T into a single layer. The result is a value of type T^A .

13.2.6 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.

We will now show that the transformer T_R is lawful (satisfies all five laws shown in Section 13.2.1), as long as both T_P and T_Q satisfy the same five laws. To shorten the notation, we 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 construction law for T_Q says that T_Q^M is a monad. The monad construction law for T_P says that T_P^S is a monad for any monad S ; in particular, for $S = T_Q^M$. Therefore, $T_P^S = T_P^{T_Q^M}$ 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 already have a monadic morphism $\phi : Q \rightsquigarrow T_Q^{\text{Id}}$ and its inverse, $\chi : T_Q^{\text{Id}} \rightsquigarrow Q$. The runner mrun_P for T_P can be applied to both ϕ and χ since they are monadic morphisms. So we obtain two new monadic morphisms,

$$\text{mrun}_P(\phi) : T_P^Q \rightsquigarrow T_P^{T_Q^{\text{Id}}} \quad ; \quad \text{mrun}_P(\chi) : T_P^{T_Q^{\text{Id}}} \rightsquigarrow T_P^Q .$$

Are these two monadic morphisms inverses of each other? To show this, we need to verify that

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

By the runner law for T_P , we have $\text{mrun}_P(f) \circ \text{mrun}_P(g) = \text{mrun}_P(f \circ g)$ for any two monadic morphisms f and g . We also have $\text{mrun}_P(\text{id}) = \text{id}$ by the same

law. So,

$$\begin{aligned}\text{mrun}_P(\phi) \circ \text{mrun}_P(\chi) &= \text{mrun}_P(\phi \circ \chi) = \text{mrun}_P(\text{id}) = \text{id} \quad , \\ \text{mrun}_P(\chi) \circ \text{mrun}_P(\phi) &= \text{mrun}_P(\chi \circ \phi) = \text{mrun}_P(\text{id}) = \text{id} \quad .\end{aligned}$$

We have indeed obtained a monadic isomorphism between T_P^Q and T_P^{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 $M \rightsquigarrow T_P^{T_Q^M}$. (A composition of monadic morphisms is again a monadic morphism by Statement 13.2.4.4.)

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

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

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

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

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

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

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

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

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

13 Computations in functor blocks. III. Monad transformers

$\text{id} : M \rightsquigarrow M,$

$$\text{lift}_R (\text{id}^{M \rightsquigarrow M}) = \text{lift}_P (\text{lift}_Q \text{id}^{M \rightsquigarrow M})$$

$$\text{identity law for lift}_Q : = \text{lift}_P (\text{id}^{T_Q^M \rightsquigarrow T_Q^M})$$

$$\text{identity law for lift}_P : = \text{id} \quad .$$

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

$$\text{lift}_R (\phi \circ \chi) = \text{lift}_P (\text{lift}_Q (\phi \circ \chi))$$

$$\text{composition law for lift}_Q : = \text{lift}_P (\text{lift}_Q \phi \circ \text{lift}_Q \chi)$$

$$\text{composition law for lift}_P : = \text{lift}_P (\text{lift}_Q \phi) \circ \text{lift}_P (\text{lift}_Q \chi)$$

$$\text{definition of lift}_R : = \text{lift}_R \phi \circ \text{lift}_R \chi \quad .$$

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 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 we can define the monadic morphism $\text{brun}_R\theta$ as

$$\begin{aligned} \text{brun}_R\theta : T_P^{T_Q^M} &\leadsto M \quad , \\ \text{brun}_R\theta &\triangleq \text{brun}_P (\text{blift}_P \circ \theta) \circ \text{brun}_Q (\text{lift}_P \circ \theta) \quad . \end{aligned}$$

Since we have defined $\text{brun}_R\theta$ as a composition of monadic morphisms, $\text{brun}_R\theta$ is a monadic morphism by Statement 13.2.4.4.

To verify the non-degeneracy law of the base runner, $\text{lift}_R \circ \text{brun}_R\theta = \text{id}$, we need to use the non-degeneracy laws for the base runners of T_P and T_Q , which are

$$\text{lift}_P \circ \text{brun}_P \chi^{P \leadsto \text{Id}} = \text{id} \quad , \quad \text{lift}_Q \circ \text{brun}_Q \psi^{Q \leadsto \text{Id}} = \text{id} \quad .$$

Then we can write

$$\begin{aligned} &\text{lift}_R \circ \text{brun}_R\theta \\ \text{expand definitions :} &= \text{lift}_Q \circ \text{lift}_P \circ \text{brun}_P (\text{blift}_P \circ \theta) \circ \text{brun}_Q (\text{lift}_P \circ \theta) \\ \text{non-degeneracy for brun}_P : &= \text{lift}_Q \circ \text{brun}_Q (\text{lift}_P \circ \theta) \\ \text{non-degeneracy for brun}_Q : &= \text{id} \quad . \end{aligned}$$

13.2.7 Stacking any number of 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:

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

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 now define a transformer stack for any number of monads P , Q , ..., Z in a similar way,

$$T_S^{M, A} \triangleq T_P^{T_Q^{\dots T_Z^M}, A} \quad . \quad (13.4)$$

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 turns out 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. (13.4). Let us temporarily denote by J the monad

$$J \triangleq T_Q \cdot \overset{T^{\text{Id}}}{\cdot} ,$$

where we used the identity monad Id in the place normally taken by a foreign monad M . 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 as

$$T_J^M \triangleq T_Q \cdot \overset{T^M}{\cdot} .$$

Since both T_P and T_J are lawful transformers, their stacking composition $T_P^{T_J^M}$ is also a lawful transformer (this was shown in the Section 13.2.6). 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.

13.3 Monad transformers via functor composition: General properties

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

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

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

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

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

13.3 Monad transformers via functor composition: General properties

We can prove many properties of both kinds of monad transformers via a single derivation if we temporarily drop the distinction between the base monad and the foreign monad. We simply assume that two different monads, L and M , have a functor composition $T^\bullet \triangleq L^{M^\bullet}$ that 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 `mrtn` and `brun`, which we will need to derive separately for each of the two kinds of transformers.

13.3.1 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

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

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

$$\text{pu}_T \triangleq \text{pu}_M \circ \text{pu}_L \quad . \tag{13.5}$$

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

It remains to implement `ftnT`. 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

13 Computations in functor blocks. III. Monad transformers

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 $L \circ L \circ M \circ M$ 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 have a function called `swap` (short notation “sw”), which can interchange the order of the layers. The type signature of `swap` is

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

which is equivalently written in a more verbose 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

$$\begin{array}{ccccc} L^{M^L M^A} & \xrightarrow{\text{sw}^{\uparrow L}} & L^{L M^M A} & \xrightarrow{\text{ftn}_L} & L^{M^M A} \\ & \searrow \text{ftn}_T \triangleq & & & \downarrow \text{ftn}_M^{\uparrow L} \\ & & & & L^{M^A} \end{array} \quad . \quad (13.6)$$

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`

13.3 Monad transformers via functor composition: General properties

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 difference between the operations `swap` and `sequence` There is a certain similarity between the `swap` operation introduced here and the `sequence` operation introduced in Chapter 11 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 stops here. The laws required for the `swap` operation to yield a monad T are different from the laws of traversable functors. In particular, if we wish M^{L^\bullet} 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 functors); 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 the applicative functor F .

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

13.3.2 Deriving the necessary 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 $f : A \Rightarrow B$,

$$f^{\uparrow L \uparrow M} \circ \text{sw} = \text{sw} \circ f^{\uparrow M \uparrow L} \quad . \quad (13.7)$$

¹The `swap` operation was used 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.

13 Computations in functor blocks. III. Monad transformers

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

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

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

Additionally, T must satisfy the laws of a monad transformer. We will now discover the laws for `swap` that make the laws for `ftnT` hold automatically, as long as `ftnT` is derived from `swap` using Eq. (13.6).

We substitute Eq. (13.6) into the left identity law for `ftnT` and simplify:

$$\begin{aligned}
 \text{id} &= \text{pu}_T \circ \underline{\text{ftn}_T} \\
 \text{replace } \text{ftn}_T \text{ using Eq. (13.6)} : &= \text{pu}_T \circ \text{sw}^{\uparrow L} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \\
 \text{replace } \text{pu}_T \text{ using Eq. (13.5)} : &= \text{pu}_M \circ \underline{\text{pu}_L \circ \text{sw}^{\uparrow L} \circ \text{ftn}_L} \circ \text{ftn}_M^{\uparrow L} \\
 \text{naturality of } \text{pu}_L : &= \text{pu}_M \circ \text{sw} \circ \underline{\text{pu}_L \circ \text{ftn}_L} \circ \text{ftn}_M^{\uparrow L} \\
 \text{left identity law for } L : &= \text{pu}_M \circ \text{sw} \circ \text{ftn}_M^{\uparrow L} \quad . \tag{13.8}
 \end{aligned}$$

How could the last expression in Eq. (13.8) 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. (13.8) 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 “outer-identity law” for `swap`,

$$\text{pu}_M \circ \text{sw} = \text{pu}_M^{\uparrow L} \quad . \tag{13.9}$$

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

13.3 Monad transformers via functor composition: General properties

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

With this law, we can finish the derivation in Eq. (13.8) as

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

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

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

$$\text{pu}_L^{\uparrow M} \circ \text{sw} = \text{pu}_L \quad . \quad (13.10)$$

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

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 :

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

Deriving the monad associativity law for T ,

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

13 Computations in functor blocks. III. Monad transformers

turns out to require *two* further laws for `swap`. Let us see why.

Substituting the definition of `ftnT` into the associativity law, we get

$$\begin{aligned} & (\text{sw}^{\uparrow L} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L})^{\uparrow M \uparrow L} \circ \text{sw}^{\uparrow L} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \\ &= \text{sw}^{\uparrow L} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \circ \text{sw}^{\uparrow L} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} . \end{aligned} \quad (13.11)$$

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 $M^{M^{L^A}} \Rightarrow L^{M^A}$ or, in another notation, $M \circ M \circ L \rightsquigarrow L \circ M$. 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 “outer-interchange” law

$$\text{ftn}_M \circ \text{sw} = \text{sw}^{\uparrow M} \circ \text{sw} \circ \text{ftn}_M^{\uparrow L} . \quad (13.12)$$

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

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

$$\text{ftn}_L^{\uparrow M} \circ \text{sw} = \text{sw} \circ \text{sw}^{\uparrow L} \circ \text{ftn}_L . \quad (13.13)$$

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

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

13.3 Monad transformers via functor composition: General properties

Both sides of the law in Eq. (13.11) involve compositions of several `flatten`s and `swap`s. The heuristic idea of the proof is to use various laws to move all `flatten`s to right of the composition, while moving all `swap`s to the left. In this way we will transform both sides of Eq. (13.11) into a similar form, hoping to prove that they are equal.

We begin with the right-hand side of Eq. (13.11) since it is simpler than the left-hand side, and look for ways of using the interchange laws. At every step of the calculation, there happens to be only one place where some law can be applied:

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

Now all `swap`s are on the left and all `flatten`s on the right of the expression.

Transform the right-hand side of Eq. (13.11) in the same way as

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

We have again managed to move all `swap`s to the left and all `flatten`s to the right of the expression.

13 Computations in functor blocks. III. Monad transformers

Comparing now the two sides of the associativity law, we see that all the `flatten`s 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

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

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 13.3.2.1 If two monads L and M are such that there exists a function

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

(called “`swap`”), which is a natural transformation satisfying four additional laws:

$$\text{outer-identity} : \text{pu}_L^{\uparrow M} \circ \text{sw} = \text{pu}_L ,$$

$$\text{inner-identity} : \text{pu}_M \circ \text{sw} = \text{pu}_M^{\uparrow L} ,$$

$$\text{outer-interchange} : \text{ftn}_L^{\uparrow M} \circ \text{sw} = \text{sw} \circ \text{sw}^{\uparrow L} \circ \text{ftn}_L ,$$

$$\text{inner-interchange} : \text{ftn}_M \circ \text{sw} = \text{sw}^{\uparrow M} \circ \text{sw} \circ \text{ftn}_M^{\uparrow L} ,$$

then the functor composition

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

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

$$\text{pu}_T \triangleq \text{pu}_M \circ \text{pu}_L , \tag{13.14}$$

$$\text{ftn}_T \triangleq \text{sw}^{\uparrow L} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} . \tag{13.15}$$

13.3.3 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,

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

13.3 Monad transformers via functor composition: General properties

we can flatten the layers using ftn_L , or ftn_M , ftn_T , 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 some L -effects in a separate functor block of $L\text{-flatMaps}$ 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 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 .

13.3.4 Deriving swap from flatten

We have shown that the flatten method of 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 the composed monad T . Does a suitable swap function exist for these examples? In other words, if a flatten function for the monad $T = L \circ M$ is already known, can we establish whether a swap function exists such that the *given* flatten function is expressed via Eq. (13.6)?

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

$$\text{sw} = \text{pu}_M^{\uparrow L \uparrow M} \circ \text{pu}_L \circ \text{ftn}_T \quad . \quad (13.16)$$

13 Computations in functor blocks. III. Monad transformers

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

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

1. If we first define ftn_T using Eq. (13.6) through a given implementation of swap and then substitute that ftn_T into Eq. (13.16), will we always recover the initially given function sw ? (Yes, assuming naturality for swap .)
2. If we first define sw using Eq. (13.16) through a given implementation of ftn_T and then substitute that sw into Eq. (13.6), will we always recover the initially given function ftn_T ? (No, not without additional laws for ftn_T .)

To answer the first question, substitute ftn_T from Eq. (13.6) into Eq. (13.16):

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

So, indeed, we always recover the initial swap function.

To answer the second question, substitute sw from Eq. (13.16) into Eq. (13.6):

$$\begin{aligned}
 & \text{sw}^{\uparrow L} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \\
 \text{substitute sw} : &= \left(\text{pu}_M^{\uparrow L \uparrow M} \circ \text{pu}_L \circ \text{ftn}_T \right)^{\uparrow L} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \\
 \text{functor composition} : &= \text{pu}_M^{\uparrow L \uparrow M \uparrow L} \circ \text{pu}_L^{\uparrow L} \circ \text{ftn}_T^{\uparrow L} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \quad . \quad (13.17)
 \end{aligned}$$

13.3 Monad transformers via functor composition: General properties

At this point, we are stuck: we can find no laws to transform the last expression. Without assuming additional laws, it *does not follow* that the right-hand side of Eq. (13.17) is equal to ftn_T . Let us now derive those additional laws.

The only sub-expression in Eq. (13.17) that we could possibly transform is the composition $\text{ftn}_T^{\uparrow L} \circ \text{ftn}_L$. So, we need to assume a law involving the expression

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

This function flattens the two layers of $(L \circ M)$ and then flattens the remaining two layers of L . Another function with the same type signature could first flatten the two *outside* layers of L and then 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 expresses a kind of “compatibility” between the monads L and T .

With this law, the right-hand side of Eq. (13.17) 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 assume a law involving the expression

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

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 :

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

13 Computations in functor blocks. III. Monad transformers

The second conjectured law is therefore

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

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

This law expresses a kind of “compatibility” between the monads M and T .

Assuming this law, we can finally complete the derivation:

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

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 13.3.4.1).

It is often hard to verify the monad laws for $L \circ M$ directly because of deeply nested type constructors, e.g. $L \circ M \circ L \circ M \circ L \circ M$. If the monad $L \circ M$ has a `swap` method (in practice, this is always the case), it is simpler to verify the laws of `swap` and then obtain the monad laws of $L \circ M$ via Theorem 13.3.2.1.

Exercise 13.3.4.1 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. (13.6),

show that the two interchange laws must hold for ftn_T :

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

13.3 Monad transformers via functor composition: General properties

Exercise 13.3.4.2 With the same assumptions as Exercise 13.3.4.1 and additionally assuming the inner and outer identity laws for swap (see Theorem 13.3.2.1), show that the monad $T^\bullet \triangleq L^{M^\bullet}$ satisfies two “pure compatibility” laws,

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

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

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

13.3.5 Laws of monad transformer liftings: Proofs

We will now derive the laws of monad transformer liftings from the laws of swap , using Eqs. (13.14)–(13.15) 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

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

The identity laws are verified quickly,

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

13 Computations in functor blocks. III. Monad transformers

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

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

The composition law for blift :

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

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

13.3.6 Laws of monad transformer runners: Proofs

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 with respect to the *foreign* monad. So, this law needs to be written differently, depending on the choice of the base monad. Let us consider separately the situations when either L or M is the base monad.

If the base monad is L , the runners are

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

13.3 Monad transformers via functor composition: General properties

The laws of runners require that $\text{mrun } \phi$ and $\text{brun } \theta$ must be monadic morphisms, i.e. the identity and composition laws must hold for $\text{mrun } \phi$ and $\text{brun } \theta$:

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

To derive these laws, we may use the identity and composition laws of monadic morphisms for ϕ and θ . We also use Eqs. (13.14)–(13.15) as definitions of the monad T . Additionally, the **monadic naturality** of swap with respect to ϕ and θ are assumed to hold,

$$\begin{aligned} \text{sw}_{L,M} \circ \phi^{\uparrow L} &= \phi \circ \text{sw}_{L,N} \quad , & \text{sw}_{L,M} \circ \theta &= \theta^{\uparrow M} \quad . \end{aligned}$$

$$\begin{array}{ccc} M^{L^A} & \xrightarrow{\text{sw}_{L,M}} & L^{M^A} \\ \phi \downarrow & & \downarrow \phi^{\uparrow L} \\ N^{L^A} & \xrightarrow{\text{sw}_{L,N}} & L^{N^A} \end{array} \qquad \begin{array}{ccc} M^{L^A} & \xrightarrow{\text{sw}_{L,M}} & L^{M^A} \\ & \searrow \theta^{\uparrow M} & \downarrow \theta \\ & & M^A \end{array}$$

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

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

13 Computations in functor blocks. III. Monad transformers

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

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

The next law is the identity law for brun :

$$\begin{aligned}
 \text{expect this to equal } \text{pu}_M : & \quad \text{pu}_{L \circ M} \circ \text{brun } \theta \\
 \text{definitions of brun and } \text{pu}_{L \circ M} : & \quad = \text{pu}_M \circ \text{pu}_L \circ \theta \\
 \text{identity law for } \theta : & \quad = \text{pu}_M .
 \end{aligned}$$

The last law to be shown is the composition law for $\text{brun } \theta$. Begin with its right-hand side since it is simpler,

$$\begin{aligned}
 & \quad (\text{brun } \theta)^{\uparrow M \uparrow L} \circ \text{brun } \theta \circ \text{ftn}_M \\
 \text{definition of brun :} & \quad = \theta^{\uparrow M \uparrow L} \circ \theta \circ \text{ftn}_M .
 \end{aligned}$$

We cannot simplify this expression any more, and yet it is still different from the left-hand side. So let us transform the left-hand side, hoping to obtain the same expression. In particular, we need to move ftn_M to the right and θ to the left:

$$\begin{aligned}
 & \quad \text{ftn}_{L \circ M} \circ \text{brun } \theta \\
 \text{definitions of } \text{ftn}_{L \circ M} \text{ and brun :} & \quad = \text{sw}_{L,M}^{\uparrow L} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \circ \theta \\
 \text{naturality of } \theta : & \quad = \text{sw}_{L,M}^{\uparrow L} \circ \text{ftn}_L \circ \theta \circ \text{ftn}_M \\
 \text{composition law for } \theta : & \quad = \text{sw}_{L,M}^{\uparrow L} \circ \theta^{\uparrow L} \circ \theta \circ \text{ftn}_M \\
 \text{functor composition :} & \quad = (\text{sw}_{L,M} \circ \theta)^{\uparrow L} \circ \theta \circ \text{ftn}_M \\
 \text{monadic naturality of } \text{sw}_{L,M} : & \quad = \theta^{\uparrow M \uparrow L} \circ \theta \circ \text{ftn}_M .
 \end{aligned}$$

13.3 Monad transformers via functor composition: General properties

We have transformed both sides of the law into the same expression.

The lifting laws for mrun are

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

Since $\text{mrun}(\phi) = \phi^{\uparrow L}$ in our case, these laws hold because they are the same as the functor laws of L .

Finally, we verify the non-degeneracy law for brun :

$$\text{expect to equal id} : \quad \text{lift} \circ \text{brun}(\theta)$$

$$\text{definitions of lift and brun} : \quad = \text{pu}_L \circ \theta$$

$$\text{identity law for } \theta : \quad = \text{id} \quad .$$

If the base monad is M , the runners have the type signatures

$$\text{mrun} \phi : L^{\bullet} \rightsquigarrow N^{\bullet} : L^{M^{\bullet}} \rightsquigarrow N^{M^{\bullet}} \quad , \quad \text{mrun} \phi = \phi \quad ;$$

$$\text{brun} \theta : M^{\bullet} \rightsquigarrow \bullet : L^{M^{\bullet}} \rightsquigarrow L^{\bullet} \quad , \quad \text{brun} \theta = \theta^{\uparrow L} \quad .$$

The laws of runners require that $\text{mrun} \phi$ and $\text{brun} \theta$ must be monadic morphisms, i.e. the identity and composition laws must hold for $\text{mrun} \phi$ and $\text{brun} \theta$:

$$\text{pu}_{L \circ M} \circ \text{mrun} \phi = \text{pu}_{N \circ M} \quad ,$$

$$\text{ftn}_{L \circ M} \circ \text{mrun} \phi = (\text{mrun} \phi)^{\uparrow M \uparrow L} \circ \text{mrun} \phi \circ \text{ftn}_{N \circ M} \quad ,$$

$$\text{pu}_{L \circ M} \circ \text{brun} \theta = \text{pu}_L \quad ,$$

$$\text{ftn}_{L \circ M} \circ \text{brun} \theta = (\text{brun} \theta)^{\uparrow M \uparrow L} \circ \text{brun} \theta \circ \text{ftn}_L \quad .$$

The monadic naturality laws for swap with respect to ϕ and χ are

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

$$\begin{array}{ccc} M^{L^A} & \xrightarrow{\text{sw}_{L,M}} & L^{M^A} \\ \phi^{\uparrow M} \downarrow & & \downarrow \phi \\ M^{N^A} & \xrightarrow{\text{sw}_{N,M}} & N^{M^A} \end{array} \quad \begin{array}{ccc} M^{L^A} & \xrightarrow{\text{sw}_{L,M}} & L^{M^A} \\ & \searrow \theta & \downarrow \theta^{\uparrow L} \\ & & L^A \end{array}$$

The first law to be proved is

$$\text{expect to equal pu}_{N \circ M} : \quad \text{pu}_{L \circ M} \circ \text{mrun} \phi$$

$$\text{definitions of mrun and pu}_{L \circ M} : \quad = \text{pu}_M \circ \text{pu}_L \circ \phi$$

$$\text{identity law for } \phi : \quad = \text{pu}_M \circ \text{pu}_N$$

$$\text{definition of pu}_{N \circ M} : \quad = \text{pu}_{N \circ M} \quad .$$

13 Computations in functor blocks. III. Monad transformers

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

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

The next law is the identity law for $\text{brun } (\theta)$:

$$\begin{aligned}
 \text{expect this to equal } \text{pu}_L : & \quad \text{pu}_{L \circ M} \circ \text{brun } \theta \\
 \text{definitions of brun and } \text{pu}_{L \circ M} : & \quad = \text{pu}_M \circ \underbrace{\text{pu}_L \circ \theta^{\uparrow L}} \\
 \text{naturality of } \text{pu}_L : & \quad = \underbrace{\text{pu}_M \circ \theta} \circ \text{pu}_L \\
 \text{identity law for } \theta : & \quad = \text{pu}_L \quad .
 \end{aligned}$$

The last law is the composition law for $\text{brun } (\theta)$. Begin with its right-hand side,

$$\begin{aligned}
 & \quad (\text{brun } \theta)^{\uparrow M \uparrow L} \circ \text{brun } \theta \circ \text{ftn}_L \\
 \text{definition of brun :} & \quad = \theta^{\uparrow L \uparrow M \uparrow L} \circ \theta^{\uparrow L} \circ \text{ftn}_L \\
 \text{functor composition :} & \quad = (\theta^{\uparrow L \uparrow M} \circ \theta)^{\uparrow L} \circ \text{ftn}_L \\
 \text{naturality of } \theta : & \quad = (\theta \circ \theta^{\uparrow L})^{\uparrow L} \circ \text{ftn}_L \quad .
 \end{aligned}$$

We now transform the left-hand side, hoping to obtain the same expression. We

13.3 Monad transformers via functor composition: General properties

need to move ftn_L to the right and θ to the left:

$$\begin{aligned}
 \text{expect to equal } \theta^{\uparrow L} \circ \theta^{\uparrow L \uparrow L} \circ \text{ftn}_L &: \text{ftn}_{L \circ M} \circ \text{brun } \theta \\
 \text{definitions of } \text{ftn}_{L \circ M} \text{ and } \text{brun} &: = \text{sw}_{L,M}^{\uparrow L} \circ \text{ftn}_L \circ \text{ftn}_M^{\uparrow L} \circ \theta^{\uparrow L} \\
 \text{composition law for } \theta &: = \text{sw}_{L,M}^{\uparrow L} \circ \text{ftn}_L \circ (\theta \circ \theta)^{\uparrow L} \\
 \text{naturality of } \text{ftn}_L &: = \text{sw}_{L,M}^{\uparrow L} \circ (\theta \circ \theta)^{\uparrow L \uparrow L} \circ \text{ftn}_L \\
 \text{functor composition} &: = (\text{sw}_{L,M} \circ \theta^{\uparrow L})^{\uparrow L} \circ \theta^{\uparrow L \uparrow L} \circ \text{ftn}_L \\
 \text{composition law for } \theta &: = \theta^{\uparrow L} \circ \theta^{\uparrow L \uparrow L} \circ \text{ftn}_L .
 \end{aligned}$$

The lifting laws for mrun are

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

Since $\text{mrun}(\phi) = \phi$ in our case, these laws are trivially satisfied.

Finally, the non-degeneracy law for brun :

$$\begin{aligned}
 \text{expect to equal id} &: \text{lift} \circ \text{brun}(\theta) \\
 \text{definitions of lift and brun} &: = \text{pu}_M^{\uparrow L} \circ \theta^{\uparrow L} \\
 \text{identity law for } \theta &: = \text{id}^{\uparrow L} = \text{id} \quad .
 \end{aligned}$$

13.3.7 Summary of results

The derivations in Section 13.3 enable us to state the following results:

Theorem 13.3.7.1 (composed-outside) For a base monad L and a foreign monad M , the functor composition $L \circ M$ is a lawful monad transformer if a `swap` function $\text{sw}_{L,M} : M \circ L \rightsquigarrow L \circ M$ exists, satisfying the conditions of Theorem 13.3.2.1 and the monadic naturality laws

$$\text{sw}_{L,M} \circ \phi^{\uparrow L} = \phi \circ \text{sw}_{L,N} \quad , \quad \text{sw}_{L,M} \circ \theta = \theta^{\uparrow M} \quad ,$$

with respect to arbitrary monadic morphisms $\phi : M \rightsquigarrow N$ and $\theta : L \rightsquigarrow \text{Id}$.

Theorem 13.3.7.2 (composed-inside) For a base monad M and a foreign monad L , the functor composition $L \circ M$ is a lawful monad transformer if a `swap` function $\text{sw}_{L,M} : M \circ L \rightsquigarrow L \circ M$ exists, satisfying the conditions of Theorem 13.3.2.1 and the monadic naturality laws

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

with respect to arbitrary monadic morphisms $\phi : L \rightsquigarrow N$ and $\theta : M \rightsquigarrow \text{Id}$.

These theorems enable us to check more easily whether a given base monad has a compositional monad transformer. It is quicker to check the laws of the `swap` function than to verify the monad transformer laws directly.

13.4 Composed-inside transformers: Linear monads

A monad M is **linear** if it is of the form $M^A \triangleq P + Q \times A$, where P and Q are fixed types, and Q is a monoid. (The polynomial $P + Q \times A$ is linear in its type parameter A .) Well-known examples of linear monads are `Option`, `Either`, and `Writer`. The general case $M^A \triangleq P + Q \times A$ represents a computation that can fail and at the same time produce a log message. So, M can be seen as a composition of `Either` and `Writer`.

A different (but also linear) monad is obtained from the composition of `Writer` and `Either`. The type constructor of this monad is $Q \times (P + A)$.

In general, composition of two linear monads $M_1^A \triangleq P_1 + Q_1 \times A$ and $M_2^A \triangleq P_2 + Q_2 \times A$ is again linear because

$$\begin{aligned} & P_1 + Q_1 \times (P_2 + Q_2 \times A) \\ \text{expand brackets :} &= \underline{P_1 + Q_1 \times P_2} + \underline{Q_1 \times Q_2 \times A} \\ \text{define new } P, Q : &= P + Q \times A \quad . \end{aligned}$$

Note that we need to define $Q \triangleq Q_1 \times Q_2$, and so Q is a monoid since, by assumption, Q_1 and Q_2 are monoids.

For a linear monad M and any foreign monad L , the functor composition $L \circ M$ is a monad. For example, the type constructor for the `OptionT` monad transformer can be defined as

```
type OptionT[L[_], A] = L[Option[A]]
```

The `Option` type constructor must be nested *inside* the foreign monad L . This is the case for all linear monads. Also, linear monads are the only known examples of monads whose transformers are composed inside the foreign monad.

13.4.1 Definitions of swap and flatten

To show that the monad transformer for the base monad $M^A \triangleq P + Q \times A$ is $T_M^{L,A} = L^{M^A}$, we will implement a suitable `swap` function having the type signature

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

for the base monad $M^A \triangleq P + Q \times A$ and an arbitrary foreign monad L . We will then prove that `swap` satisfies all the required laws stated in Theorem 13.3.7.2. This will guarantee that $T_M^{L,A} = L^{M^A}$ is a lawful monad transformer.

Expanding the definition of the type constructor M^\bullet , we can write the type signature of the `swap` function as

$$\text{sw}_{L,M} : P + Q \times L^A \Rightarrow L^{P+Q \times A} .$$

We can map P to L^P by applying pu_L . We can also map $Q \times L^A \Rightarrow L^{Q \times A}$ since L is a functor,

$$q \times l \Rightarrow (a \Rightarrow q \times a)^{\uparrow L} l .$$

It remains to unite these two functions. In the matrix notation, we write

$$\text{sw}_{L,M} = \left\{ \begin{array}{c} P \\ Q \times L^A \end{array} \left\| \begin{array}{c} (x^{\cdot P} \Rightarrow x + \mathbb{0}^{\cdot Q \times A}) ; \text{pu}_L \\ q \times l \Rightarrow (a^{\cdot A} \Rightarrow \mathbb{0}^{\cdot P} + q \times a)^{\uparrow L} l \end{array} \right. \right\} .$$

In Scala, the code is

```
type M[A, P, Q] = Either[P, (Q, A)]
def swap[L[_]: Monad, A, P, Q]: M[L[A]] => L[M[A]] = {
  case Left(p) => Monad[L].pure(Left(p))
  case Right((q, la)) => la.map(a => Right((q, a)))
}
```

Given this `swap` function, we define the `flatten` method (short notation ftn_T) for the transformed monad T by the standard formula

$$\text{ftn}_T = \text{sw}^{\uparrow L} ; \text{ftn}_L ; \text{ftn}_M^{\uparrow L} .$$

The pure method for T (short notation pu_T) is $\text{pu}_T = \text{pu}_M ; \text{pu}_L$. In Scala:

```
def pure[L[_]: Monad, A, P, Q: Monoid](x: A): L[M[A]] =
  Monad[L].pure(Right((Monoid[Q].empty, x)))
def flatten[L[_]: Monad, A, P, Q: Monoid](tt: L[M[L[M[A]]]]): L[M[A]] =
  tt.map(swap).flatten.map(_.flatten) // Assuming suitable implicits in scope.
```

13 Computations in functor blocks. III. Monad transformers

Now we will verify that the laws of `swap` hold. We will need to use the code for the methods `fmapM`, `ftnM`, and `puM` of the monad M :

$$\begin{aligned} \text{fmap}_M f^{A \Rightarrow B} &= f^{\uparrow M} = \left[\begin{array}{c|c} & \begin{array}{c} P \quad Q \times B \end{array} \\ \hline \begin{array}{c} P \\ Q \times A \end{array} & \begin{array}{c} \text{id} \quad \mathbb{0} \\ \mathbb{0} \quad q \times a \Rightarrow q \times f(a) \end{array} \end{array} \right] , \\ \text{pu}_M a^{A \Rightarrow B} &= \mathbb{0}^P + q_0 \times a \quad , \quad \text{pu}_M = \left[\begin{array}{c|c} & \begin{array}{c} P \quad Q \times A \end{array} \\ \hline A & \begin{array}{c} \mathbb{0} \quad a \Rightarrow q_0 \times a \end{array} \end{array} \right] , \\ \text{ftn}_M^{M^A \Rightarrow M^A} &= \left[\begin{array}{c|c} & \begin{array}{c} P \quad Q \times A \end{array} \\ \hline \begin{array}{c} P \\ Q \times P \\ Q \times Q \times A \end{array} & \begin{array}{c} \text{id} \quad \mathbb{0} \\ q \times p \Rightarrow p \quad \mathbb{0} \\ \mathbb{0} \quad q_1 \times q_2 \times a \Rightarrow (q_1 \oplus q_2) \times a \end{array} \end{array} \right] . \end{aligned}$$

13.4.2 Laws of `swap`

We do not need to verify naturality since `swap` is defined as a fully parametric function.

The inner-identity law We need to show that $\text{pu}_L^{\uparrow M} \circ \text{sw} = \text{pu}_L$:

$$\begin{aligned} \text{pu}_L^{\uparrow M} \circ \text{sw} &= \left[\begin{array}{c|c} \text{id} \quad \mathbb{0} \\ \hline \mathbb{0} \quad q \times a \Rightarrow q \times \text{pu}_L a \end{array} \right] \circ \left[\begin{array}{c|c} (x^{P \Rightarrow x + \mathbb{0}^{Q \times A}}) \circ \text{pu}_L \\ \hline q \times l \Rightarrow (x^{A \Rightarrow \mathbb{0}^P + q \times x})^{\uparrow L} l \end{array} \right] \\ \text{composition :} &= \left[\begin{array}{c|c} (x^{P \Rightarrow x + \mathbb{0}^{Q \times A}}) \circ \text{pu}_L \\ \hline q \times a \Rightarrow (a^{A \Rightarrow \mathbb{0}^P + q \times a})^{\uparrow L} (\text{pu}_L a) \end{array} \right] \\ \text{pu}_L \text{'s naturality :} &= \left[\begin{array}{c|c} P \\ \hline Q \times A \end{array} \left[\begin{array}{c} x^{P \Rightarrow \text{pu}_L(x + \mathbb{0}^{Q \times A})} \\ q \times a \Rightarrow \text{pu}_L(\mathbb{0}^P + q \times a) \end{array} \right] \right] \\ \text{matrix notation :} &= \text{pu}_L \quad . \end{aligned}$$

13.4 Composed-inside transformers: Linear monads

The outer-identity law We need to show that $\text{pu}_M \circ \text{sw} = \text{pu}_M^{\uparrow L}$:

$$\text{pu}_M \circ \text{sw} = \left\| \begin{array}{c} \mathbb{0} \\ l : L^A \Rightarrow q_0 \times l \end{array} \right\| \circ \left\| \begin{array}{c} (x : P \Rightarrow x + \mathbb{0} : Q \times A) \circ \text{pu}_L \\ q \times l \Rightarrow (x : A \Rightarrow \mathbb{0} : P + q \times x)^{\uparrow L} l \end{array} \right\|$$

composition : $= l : L^A \Rightarrow (x : A \Rightarrow \mathbb{0} : P + q_0 \times x)^{\uparrow L} l$

definition of pu_M : $= l \Rightarrow \text{pu}_M^{\uparrow L} l = \text{pu}_M$.

The inner-interchange law Show that $\text{ftn}_L^{\uparrow M} \circ \text{sw} = \text{sw} \circ \text{sw}^{\uparrow L} \circ \text{ftn}_L$:

$$\begin{aligned} \text{ftn}_L^{\uparrow M} \circ \text{sw} &= \left\| \begin{array}{c} \text{id} \\ \mathbb{0} \\ q \times l \Rightarrow q \times \text{ftn}_L l \end{array} \right\| \circ \left\| \begin{array}{c} (x : P \Rightarrow x + \mathbb{0} : Q \times A) \circ \text{pu}_L \\ q \times l \Rightarrow (a \Rightarrow \mathbb{0} : P + q \times a)^{\uparrow L} l \end{array} \right\| \\ &= \left\| \begin{array}{c} (x : P \Rightarrow x + \mathbb{0} : Q \times A) \circ \text{pu}_L \\ q \times l \Rightarrow (a \Rightarrow \mathbb{0} : P + q \times a)^{\uparrow L} (\text{ftn}_L l) \end{array} \right\| , \quad (13.18) \\ \text{sw} \circ \text{sw}^{\uparrow L} \circ \text{ftn}_L &= \left\| \begin{array}{c} (x : P \Rightarrow x + \mathbb{0} : Q \times A) \circ \text{pu}_L \\ q \times l \Rightarrow (a \Rightarrow \mathbb{0} : P + q \times a)^{\uparrow L} l \end{array} \right\| \circ \text{sw}^{\uparrow L} \circ \text{ftn}_L \\ &= \left\| \begin{array}{c} (x : P \Rightarrow x + \mathbb{0} : Q \times A) \circ \text{pu}_L \circ \text{sw}^{\uparrow L} \circ \text{ftn}_L \\ (q \times l \Rightarrow (a \Rightarrow \mathbb{0} : P + q \times a)^{\uparrow L} l) \circ \text{sw}^{\uparrow L} \circ \text{ftn}_L \end{array} \right\| . \end{aligned}$$

It is quicker to simplify each expression in the last column separately and then to compare with the column in Eq. (13.18). Simplify the upper expression:

$$\begin{aligned} &(x : P \Rightarrow x + \mathbb{0} : Q \times A) \circ \text{pu}_L \circ \text{sw}^{\uparrow L} \circ \text{ftn}_L \\ \text{naturality of } \text{pu}_L : &= (x : P \Rightarrow x + \mathbb{0} : Q \times A) \circ \text{sw} \circ \text{pu}_L \circ \text{ftn}_L \\ \text{identity law of } L : &= (x : P \Rightarrow x + \mathbb{0} : Q \times A) \circ \text{sw} \\ \text{definition of } \text{sw} : &= (x : P \Rightarrow x + \mathbb{0} : Q \times A) \circ \text{pu}_L . \end{aligned}$$

This equals the upper expression in Eq. (13.18). Simplify the lower expression;

$$\begin{aligned} &(q \times l \Rightarrow (a \Rightarrow \mathbb{0} : P + q \times a)^{\uparrow L} l) \circ \text{sw}^{\uparrow L} \circ \text{ftn}_L \\ \text{definition of } \triangleright : &= q \times l \Rightarrow l \triangleright (a \Rightarrow \mathbb{0} : P + q \times a)^{\uparrow L} \circ \text{sw}^{\uparrow L} \circ \text{ftn}_L . \quad (13.19) \end{aligned}$$

13 Computations in functor blocks. III. Monad transformers

Simplify the expression $(a \Rightarrow \mathbb{0}^P + q \times a)^{\uparrow L} \mathbin{\circ} \text{sw}^{\uparrow L}$ separately:

$$\begin{aligned}
 & (a \Rightarrow \mathbb{0}^P + q \times a)^{\uparrow L} \mathbin{\circ} \text{sw} \\
 \text{composition :} &= a \Rightarrow \underline{\text{sw} (\mathbb{0}^P + q \times a)} \\
 \text{definition of sw :} &= \underline{a \Rightarrow (x \Rightarrow \mathbb{0}^P + q \times x)^{\uparrow L} a} \\
 \text{omit argument :} &= (x \Rightarrow \mathbb{0}^P + q \times x)^{\uparrow L} . \tag{13.20}
 \end{aligned}$$

Then we continue simplifying Eq. (13.19):

$$\begin{aligned}
 & q \times l \Rightarrow l \triangleright \underline{(a \Rightarrow \mathbb{0}^P + q \times a)^{\uparrow L} \mathbin{\circ} \text{sw}^{\uparrow L} \mathbin{\circ} \text{ftn}_L} \\
 \text{use Eq. (13.20) :} &= q \times l \Rightarrow l \triangleright \underline{(x \Rightarrow \mathbb{0}^P + q \times x)^{\uparrow L \uparrow L} \mathbin{\circ} \text{ftn}_L} \\
 \text{naturality of ftn}_L : &= q \times l \Rightarrow \underline{l \triangleright \text{ftn}_L} \mathbin{\circ} (x \Rightarrow \mathbb{0}^P + q \times x)^{\uparrow L} \\
 \text{definition of } \triangleright : &= q \times l \Rightarrow (x \Rightarrow \mathbb{0}^P + q \times x)^{\uparrow L} (\text{ftn}_L l) .
 \end{aligned}$$

This equals the lower expression in Eq. (13.18) after renaming x to a .

The outer-interchange law Show that $\text{ftn}_M \mathbin{\circ} \text{sw} = \text{sw}^{\uparrow M} \mathbin{\circ} \text{sw} \mathbin{\circ} \text{ftn}_M^{\uparrow L}$. The left-hand side is written using the matrices for ftn_M and sw :

$$\begin{aligned}
 & \text{ftn}_M \mathbin{\circ} \text{sw} \\
 = & \left\| \begin{array}{cc} \text{id} & \mathbb{0} \\ q \times p \Rightarrow p & \mathbb{0} \\ \mathbb{0} & q_1 \times q_2 \times a \Rightarrow (q_1 \oplus q_2) \times a \end{array} \right\| \mathbin{\circ} \left\| \begin{array}{c} (x^P \Rightarrow x + \mathbb{0}) \mathbin{\circ} \text{pu}_L \\ q \times l \Rightarrow (x \Rightarrow \mathbb{0} + q \times x)^{\uparrow L} l \end{array} \right\| \\
 = & \left\| \begin{array}{c} (x^P \Rightarrow x + \mathbb{0}) \mathbin{\circ} \text{pu}_L \\ (q \times p \Rightarrow p + \mathbb{0}) \mathbin{\circ} \text{pu}_L \\ q_1 \times q_2 \times a \Rightarrow (x \Rightarrow \mathbb{0} + (q_1 \oplus q_2) \times x)^{\uparrow L} a \end{array} \right\| . \tag{13.21}
 \end{aligned}$$

We cannot simplify this any more, so we hope to transform the right-hand side, $\text{sw}^{\uparrow M} \mathbin{\circ} \text{sw} \mathbin{\circ} \text{ftn}_M^{\uparrow L}$, to the same column expression. Begin by writing the matrix for

13.4 Composed-inside transformers: Linear monads

$\text{sw}^{\uparrow M}$, expanding the rows for the input type $M^{M^{L^A}}$:

$$\begin{aligned} \text{sw}^{\uparrow M} &= \left\| \begin{array}{c} P \\ Q \times P \\ Q \times Q \times L^A \end{array} \right\| \left\| \begin{array}{cc} \text{id} & \mathbb{0} \\ \mathbb{0} & q \times p \Rightarrow q \times \text{sw} (p + \mathbb{0}) \\ \mathbb{0} & q_1 \times q_2 \times l \Rightarrow q_1 \times \text{sw} (\mathbb{0} + q_2 \times l) \end{array} \right\| \\ &= \left\| \begin{array}{c} P \\ Q \times P \\ Q \times Q \times L^A \end{array} \right\| \left\| \begin{array}{cc} \text{id} & \mathbb{0} \\ \mathbb{0} & q \times p \Rightarrow q \times \text{pu}_L (p + \mathbb{0}) \\ \mathbb{0} & q_1 \times q_2 \times l \Rightarrow q_1 \times (x \Rightarrow \mathbb{0} + q_2 \times x)^{\uparrow L} l \end{array} \right\|. \end{aligned}$$

Then compute the composition $\text{sw}^{\uparrow M} \circ \text{sw}$ as

$$\begin{aligned} &\text{sw}^{\uparrow M} \circ \text{sw} \\ &= \left\| \begin{array}{cc} \text{id} & \mathbb{0} \\ \mathbb{0} & q \times p \Rightarrow q \times \text{pu}_L (p + \mathbb{0}) \\ \mathbb{0} & q_1 \times q_2 \times l \Rightarrow q_1 \times (x \Rightarrow \mathbb{0} + q_2 \times x)^{\uparrow L} l \end{array} \right\| \circ \left\| \begin{array}{c} (x^{:P} \Rightarrow x + \mathbb{0}) \circ \text{pu}_L \\ q \times l \Rightarrow (x \Rightarrow \mathbb{0} + q \times x)^{\uparrow L} l \end{array} \right\| \\ &= \left\| \begin{array}{c} (x^{:P} \Rightarrow x + \mathbb{0}) \circ \text{pu}_L \\ q \times p \Rightarrow (x^{:M^A} \Rightarrow \mathbb{0}^{:P} + q \times x)^{\uparrow L} (\text{pu}_L (p + \mathbb{0})) \\ q_1 \times q_2 \times l \Rightarrow (x^{:M^A} \Rightarrow \mathbb{0}^{:P} + q_1 \times x)^{\uparrow L} (x \Rightarrow \mathbb{0} + q_2 \times x)^{\uparrow L} l \end{array} \right\| \\ &= \left\| \begin{array}{c} (x^{:P} \Rightarrow x + \mathbb{0}) \circ \text{pu}_L \\ q \times p \Rightarrow \text{pu}_L (\mathbb{0}^{:P} + q \times (p + \mathbb{0})) \\ q_1 \times q_2 \times l \Rightarrow (x^{:M^A} \Rightarrow \mathbb{0} + q_1 \times (\mathbb{0} + q_2 \times x))^{\uparrow L} l \end{array} \right\|. \end{aligned}$$

13 Computations in functor blocks. III. Monad transformers

Now we need to post-compose $\text{ftn}_M^{\uparrow L}$ with this column:

$$\begin{aligned}
 \text{sw}^{\uparrow M} \circ \text{sw} \circ \text{ftn}_M^{\uparrow L} &= \left\| \begin{array}{l} (x^{:P} \Rightarrow x + \mathbb{0}) \circ \text{pu}_L \circ \text{ftn}_M^{\uparrow L} \\ (q \times p \Rightarrow \mathbb{0}^{:P} + q \times (p + \mathbb{0})) \circ \text{pu}_L \circ \text{ftn}_M^{\uparrow L} \\ q_1 \times q_2 \times l \Rightarrow l \triangleright (x^{:M^A} \Rightarrow \mathbb{0} + q_1 \times (\mathbb{0} + q_2 \times x))^{\uparrow L} \circ \text{ftn}_M^{\uparrow L} \end{array} \right\| \\
 \text{pu}_L \text{'s naturality :} &= \left\| \begin{array}{l} (x^{:P} \Rightarrow x + \mathbb{0}) \circ \text{ftn}_M \circ \text{pu}_L \\ (q \times p \Rightarrow \mathbb{0}^{:P} + q \times (p + \mathbb{0})) \circ \text{ftn}_M \circ \text{pu}_L \\ q_1 \times q_2 \times l \Rightarrow l \triangleright (x^{:M^A} \Rightarrow \text{ftn}_M(\mathbb{0} + q_1 \times (\mathbb{0} + q_2 \times x)))^{\uparrow L} \end{array} \right\| \\
 \text{compute } \text{ftn}_M(\dots) : &= \left\| \begin{array}{l} (x^{:P} \Rightarrow \text{ftn}_M(x + \mathbb{0})) \circ \text{pu}_L \\ (q \times p \Rightarrow \text{ftn}_M(\mathbb{0}^{:P} + q \times (p + \mathbb{0}))) \circ \text{pu}_L \\ q_1 \times q_2 \times l \Rightarrow l \triangleright (x^{:M^A} \Rightarrow \mathbb{0} + (q_1 \oplus q_2) \times x)^{\uparrow L} \end{array} \right\| \\
 \text{compute } \text{ftn}_M(\dots) : &= \left\| \begin{array}{l} (x^{:P} \Rightarrow x + \mathbb{0}) \circ \text{pu}_L \\ (q \times p \Rightarrow p + \mathbb{0}) \circ \text{pu}_L \\ q_1 \times q_2 \times l \Rightarrow (x^{:M^A} \Rightarrow \mathbb{0} + (q_1 \oplus q_2) \times x)^{\uparrow L} \end{array} \right\| .
 \end{aligned}$$

After renaming l to a , this is the same as the column in Eq. (13.21).

Monadic naturality laws Verify the laws of Theorem 13.3.7.2,

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

for arbitrary monadic morphisms $\phi : L \rightsquigarrow N$ and $\theta : M \rightsquigarrow \text{Id}$.

Begin with the left-hand side of the first law,

$$\begin{aligned}
 &\text{sw}_{L,M} \circ \phi \\
 \text{definition of } \text{sw}_{L,M} : &= \left\| \begin{array}{l} (x^{:P} \Rightarrow x + \mathbb{0}^{:Q \times A}) \circ \text{pu}_L \\ q \times l \Rightarrow (a \Rightarrow \mathbb{0}^{:P} + q \times a)^{\uparrow L} \end{array} \right\| \circ \phi \\
 \text{compose with } \phi : &= \left\| \begin{array}{l} (x^{:P} \Rightarrow x + \mathbb{0}^{:Q \times A}) \circ \text{pu}_L \circ \phi \\ q \times l \Rightarrow l \triangleright (a \Rightarrow \mathbb{0}^{:P} + q \times a)^{\uparrow L} \circ \phi \end{array} \right\| \\
 \text{naturality of } \phi : &= \left\| \begin{array}{l} (x^{:P} \Rightarrow x + \mathbb{0}^{:Q \times A}) \circ \text{pu}_N \\ q \times l \Rightarrow l \triangleright \phi \circ (a \Rightarrow \mathbb{0}^{:P} + q \times a)^{\uparrow N} \end{array} \right\| .
 \end{aligned}$$

13.4 Composed-inside transformers: Linear monads

The right-hand side is

$$\begin{aligned}
 & \phi^{\uparrow M} \circ \text{sw}_{N,M} \\
 &= \left\| \begin{array}{cc} \text{id} & \mathbb{0} \\ \mathbb{0} & q \times l \Rightarrow q \times \phi(l) \end{array} \right\| \circ \left\| \begin{array}{c} (x^{:P} \Rightarrow x + \mathbb{0}^{:Q \times A}}) \circ \text{pu}_N \\ q \times n \Rightarrow n \triangleright (a \Rightarrow \mathbb{0}^{:P} + q \times a)^{\uparrow N} \end{array} \right\| \\
 \text{composition : } &= \left\| \begin{array}{c} (x^{:P} \Rightarrow x + \mathbb{0}^{:Q \times A}}) \circ \text{pu}_N \\ q \times l \Rightarrow n \triangleright \phi \circ (a \Rightarrow \mathbb{0}^{:P} + q \times a)^{\uparrow N} \end{array} \right\|.
 \end{aligned}$$

Both sides of the first law are now shown to be equal.

The left-hand side of the second law is

$$\begin{aligned}
 & \text{sw}_{L,M} \circ \theta^{\uparrow L} \\
 \text{compose with } \theta^{\uparrow L} : &= \left\| \begin{array}{c} (x^{:P} \Rightarrow x + \mathbb{0}^{:Q \times A}}) \circ \underline{\text{pu}_L \circ \theta^{\uparrow L}} \\ q \times l \Rightarrow l \triangleright (a \Rightarrow \mathbb{0}^{:P} + q \times a)^{\uparrow L} \circ \theta^{\uparrow L} \end{array} \right\| \\
 \text{naturality of } \text{pu}_L : &= \left\| \begin{array}{c} (x^{:P} \Rightarrow x + \mathbb{0}^{:Q \times A}}) \circ \theta \circ \text{pu}_L \\ q \times l \Rightarrow l \triangleright (a \Rightarrow \theta(\mathbb{0}^{:P} + q \times a))^{\uparrow L} \end{array} \right\|. \quad (13.22)
 \end{aligned}$$

We expect this to equal the right-hand side, which we write as

$$\begin{aligned}
 & m^{:M^{L^A}} \Rightarrow \theta(m) \\
 \text{matrix notation : } &= \left\| \begin{array}{c} x^{:P} \Rightarrow \theta(x + \mathbb{0}^{:Q \times L^A}) \\ q \times l \Rightarrow \theta(\mathbb{0}^{:P} + q \times l) \end{array} \right\|. \quad (13.23)
 \end{aligned}$$

Now consider each line in Eq. (13.22) separately. The upper line can be transformed as

$$\begin{aligned}
 & (x^{:P} \Rightarrow x + \mathbb{0}^{:Q \times A}}) \circ \underline{\theta \circ \text{pu}_L} \\
 \text{naturality of } \theta : &= (x^{:P} \Rightarrow x + \mathbb{0}^{:Q \times A}}) \circ \underline{\text{pu}_L^{\uparrow M} \circ \theta} \\
 \text{definition of } \uparrow^M : &= x^{:P} \Rightarrow \left\| \begin{array}{cc} x & \mathbb{0} \\ \mathbb{0} & q \times l \Rightarrow q \times \text{pu}_L l \end{array} \right\| \triangleright \left\| \begin{array}{c} \text{id} & \mathbb{0} \\ \mathbb{0} & q \times l \Rightarrow q \times \text{pu}_L l \end{array} \right\| \circ \theta \\
 \text{matrix product : } &= x^{:P} \Rightarrow (x + \mathbb{0}^{:Q \times L^A}) \triangleright \theta.
 \end{aligned}$$

13 Computations in functor blocks. III. Monad transformers

This is now equal to the upper line of Eq. (13.23).

To proceed with the proof for the lower line of Eq. (13.22), we need to evaluate the monadic morphism $\theta : M^A \Rightarrow A$ on a specific value of type M^A of the form $0 + q \times a$. We note that the value $\theta(0 + q \times a)$ is of type A and must be computed in the same way for all types A , because θ is a natural transformation. It seems clear that the result cannot depend on the value q^Q since Q is a type not related to A . In other words, we expect that $\theta(0 + q \times a) = a$ as a consequence of naturality of θ . To prove this formally, we use the trick of starting with a unit type, 1 , and mapping it to a within the naturality law. For any values q^Q, a^A , we define

$$\begin{aligned} m : P + Q \times A &= M^A, & m &\triangleq \mathbb{0}^P + q \times a, \\ m_1 : P + Q \times \mathbb{1} &= M^{\mathbb{1}}, & m_1 &\triangleq \mathbb{0}^P + q \times 1. \end{aligned}$$

We can compute m from m_1 if replace 1 by a under the functor M . To write this as a formula, define the function $f : \mathbb{1} \Rightarrow A$ as $f \triangleq (_ \Rightarrow a)$ using the fixed value a . Then we have $m = f^{\uparrow M} m_1$. Now we apply both sides of the naturality law $f^{\uparrow M} \circ \theta = \theta \circ f$ to the value m_1 :

$$m_1 \triangleright f^{\uparrow M} \circ \theta = m_1 \triangleright \theta \circ f.$$

Simplify the left-hand side to

$$m_1 \triangleright f^{\uparrow M} \circ \theta = \theta(f^{\uparrow M} m_1) = \theta(m) = \theta(\mathbb{0}^P + q \times a).$$

Simplify the right-hand side to

$$m_1 \triangleright \theta \circ f = f(\theta(m_1)) = a,$$

since f always returns a . Therefore

$$\theta(\mathbb{0}^P + q \times a) = a. \quad (13.24)$$

We can now compute the second line in Eq. (13.22) as

$$\begin{aligned} q \times l &\Rightarrow l \triangleright (a \Rightarrow \theta(\mathbb{0}^P + q \times a))^{\uparrow L} \\ \text{use Eq. (13.24) :} &= q \times l \Rightarrow l \triangleright \underline{(a \Rightarrow a)}^{\uparrow L} \\ \text{identity law :} &= q \times l \Rightarrow l. \end{aligned}$$

The second line in Eq. (13.23) is the same function, $q \times l \Rightarrow l$.

This concludes the proof of the swap laws for linear monads. It follows that linear monads have monad transformers of composed-inside kind.

13.4.3 Composition of transformers for linear monads

We have just shown that the “Either/Writer” monad $M^A \triangleq P + Q \times A$ has the `swap` function that satisfies the laws necessary for a composed-inside transformer. The other type of linear monad is the “Writer/Either” monad $W^A \triangleq Q \times (P + A)$. We need to show separately that the monad W has a lawful `swap` function? Actually, this follows from the stacking property of monad transformers (see Sections 13.2.6–13.2.7). The monad W is a functor composition of the `Writer` monad $Q \times A$ with the `Either` monad $P + A$, which is the same as applying the `Either` monad’s transformer to the `Writer` monad. Because of the transformer stacking property, the monad transformer of W works as composed-inside.

We can show in general that the functor composition of any two linear monads has a composed-inside transformer. Suppose M_1 and M_2 are linear monads, so their transformers are of the composed-inside kind:

$$T_{M_1}^N = N \circ M_1 \quad , \quad T_{M_2}^N = N \circ M_2 \quad .$$

The functor composition of M_1 and M_2 can be seen as a monad stack,

$$M_1 \circ M_2 = T_{M_2}^{M_1} \quad .$$

What is the transformer for the monad $M_1 \circ M_2$? For any foreign monad N , we have the transformer stack

$$T_{M_2}^{T_{M_1}^N} = T_{M_2}^{N \circ M_1} = N \circ M_1 \circ M_2 \quad .$$

Since this is a transformer stack, it is a lawful monad transformer, as we have seen in Section 13.2.6. So, this is the monad transformer for $M_1 \circ M_2$, and it is of the composed-inside kind.

13.5 Composed-outside transformers: Rigid monads

Section 13.4 shows that the composed-inside monad transformers are available only for a limited subset of all monads, namely the monads of the form $M^A = P + Q \times A$ and $M^A = Q \times (P + A)$, which I call “linear”. It turns out that the composed-outside transformers are available for a significantly wider range of monads. I call those monads “rigid”, because one of their general properties is having a single “shape” (Section ???). (There does not seem to be an already accepted name for these monads.)

Definition: A monad R is **rigid** if it has a composed-outside monad transformer, $T_R^M = R \circ M$, where M is a foreign monad.

This definition does not explain what monads are rigid or how to recognize a non-rigid monad. Two simple examples of rigid monads are the `Reader` monad and the `Search` monad,

$$\begin{aligned} \text{(the Reader monad) : } & R^A \triangleq Z \Rightarrow A \quad , \\ \text{(the Search monad) : } & S^A \triangleq (A \Rightarrow Z) \Rightarrow A \quad , \end{aligned}$$

where Z is a fixed type. These monads have composed-outside transformers:

$$\begin{aligned} \text{(the ReaderT transformer) : } & T_R^{M,A} \triangleq Z \Rightarrow M^A \quad , \\ \text{(the SearchT transformer) : } & T_S^{M,A} \triangleq (M^A \Rightarrow Z) \Rightarrow M^A \quad . \end{aligned}$$

To build intuition for rigid monads, we will look at some general constructions that create new rigid monads or combine existing rigid monads into new ones. In this section, we will prove that the following four constructions give rigid monads:

1. Choice: $C^A \triangleq H^A \Rightarrow A$ is a rigid monad if H is any contrafunctor.
2. Composition: $P \circ R$ is a rigid monad if P and R are rigid monads.
3. Product: $P^A \times R^A$ is a rigid monad if P and R are rigid monads.
4. Selector: $S^A \triangleq F^{A \Rightarrow R^Q} \Rightarrow R^A$ is a rigid monad for any rigid monad R , any functor F , and any fixed type Q .

I do not know whether these four constructions are the only possible ways of creating new rigid monads. Below I will also mention other open questions I have about rigid monads.

13.5.1 Rigid monad construction 1: choice

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

This monad chooses a value of type A given a contrafunctor H that may *consume* values of type A (and presumably could 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.

13.5 Composed-outside transformers: Rigid monads

Different choices of the contrafunctor H give specific examples of rigid monads, such as $R^A \triangleq \mathbb{1}$ (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, if possible, satisfies the predicate. If no such value can be found, *some* value of type A is still returned.

A closely related monad is the search-with-failure monad, $R^A \triangleq (A \Rightarrow \text{Bool}) \Rightarrow \mathbb{1} + A$. This (non-rigid) monad will return an empty value $1 + \mathbb{0}^A$ if no value satisfying the predicate was found. There is a natural transformation from the search monad to the search-with-failure monad, implemented by checking whether the value returned by the search monad does actually satisfies 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 will first give a self-contained proof that T is a monad. To verify the laws of monad transformers for T , we will derive the `swap` function and verify its laws.

Statement 13.5.1.1 $T^\bullet \triangleq R^{M^\bullet} \triangleq H^{M^\bullet} \Rightarrow M^\bullet$ is a monad if M is any monad and H is any contrafunctor. (If we set $M^A \triangleq A$, this also proves 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.

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

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

13 Computations in functor blocks. III. Monad transformers

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 pure_T and $\tilde{\otimes}_T$ the flipped Kleisli operations:

$$\begin{aligned} \text{pure}_T : H^{M^A} &\Rightarrow A \Rightarrow M^A \\ f : H^{M^B} \Rightarrow A \Rightarrow M^B \quad \tilde{\otimes}_T g : H^{M^C} \Rightarrow B \Rightarrow M^C &: H^{M^C} \Rightarrow A \Rightarrow M^C \quad . \end{aligned}$$

To define the operations pure_T and $\tilde{\otimes}_T$, 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

$$\begin{aligned} \text{pure}_T q &= q \Rightarrow \text{pu}_M \quad (\text{the argument } q \text{ is unused}), \\ f \tilde{\otimes}_T g &= q \Rightarrow (f p) \diamond_M (g q) \quad \text{where} \\ p : H^{M^B} &= (\text{flm}_M (g q))^{\downarrow H} q \quad . \end{aligned}$$

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 $p : H^{M^B}$.

Written as a single expression, the definition of $\tilde{\otimes}_T$ is

$$f \tilde{\otimes}_T g = q \Rightarrow f \left((\text{flm}_M (g q))^{\downarrow H} q \right) \diamond_M (g q) \quad . \quad (13.25)$$

Checking the left identity law:

$$\begin{aligned} &\text{pure}_T \tilde{\otimes}_T g \\ \text{definition of } \tilde{\otimes}_T : &= q \Rightarrow \text{pure}_T \left((\text{flm}_M (g q))^{\downarrow H} q \right) \diamond_M (g q) \\ \text{definition of } \text{pure}_T : &= q \Rightarrow \underline{\text{pu}_M \diamond_M g q} \\ \text{left identity law for } M : &= q \Rightarrow g q \\ \text{function expansion :} &= g \end{aligned}$$

13.5 Composed-outside transformers: Rigid monads

Checking the right identity law:

$$\begin{aligned}
 & f \tilde{\circ}_T \tilde{\text{pu}}_T \\
 \text{definition of } \tilde{\circ}_T : & = q \Rightarrow f \left(\left(\text{flm}_M (\tilde{\text{pu}}_T q) \right)^{\downarrow H} q \right) \diamond_M (\tilde{\text{pu}}_T q) \\
 \text{definition of } \tilde{\text{pu}}_T : & = q \Rightarrow f \left(\left(\text{flm}_M (\text{pu}_M) \right)^{\downarrow H} q \right) \diamond_M \underline{\text{pu}_M} \\
 \text{right identity law for } M : & = q \Rightarrow f \left(\underline{(\text{id})^{\downarrow H} q} \right) \\
 \text{identity law for } H : & = q \Rightarrow f q \\
 \text{function expansion :} & = f
 \end{aligned}$$

Checking the associativity law: $(f \tilde{\circ}_T g) \tilde{\circ}_T h$ must equal $f \tilde{\circ}_T (g \tilde{\circ}_T h)$. We have

$$\begin{aligned}
 & (f \tilde{\circ}_T g) \tilde{\circ}_T h \\
 & = (s \Rightarrow f \left((\text{flm}_M (g s))^{\downarrow H} s \right) \diamond_M (g s)) \tilde{\circ}_T h \\
 & = q \Rightarrow f \left((\text{flm}_M (g r))^{\downarrow H} r \right) \diamond_M (g r) \diamond_M (h q) \quad \text{where} \\
 & \quad r \triangleq (\text{flm}_M (h q))^{\downarrow H} q \quad ;
 \end{aligned}$$

while

$$\begin{aligned}
 & f \tilde{\circ}_T (g \tilde{\circ}_T h) \\
 & = f \tilde{\circ}_T (q \Rightarrow g \left((\text{flm}_M (h q))^{\downarrow H} q \right) \diamond_M (h q)) \\
 & = q \Rightarrow f \left((\text{flm}_M k)^{\downarrow H} q \right) \diamond_M u \quad \text{where} \\
 & \quad r \triangleq (\text{flm}_M (h q))^{\downarrow H} q \quad \text{and} \\
 & \quad u \triangleq (g r) \diamond_M (h q) \quad .
 \end{aligned}$$

It remains to show that the following two expressions are equal,

$$\begin{aligned}
 & f \left((\text{flm}_M (g r))^{\downarrow H} r \right) \diamond_M (g r) \diamond_M (h q) \quad \text{and} \\
 & f \left((\text{flm}_M ((g r) \diamond_M (h q)))^{\downarrow H} q \right) \diamond_M (g r) \diamond_M (h q), \quad \text{where} \\
 & \quad r \triangleq (\text{flm}_M (h q))^{\downarrow H} q \quad .
 \end{aligned}$$

These two expressions differ only by the following sub-expressions,

$$(\text{flm}_M (g r))^{\downarrow H} r$$

13 Computations in functor blocks. III. Monad transformers

and

$$(\text{flm}_M ((g\ r) \diamond_M (h\ q)))^{\downarrow H} q \quad ,$$

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

$$\begin{aligned} & (\text{flm}_M (g\ r))^{\downarrow H} (\text{flm}_M (h\ q))^{\downarrow H} q \quad \text{and} \\ & (\text{flm}_M ((g\ r) \diamond_M (h\ q)))^{\downarrow H} q \quad . \end{aligned}$$

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:

$$\begin{aligned} & \underline{(\text{flm}_M (g\ r))^{\downarrow H} (\text{flm}_M (h\ q))^{\downarrow H} q} \\ \text{composition law for } H : & = (\text{flm}_M (g\ r) \mathbin{\text{;}} \text{flm}_M (h\ q))^{\downarrow H} q \\ \text{associativity law for } M : & = \underline{(\text{flm}_M ((g\ r) \mathbin{\text{;}} \text{flm}_M (h\ q)))^{\downarrow H} q} \\ \text{definition of } \diamond_M \text{ via } \text{flm}_M : & = (\text{flm}_M ((g\ r) \diamond_M (h\ q)))^{\downarrow H} q \quad . \end{aligned}$$

This proves the associativity law for $\tilde{\delta}_T$.

Statement 13.5.1.2 The monad instance for T defined in Statement 3 can be defined equivalently as

$$\begin{aligned} \text{pu}_T(a^{:A}) : H^{M^A} &\Rightarrow M^A \quad , \\ \text{pu}_T(a) &\triangleq (- \Rightarrow \text{pu}_M a) \quad ; \\ \text{flm}_T(f^{:A \Rightarrow H^{M^B} \Rightarrow M^B}) : (H^{M^A} \Rightarrow M^A) &\Rightarrow H^{M^B} \Rightarrow M^B \quad , \\ \text{flm}_T f &\triangleq t^{:R^{M^A}} \Rightarrow q^{:H^{M^B}} \Rightarrow (\text{flm}_M(x^{:A} \Rightarrow f\ x\ q))^{\uparrow R} t\ q \quad . \end{aligned} \quad (13.26)$$

Expressed through R ’s `flatMap` method, which is given by

$$\text{flm}_{Rg}^{:A \Rightarrow R^B} = t^{:R^A} \Rightarrow q^{:H^B} \Rightarrow (x^{:A} \Rightarrow g\ x\ q)^{\uparrow R} t\ q \quad , \quad (13.27)$$

the method flm_T can be written as

$$\text{flm}_T f = \text{flm}_R (y \Rightarrow q \Rightarrow \text{flm}_M(x \Rightarrow f\ x\ q)\ y) \quad . \quad (13.28)$$

13.5 Composed-outside transformers: Rigid monads

Proof The definition of \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:

$$\begin{aligned} f : A \Rightarrow H^{M^B} \Rightarrow M^B \diamond_T g : B \Rightarrow H^{M^C} \Rightarrow M^C & : A \Rightarrow H^{M^C} \Rightarrow M^C \quad ; \\ f \diamond_T g = t \Rightarrow q \Rightarrow \left(\tilde{f} \left((\text{flm}_M (b \Rightarrow g \, b \, q)) \downarrow^H q \right) \diamond_M (b \Rightarrow g \, b \, q) \right) t & , \end{aligned}$$

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 $x \diamond_M y = x \mathbin{\circ} \text{flm}_M y$ to find

$$f \diamond_T g = t \Rightarrow q \Rightarrow \left(\tilde{f} \left(p \downarrow^H q \right) \mathbin{\circ} p \right) t \quad \text{where } p = \text{flm}_M (x \Rightarrow g \, x \, q) \quad .$$

To obtain an implementation of flm_T , express flm_T through \diamond_T as

$$\text{flm}_T g : A \Rightarrow T^B = \text{id} : T^A \Rightarrow T^A \diamond_T g \quad .$$

Now we need to substitute $f : T^A \Rightarrow T^A = \text{id}$ into $f \diamond_T g$. Noting that \tilde{f} will then become

$$\tilde{f} = (h \Rightarrow k \Rightarrow \text{id} \, k \, h) = (h \Rightarrow k \Rightarrow k \, h) \quad ,$$

we get

$$\begin{aligned} \text{flm}_T g : A \Rightarrow T^B &= \text{id} \mathbin{\circ} \text{flm}_T g \\ \text{definition of } \diamond_T : &= t : T^A \Rightarrow q : H^{M^B} \Rightarrow \left(\tilde{f} \left(p \downarrow^H q \right) \mathbin{\circ} p \right) t \\ &\quad \text{where } p \triangleq \text{flm}_M (x \Rightarrow g \, x \, q) \\ \text{substitute } f = \text{id} : &= t \Rightarrow q \Rightarrow \left((h \Rightarrow k \Rightarrow k \, h) \left(p \downarrow^H q \right) \mathbin{\circ} p \right) t \\ \text{apply } k \text{ to } p \downarrow^H q : &= t \Rightarrow q \Rightarrow \left((k \Rightarrow k \left(p \downarrow^H q \right)) \mathbin{\circ} p \right) t \\ \text{definition of } \mathbin{\circ} : &= t \Rightarrow q \Rightarrow p \left(t \left(p \downarrow^H q \right) \right) \quad . \end{aligned}$$

By definition of the functor $R^A \triangleq H^A \Rightarrow A$, we raise any function $p : A \Rightarrow B$ into R as

$$\begin{aligned} p \uparrow^R : (H^A \Rightarrow A) &\Rightarrow H^B \Rightarrow B \quad , \\ p \uparrow^R r : H^A \Rightarrow A &\triangleq p \downarrow^H \mathbin{\circ} r \mathbin{\circ} p \\ &= q : H^B \Rightarrow p \left(r \left(p \downarrow^H q \right) \right) \quad . \end{aligned}$$

13 Computations in functor blocks. III. Monad transformers

Finally, renaming g to f , we obtain the desired code,

$$\text{flm}_T f = t \Rightarrow q \Rightarrow p^{\uparrow R} t q \quad \text{where } p \triangleq \text{flm}_M (x \Rightarrow f x q) \quad .$$

To express flm_T via flm_R , we just need to choose the value of g such that Eq. (13.27) becomes equal to Eq. (13.26). Comparing these two expressions, we find that we need

$$\text{flm}_M (x \Rightarrow f x q) = (y \Rightarrow g y q) \quad .$$

This is achieved if we define $g y q = \text{flm}_M (x^{:A} \Rightarrow f x q) y$, or equivalently

$$g = y \Rightarrow q \Rightarrow \text{flm}_M (x \Rightarrow f x q) y \quad .$$

This gives the desired Eq. (13.28).

Statement 13.5.1.3 The monad T has the methods `flatten` and `swap` defined by

$$\text{ftn}_T = t^{:T^{T^A}} \Rightarrow q^{:H^{M^A}} \Rightarrow q \triangleright \left(t \triangleright (\text{flm}_M (x^{:R^{M^A}} \Rightarrow x q))^{\uparrow R} \right) \quad , \quad (13.29)$$

$$\text{sw}_{R,M} = m^{:M^{R^A}} \Rightarrow q^{:H^{M^A}} \Rightarrow (r^{:R^A} \Rightarrow r(\text{pu}_M^{\downarrow H} q))^{\uparrow M} m \quad . \quad (13.30)$$

These functions are computationally equivalent (can be derived from each other). In the \triangleright -notation, the formula for $\text{sw}_{R,M}$ is

$$q \triangleright (m \triangleright \text{sw}_{R,M}) = m \triangleright (r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ r)^{\uparrow M} \quad . \quad (13.31)$$

Proof Using Eq. (13.26) and the relationship $\text{ftn}_T = \text{flm}_T \text{id}^{:T^A \Rightarrow T^A}$, we find

$$\begin{aligned} \text{ftn}_T t^{:T^{T^A}} &= \underline{\text{flm}_T}(\text{id}) t \\ \text{use Eq. (13.26)} : &= q \Rightarrow (\text{flm}_M (x^{:A} \Rightarrow f x q))^{\uparrow R} \underline{t q} \\ \text{definition of } \triangleright : &= q \Rightarrow \underline{q \triangleright (t \triangleright (\text{flm}_M (x^{:A} \Rightarrow f x q))^{\uparrow R})} \quad . \end{aligned}$$

Using Eq. (13.28) instead of Eq. (13.26), we get

$$\begin{aligned} \text{ftn}_T &= \text{flm}_T(\text{id}) \\ &= \text{flm}_R (y \Rightarrow q \Rightarrow \text{flm}_M (x \Rightarrow x q) y) \quad . \end{aligned}$$

Deriving the formulas for $\text{sw}_{R,M}$ We start with ftn_T as just obtained and substitute into Eq. (13.16):

$$\begin{aligned}
 \text{sw}_{R,M}(m) &= m \triangleright \text{pu}_M^{\uparrow R \uparrow M} \circ \text{pu}_R \circ \text{ftn}_T \\
 \text{use Eq. (13.28) : } &= m \triangleright \text{pu}_M^{\uparrow R \uparrow M} \circ \underline{\text{pu}_R \circ \text{flm}_R} (y \Rightarrow q \Rightarrow \text{flm}_M(x \Rightarrow x q) y) \\
 \text{left identity law of } R : &= m \triangleright \text{pu}_M^{\uparrow R \uparrow M} \circ (y \Rightarrow q \Rightarrow \text{flm}_M(x \Rightarrow x q) y) \\
 \text{function composition : } &= q \Rightarrow \text{flm}_M(x \Rightarrow x q) (m \triangleright \text{pu}_M^{\uparrow R \uparrow M}) \\
 \triangleright \text{ notation : } &= q \Rightarrow m \triangleright \text{pu}_M^{\uparrow R \uparrow M} \circ \text{flm}_M (x \Rightarrow x q) \\
 \text{express flm}_M \text{ via ftm}_M : &= q \Rightarrow m \triangleright \underline{\text{pu}_M^{\uparrow R \uparrow M} \circ (x \Rightarrow x q)^{\uparrow M}} \circ \text{ftn}_M \\
 \text{composition law of } M : &= q \Rightarrow m \triangleright (\text{pu}_M^{\uparrow R} \circ (x \Rightarrow x q))^{\uparrow M} \circ \text{ftn}_M \quad . \quad (13.32)
 \end{aligned}$$

It appears that simplifying this expression requires to rewrite the function $\text{pu}_M^{\uparrow R} \circ (x \Rightarrow x q)$. To proceed further, we need to use the definition of raising a function $f: A \Rightarrow B$ to the functor R ,

$$f^{\uparrow R} \triangleq r \cdot R^A \Rightarrow f^{\downarrow H} \circ r \circ f \quad ,$$

so we can write

$$\begin{aligned}
 &\text{pu}_M^{\uparrow R} \circ (x \Rightarrow x q) \\
 \text{function composition : } &= r \Rightarrow \underline{\text{pu}_M^{\uparrow R} r} q \\
 \text{definition of } \uparrow^R : &= r \Rightarrow q \triangleright \underline{\text{pu}_M^{\downarrow H} \circ r} \circ \text{pu}_M \\
 \text{forward composition : } &= (r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ r) \circ \text{pu}_M \quad . \quad (13.33)
 \end{aligned}$$

In deriving Eq. (13.33), we used the general property of the forward composition,

$$x \Rightarrow y \triangleright f(x, y) \circ g = (x \Rightarrow y \triangleright f(x, y)) \circ g \quad ,$$

where g must not depend on x or y . We can now rewrite Eq. (13.32) as

$$\begin{aligned}
 &q \Rightarrow m \triangleright (\text{pu}_M^{\uparrow R} \circ (x \Rightarrow x q))^{\uparrow M} \circ \text{ftn}_M \\
 \text{use Eq. (13.33) : } &= q \Rightarrow m \triangleright ((r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ r) \circ \text{pu}_M)^{\uparrow M} \circ \text{ftn}_M \\
 \text{functor composition for } M : &= q \Rightarrow m \triangleright ((r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ r))^{\uparrow M} \circ \underline{\text{pu}_M^{\uparrow M} \circ \text{ftn}_M} \\
 \text{identity law of } M : &= q \Rightarrow \underline{m \triangleright ((r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ r))^{\uparrow M}} \\
 \triangleright \text{ notation : } &= q \Rightarrow (r \Rightarrow r(\text{pu}_M^{\downarrow H} q))^{\uparrow M} m \quad .
 \end{aligned}$$

13 Computations in functor blocks. III. Monad transformers

The last expression coincides with Eq. (13.30).

The formula (13.31) follows by applying Eq. (13.30) to the arguments m and q . To make the computation clearer, we rename the bound variables m and q inside Eq. (13.30) to m_1 and q_1 :

$$\begin{aligned}
 & q \triangleright (m \triangleright (m_1 \Rightarrow q_1 \Rightarrow (r \Rightarrow r(\text{pu}_M^{\downarrow H} q_1))^{\uparrow M} m_1)) \\
 \text{apply to argument } m : & = q \triangleright (q_1 \Rightarrow m \triangleright (r \Rightarrow r(\text{pu}_M^{\downarrow H} q_1))^{\uparrow M}) \\
 \text{apply to argument } q : & = m \triangleright (r \Rightarrow r(\text{pu}_M^{\downarrow H} q))^{\uparrow M} \\
 \triangleright \text{ notation} : & = m \triangleright (r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ r)^{\uparrow M} .
 \end{aligned}$$

Deriving ftn_T from $\text{sw}_{R,M}$ Given the swap function defined by Eq. (13.30), we can recover the original ftn_T function from Eq. (13.29) via the standard formula (13.6), $\text{ftn}_T = \text{sw}^{\uparrow R} \circ \text{ftn}_R \circ \text{ftn}_M^{\uparrow R}$:

$$\begin{aligned}
 & \text{sw}^{\uparrow R} \circ \text{ftn}_R \circ \text{ftn}_M^{\uparrow R} \\
 \text{naturality of } \text{ftn}_R : & = \text{sw}^{\uparrow R} \circ \text{ftn}_M^{\uparrow R \uparrow R} \circ \text{ftn}_R \\
 \text{composition under } R : & = (\text{sw} \circ \text{ftn}_M^{\uparrow R})^{\uparrow R} \circ \text{ftn}_R \\
 \text{relating } \text{flm}_R \text{ and } \text{ftn}_R : & = \text{flm}_R(\text{sw} \circ \text{ftn}_M^{\uparrow R}) \\
 \text{use Eq. (13.27)} : & = t \Rightarrow q \Rightarrow (x^{\cdot A} \Rightarrow (\text{sw} \circ \text{ftn}_M^{\uparrow R}) x q)^{\uparrow R} t q . \quad (13.34)
 \end{aligned}$$

To proceed, we need to transform $\text{sw} \circ \text{ftn}_M^{\uparrow R}$ in some way:

$$\begin{aligned}
 & \text{sw} \circ \text{ftn}_M^{\uparrow R} \\
 \text{definitions} : & = (m \Rightarrow q \Rightarrow m \triangleright ((r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ r))^{\uparrow M}) \circ (r \Rightarrow \text{ftn}_M^{\downarrow H} \circ r \circ \text{ftn}_M) \\
 \text{composition} : & = m \Rightarrow \text{ftn}_M^{\downarrow H} \circ (q \Rightarrow m \triangleright ((r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ r))^{\uparrow M}) \circ \text{ftn}_M \\
 \text{expansion} : & = m \Rightarrow (q \Rightarrow q \triangleright \text{ftn}_M^{\downarrow H}) \circ (q \Rightarrow m \triangleright ((r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ r))^{\uparrow M}) \circ \text{ftn}_M \\
 \text{composition} : & = m \Rightarrow (q \Rightarrow m \triangleright ((r \Rightarrow q \triangleright \text{ftn}_M^{\downarrow H} \triangleright \text{pu}_M^{\downarrow H} \circ r))^{\uparrow M}) \circ \text{ftn}_M . \quad (13.35)
 \end{aligned}$$

13.5 Composed-outside transformers: Rigid monads

We can transform the sub-expression $(r \Rightarrow q \triangleright \text{ftn}_M^{\downarrow H} \triangleright \text{pu}_M^{\downarrow H} \circ r)$ to

$$\begin{aligned}
 &\triangleright \text{ notation : } r \Rightarrow q \triangleright \underline{\text{ftn}_M^{\downarrow H} \circ \text{pu}_M^{\downarrow H} \circ r} \\
 &\text{composition law of } H : = r \Rightarrow q \triangleright (\underline{\text{pu}_M \circ \text{ftn}_M})^{\downarrow H} \circ r \\
 &\text{left identity law of } M : = r \Rightarrow \underline{q \triangleright r} \\
 &\triangleright \text{ notation : } = r \Rightarrow r(q) \quad . \tag{13.36}
 \end{aligned}$$

Using this simplification, we continue transforming Eq. (13.35) as

$$\begin{aligned}
 m &\Rightarrow (q \Rightarrow m \triangleright ((r \Rightarrow q \triangleright \underline{\text{ftn}_M^{\downarrow H} \circ \text{pu}_M^{\downarrow H} \circ r})^{\uparrow M}) \circ \text{ftn}_M) \\
 \text{use Eq. (13.36) : } &= m \Rightarrow \underline{(q \Rightarrow m \triangleright (r \Rightarrow r(q))^{\uparrow M}) \circ \text{ftn}_M} \\
 \text{composition : } &= m \Rightarrow q \Rightarrow m \triangleright \underline{(r \Rightarrow r(q))^{\uparrow M} \circ \text{ftn}_M} \\
 \text{relating } \text{flm}_M \text{ and } \text{ftn}_M : &= m \Rightarrow q \Rightarrow \text{flm}_M (r \Rightarrow r(q)) m \quad .
 \end{aligned}$$

Substituting this instead of $\text{sw} \circ \text{ftn}_M^{\uparrow R}$ into Eq. (13.34), we get

$$\begin{aligned}
 &t \Rightarrow q \Rightarrow (x \Rightarrow \underline{\text{sw} \circ \text{ftn}_M^{\uparrow R}} x q)^{\uparrow R} t q \\
 &= t \Rightarrow q \Rightarrow (x \Rightarrow \text{flm}_M (r \Rightarrow r(q)) x)^{\uparrow R} t q \quad .
 \end{aligned}$$

The last expression is the same as Eq. (13.29).

Statement 13.5.1.4 Without assuming the monad laws for the function ftn_T , the laws in Theorem 13.3.2.1 hold for the swap function defined by Eq. (13.30).

Proof After replacing the base monad L by R , the required laws are

$$\begin{aligned}
 \text{pu}_R^{\uparrow M} \circ \text{sw} &= \text{pu}_R \quad , \quad \text{pu}_M \circ \text{sw} = \text{pu}_M^{\uparrow R} \quad , \\
 \text{ftn}_R^{\uparrow M} \circ \text{sw} &= \text{sw} \circ \text{sw}^{\uparrow R} \circ \text{ftn}_R \quad , \quad \text{ftn}_M \circ \text{sw} = \text{sw}^{\uparrow M} \circ \text{sw} \circ \text{ftn}_M^{\uparrow R} \quad .
 \end{aligned}$$

Proof of the inner-identity law Compute

$$\begin{aligned}
 &\text{pu}_R^{\uparrow M} \circ \underline{\text{sw}} \\
 \text{use Eq. (13.30) : } &= (m \Rightarrow \underline{m \triangleright \text{pu}_R^{\uparrow M}}) \circ (m \Rightarrow q \Rightarrow \underline{m \triangleright (r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ r)^{\uparrow M}}) \\
 \text{function composition : } &= m \Rightarrow q \Rightarrow m \triangleright \underline{\text{pu}_R^{\uparrow M} \circ (r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ r)^{\uparrow M}} \\
 \text{functor law of } M : &= m \Rightarrow q \Rightarrow m \triangleright (\text{pu}_R \circ (r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ r))^{\uparrow M} \quad . \tag{13.37}
 \end{aligned}$$

13 Computations in functor blocks. III. Monad transformers

To proceed, we simplify the expression $\text{pu}_R \circ (r \Rightarrow \dots)$:

$$\begin{aligned}
 & \text{pu}_R \circ (r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ r) \\
 \text{argument expansion : } &= (m \Rightarrow m \triangleright \text{pu}_R) \circ (r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \triangleright r) \\
 \text{function composition : } &= m \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \triangleright (m \triangleright \text{pu}_R) \quad . \quad (13.38)
 \end{aligned}$$

We now have to use the definition of pu_R , which is $\text{pu}_R = x \Rightarrow y \Rightarrow x$, or in the forwarding notation,

$$y \triangleright (x \triangleright \text{pu}_R) = x \quad . \quad (13.39)$$

With this simplification at hand, we continue from Eq. (13.38) to

$$\begin{aligned}
 & m \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \triangleright (m \triangleright \text{pu}_R) \\
 \text{use Eq. (13.39) : } &= m \Rightarrow m = \text{id} \quad .
 \end{aligned}$$

Therefore, Eq. (13.37) becomes

$$\begin{aligned}
 & m \Rightarrow q \Rightarrow m \triangleright (\text{pu}_R \circ (r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ r))^{\uparrow M} \\
 &= (m \Rightarrow q \Rightarrow m \triangleright \text{id}) \\
 &= (m \Rightarrow q \Rightarrow m) = \text{pu}_R \quad .
 \end{aligned}$$

This proves the inner-identity law.

Proof of the outer-identity law The left-hand side of this law is

$$\begin{aligned}
 & \text{pu}_M \circ \text{sw} \\
 \text{Eq. (13.30) : } &= (m \Rightarrow m \triangleright \text{pu}_M) \circ (m \Rightarrow q \Rightarrow m \triangleright (r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ r))^{\uparrow M} \\
 \text{function composition : } &= m \Rightarrow q \Rightarrow m \triangleright \text{pu}_M \circ (r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ r)^{\uparrow M} \\
 \text{naturality of pu}_M : &= m \Rightarrow q \Rightarrow m \triangleright (r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ r) \circ \text{pu}_M \\
 \triangleright \text{ notation : } &= m \Rightarrow q \Rightarrow m \triangleright (r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ r \circ \text{pu}_M) \\
 \text{apply function to } m : &= m \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ m \circ \text{pu}_M \\
 \text{argument expansion : } &= m \Rightarrow \text{pu}_M^{\downarrow H} \circ m \circ \text{pu}_M \\
 \text{definition of } \uparrow^R : &= \text{pu}_M^{\uparrow R} \quad .
 \end{aligned}$$

This is equal to the right-hand side of the law.

13.5 Composed-outside transformers: Rigid monads

Proof of the inner-interchange law The law is written as

$$\text{ftn}_R^{\uparrow M} \circ \text{sw} = \text{sw} \circ \text{sw}^{\uparrow R} \circ \text{ftn}_R \quad . \quad (13.40)$$

We will apply both sides of the law to arbitrary $m^{M^{RR^A}}$ and $q^{H^{M^A}}$, and transform both sides to the same expression.

Below, we will need a simplified formula for ftn_R derived from Eq. (13.27):

$$\begin{aligned} \text{ftn}_R &= \text{flm}_R(\text{id}) \\ \text{use Eq. (13.27)} : &= t \Rightarrow q \Rightarrow (x \Rightarrow x q)^{\uparrow R} t q \\ \text{definition of } \uparrow^R : &= t \Rightarrow q \Rightarrow \underline{(r \Rightarrow (x \Rightarrow x q)^{\downarrow H} \circ r \circ (x \Rightarrow x))} t q \\ \text{apply to argument} : &= t \Rightarrow q \Rightarrow ((x \Rightarrow q \triangleright x)^{\downarrow H} \circ t \circ (x \Rightarrow x q)) q \\ \text{use } \triangleright \text{ notation} : &= t \Rightarrow q \Rightarrow \underline{q \triangleright (q \triangleright (x \Rightarrow q \triangleright x))^{\downarrow H} \circ t} \quad . \quad (13.41) \end{aligned}$$

We first apply the left-hand side of the law (13.40) to m and q :

$$\begin{aligned} &q \triangleright (m \triangleright \text{ftn}_R^{\uparrow M} \circ \text{sw}) \\ \triangleright \text{ notation} : &= q \triangleright (m \triangleright \text{ftn}_R^{\uparrow M} \triangleright \text{sw}) \\ \text{use Eq. (13.31)} : &= m \triangleright \underline{\text{ftn}_R^{\uparrow M} \circ (r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ r)^{\uparrow M}} \\ \text{composition law for } M : &= m \triangleright (\text{ftn}_R \circ (r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ r))^{\uparrow M} \quad . \end{aligned}$$

We now need to simplify the sub-expression under $(\dots)^{\uparrow M}$:

$$\begin{aligned} &\underline{\text{ftn}_R \circ (r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \triangleright r)} \\ \text{function composition} : &= r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \triangleright \underline{\text{ftn}_R(r)} \\ \text{use Eq. (13.41)} : &= r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \triangleright (q \triangleright \underline{\text{pu}_M^{\downarrow H} \triangleright (x \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \triangleright x)^{\downarrow H} \circ r}) \\ \text{composition law for } H : &= r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ (q \triangleright (x \Rightarrow q \triangleright \underline{\text{pu}_M^{\downarrow H} \circ x \circ \text{pu}_M})^{\downarrow H} \circ r) \\ \text{definition of } \uparrow^R : &= r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ (q \triangleright (x \Rightarrow q \triangleright \text{pu}_M^{\uparrow R}(x))^{\downarrow H} \circ r) \quad . \end{aligned}$$

The left-hand side of the law (13.40) then becomes

$$m \triangleright (r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ (q \triangleright (x \Rightarrow q \triangleright \text{pu}_M^{\uparrow R}(x))^{\downarrow H} \circ r))^{\uparrow M} \quad .$$

13 Computations in functor blocks. III. Monad transformers

Now apply the right-hand side of the law (13.40) to m and q :

$$\begin{aligned}
 & q \triangleright (m \triangleright \text{sw} \circ \text{sw}^{\uparrow R} \circ \text{ftn}_R) \\
 \text{definition of } \uparrow^R : &= q \triangleright (\underline{m \triangleright \text{sw} \triangleright (x \Rightarrow \text{sw}^{\downarrow H} \circ x \circ \text{sw})} \triangleright \text{ftn}_R) \\
 \text{apply to arguments :} &= q \triangleright (\underline{\text{ftn}_R(\text{sw}^{\downarrow H} \circ \text{sw}(m) \circ \text{sw})}) \\
 \text{use Eq. (13.41) :} &= q \triangleright (q \triangleright \underline{(x \Rightarrow q \triangleright x)}^{\downarrow H} \circ \text{sw}^{\downarrow H} \circ \text{sw}(m) \circ \text{sw}) \\
 \text{composition law of } H : &= q \triangleright (q \triangleright \underline{(\text{sw} \circ (x \Rightarrow q \triangleright x))}^{\downarrow H} \circ \text{sw}(m) \circ \text{sw}) \quad . \quad (13.42)
 \end{aligned}$$

To proceed, we simplify the sub-expression $\text{sw}(m) \circ \text{sw}$ separately by computing the function compositions:

$$\begin{aligned}
 & \text{sw}(m) \circ \text{sw} \\
 &= (q_1 \Rightarrow m \triangleright (r \Rightarrow q_1 \triangleright \text{pu}_M^{\downarrow H} \circ r))^{\uparrow M} \circ (y \Rightarrow q_2 \Rightarrow y \triangleright (r \Rightarrow q_2 \triangleright \text{pu}_M^{\downarrow H} \circ r))^{\uparrow M} \\
 &= q_1 \Rightarrow q_2 \Rightarrow (m \triangleright (r \Rightarrow q_1 \triangleright \text{pu}_M^{\downarrow H} \circ r))^{\uparrow M} \triangleright (r \Rightarrow q_2 \triangleright \text{pu}_M^{\downarrow H} \circ r)^{\uparrow M} \\
 &= q_1 \Rightarrow q_2 \Rightarrow m \triangleright ((r \Rightarrow q_1 \triangleright \text{pu}_M^{\downarrow H} \circ r) \circ (r \Rightarrow q_2 \triangleright \text{pu}_M^{\downarrow H} \circ r))^{\uparrow M} \quad .
 \end{aligned}$$

Using this formula, we can write, for any z of a suitable type,

$$\begin{aligned}
 & q \triangleright (z \triangleright \text{sw}(m) \circ \text{sw}) = m \triangleright ((r \Rightarrow z \triangleright \text{pu}_M^{\downarrow H} \circ r) \circ (r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ r))^{\uparrow M} \\
 \text{function composition :} &= m \triangleright (r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ (z \triangleright \text{pu}_M^{\downarrow H} \circ r))^{\uparrow M} \quad . \quad (13.43)
 \end{aligned}$$

Now we can substitute this into Eq. (13.42):

$$\begin{aligned}
 & q \triangleright (q \triangleright (\text{sw} \circ (x \Rightarrow q \triangleright x))^{\downarrow H} \triangleright \text{sw}(m) \circ \text{sw}) \\
 \text{use Eq. (13.43) :} &= m \triangleright (r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ (q \triangleright (\text{sw} \circ (x \Rightarrow q \triangleright x))^{\downarrow H} \triangleright \text{pu}_M^{\downarrow H} \circ r))^{\uparrow M} \\
 H\text{'s composition :} &= m \triangleright (r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ (q \triangleright (\underline{\text{pu}_M \circ \text{sw} \circ (x \Rightarrow q \triangleright x)}^{\downarrow H} \circ r))^{\uparrow M}) \\
 \text{outer-identity :} &= m \triangleright (r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ (q \triangleright (\underline{\text{pu}_M^{\uparrow R} \circ (x \Rightarrow q \triangleright x)}^{\downarrow H} \circ r))^{\uparrow M}) \\
 \text{composition :} &= m \triangleright (r \Rightarrow q \triangleright \text{pu}_M^{\downarrow H} \circ (q \triangleright (x \Rightarrow q \triangleright \text{pu}_M^{\uparrow R} \circ x)^{\downarrow H} \circ r))^{\uparrow M} \quad .
 \end{aligned}$$

We arrived at the same expression as the left-hand side of the law.

Proof of the outer-interchange law The law is written as

$$\text{ftn}_M \circ \text{sw} = \text{sw}^{\uparrow M} \circ \text{sw} \circ \text{ftn}_M^{\uparrow R} \quad . \quad (13.44)$$

13.5 Composed-outside transformers: Rigid monads

We will apply both sides of the law to arbitrary $m:M^{M^{R^A}}$ and $q:H^{M^A}$, and transform both sides to the same expression. We begin with the more complicated right-hand side:

$$\begin{aligned}
 & q \triangleright (m \triangleright \underline{\text{sw}}^{\uparrow M} \circ \underline{\text{sw}} \circ \text{ftn}_M^{\uparrow R}) \\
 \triangleright \text{notation} : &= q \triangleright ((m \triangleright \text{sw}^{\uparrow M} \triangleright \text{sw}) \triangleright \underline{\text{ftn}}_M^{\uparrow R}) \\
 \text{definition of } \uparrow^R : &= q \triangleright (\text{ftn}_M^{\downarrow H} \circ (m \triangleright \text{sw}^{\uparrow M} \triangleright \text{sw}) \circ \text{ftn}_M) \\
 \triangleright \text{notation} : &= (q \triangleright \underline{\text{ftn}}_M^{\downarrow H} \triangleright (m \triangleright \underline{\text{sw}}^{\uparrow M} \triangleright \underline{\text{sw}})) \triangleright \text{ftn}_M \\
 \text{use Eq. (13.31)} : &= (m \triangleright \underline{\text{sw}}^{\uparrow M} \triangleright (r \Rightarrow q \triangleright \underline{\text{ftn}}_M^{\downarrow H} \triangleright \underline{\text{pu}}_M^{\downarrow H} \circ r)^{\uparrow M}) \triangleright \text{ftn}_M \\
 \text{composition for } H \text{ and } M : &= m \triangleright (\text{sw} \circ (r \Rightarrow q \triangleright (\underline{\text{pu}}_M \circ \text{ftn}_M)^{\downarrow H} \circ r))^{\uparrow M} \circ \text{ftn}_M \\
 \text{left identity law of } M : &= m \triangleright (\text{sw} \circ (r \Rightarrow q \triangleright r))^{\uparrow M} \circ \text{ftn}_M \quad . \quad (13.45)
 \end{aligned}$$

Let us simplify the sub-expression $\text{sw} \circ (r \Rightarrow q \triangleright r)$ separately:

$$\begin{aligned}
 & \underline{\text{sw}} \circ (r \Rightarrow q \triangleright r) = (x \Rightarrow x \triangleright \text{sw}) \circ (r \Rightarrow q \triangleright r) \\
 \text{function composition} : &= (x \Rightarrow \underline{q \triangleright (x \triangleright \text{sw})}) \\
 \text{use Eq. (13.31)} : &= \underline{x \Rightarrow x \triangleright (r \Rightarrow q \triangleright \underline{\text{pu}}_M^{\downarrow H} \circ r)^{\uparrow M}} \\
 \text{expand argument} : &= (r \Rightarrow q \triangleright \underline{\text{pu}}_M^{\downarrow H} \circ r)^{\uparrow M} \quad . \quad (13.46)
 \end{aligned}$$

Substituting this expression into Eq. (13.45), we get

$$\begin{aligned}
 & m \triangleright (\underline{\text{sw} \circ (r \Rightarrow q \triangleright r)})^{\uparrow M} \circ \text{ftn}_M \\
 \text{use Eq. (13.46)} : &= m \triangleright (r \Rightarrow q \triangleright \underline{\text{pu}}_M^{\downarrow H} \circ r)^{\uparrow M \uparrow M} \circ \text{ftn}_M \\
 \text{naturality of } \text{ftn}_M : &= m \triangleright \text{ftn}_M \circ (r \Rightarrow q \triangleright \underline{\text{pu}}_M^{\downarrow H} \circ r)^{\uparrow M} \quad .
 \end{aligned}$$

Now write the left-hand side of the law:

$$\begin{aligned}
 & q \triangleright (m \triangleright \underline{\text{ftn}}_M \circ \underline{\text{sw}}) = q \triangleright (m \triangleright \underline{\text{ftn}}_M \triangleright \underline{\text{sw}}) \\
 \text{use Eq. (13.31)} : &= m \triangleright \text{ftn}_M \triangleright (r \Rightarrow q \triangleright \underline{\text{pu}}_M^{\downarrow H} \circ r)^{\uparrow M} \quad .
 \end{aligned}$$

This is equal to the right-hand side we just obtained.

Statement 13.5.1.5 The monadic naturality laws in Theorem 13.3.7.1 hold for the swap function defined by Eq. (13.30) and the base monad $L \triangleq R$.

13 Computations in functor blocks. III. Monad transformers

Proof The monadic naturality laws are

$$\text{sw}_{R,M} \circ \phi^{\uparrow R} = \phi \circ \text{sw}_{R,N} \quad , \quad \text{sw}_{R,M} \circ \theta = \theta^{\uparrow M} \quad ,$$

where $\phi : M \rightsquigarrow N$ and $\theta : R \rightsquigarrow \text{Id}$ are arbitrary monadic morphisms. Eq. (13.32)

$$\text{sw}(m) = q \Rightarrow m \triangleright (\text{pu}_M^{\uparrow R} \circ (x \Rightarrow x q))^{\uparrow M} \circ \text{ftn}_M$$

To verify the first law, apply both sides to arbitrary m and q . The left-hand side:

$$\begin{aligned} q \triangleright (m \triangleright \text{sw}_{R,M} \circ \phi^{\uparrow R}) &= q \triangleright (m \triangleright \text{sw}_{R,M} \triangleright \phi^{\uparrow R}) \\ \text{definition of } \uparrow^R &:= q \triangleright (\phi^{\downarrow H} \circ (m \triangleright \text{sw}_{R,M})) \circ \phi \\ \text{notation} &:= \underline{(q \triangleright \phi^{\downarrow H})} \triangleright (m \triangleright \text{sw}_{R,M}) \triangleright \phi \\ \text{use Eq. (13.31)} &:= m \triangleright (r \Rightarrow q \triangleright \phi^{\downarrow H} \triangleright \text{pu}_M^{\downarrow H} \circ r)^{\uparrow M} \triangleright \phi \\ \text{composition law for } H &:= m \triangleright (r \Rightarrow q \triangleright \underline{(\text{pu}_M^{\downarrow H} \circ \phi)^{\downarrow H}} \circ r)^{\uparrow M} \circ \phi \\ \text{identity law for } \phi &:= m \triangleright (r \Rightarrow q \triangleright \text{pu}_N^{\downarrow H} \circ r)^{\uparrow M} \circ \phi \\ \text{naturality of } \phi &:= m \triangleright \underline{\phi \circ (r \Rightarrow q \triangleright \text{pu}_N^{\downarrow H} \circ r)^{\uparrow N}} \quad . \end{aligned}$$

The right-hand side, when applied to m and q , gives the same expression:

$$\begin{aligned} q \triangleright (m \triangleright \phi \circ \text{sw}_{R,N}) &= q \triangleright (m \triangleright \phi \triangleright \text{sw}_{R,N}) \\ \text{use Eq. (13.31)} &:= m \triangleright \phi \triangleright (r \Rightarrow q \triangleright \text{pu}_N^{\downarrow H} \circ r)^{\uparrow N} \quad . \end{aligned}$$

To argue that the second law holds,³ apply the left-hand side to m and q :

$$\begin{aligned} q \triangleright (m \triangleright \text{sw}_{R,M} \circ \theta) &= q \triangleright (m \triangleright \text{sw}_{R,M} \triangleright \theta) \\ &= q \triangleright ((q_1 \Rightarrow m \triangleright (r \Rightarrow q_1 \triangleright \text{pu}_M^{\downarrow H} \circ r)^{\uparrow M}) \triangleright \theta) \quad . \end{aligned} \tag{13.47}$$

This expression cannot be simplified any further; and neither can the right-hand side $q \triangleright (m \triangleright \theta^{\uparrow M})$. We need more detailed information about the function θ .

The type of θ is

$$\theta : \forall A. (H^A \Rightarrow A) \Rightarrow A \quad .$$

To implement a function of this type, we need to write code that takes an argument of type $H^A \Rightarrow A$ and returns a value of type A . Since the type A is arbitrary,

³I could not find a fully rigorous proof of this law. Below I will indicate the part that lacks rigor.

13.5 Composed-outside transformers: Rigid monads

the code of θ cannot store a fixed value of type A to use as the return value. The only possibility to implement a function θ with the required type signature seems to be by substituting a value of type H^A into the given argument of type $H^A \Rightarrow A$, which will return the result of type A . So,⁴ we need to produce a value of type H^A for an arbitrary type A , that is, a value of type $\forall A. H^A$. Using the contravariant Yoneda identity, we can simplify this type expression to the type H^1 :

$$\begin{aligned} \forall A. H^A &\cong \forall A. \underline{1} \Rightarrow H^A \\ \text{use identity } (A \Rightarrow \underline{1}) \cong \underline{1} : &= \forall A. (A \Rightarrow \underline{1}) \Rightarrow H^A \\ \text{contravariant Yoneda identity :} &= H^1 . \end{aligned}$$

So, we can construct a θ if we store a value h_1 of type H^1 and compute $h : H^A$ as

$$h^{H^A} = h_1^{H^1} \triangleright (a^{A \Rightarrow \underline{1}} \Rightarrow 1)^{\downarrow H} .$$

Given a fixed value $h_1 : H^1$, the code of θ is

$$(r^{H^A \Rightarrow A}) \triangleright \theta \triangleq h_1 \triangleright (_ \Rightarrow 1)^{\downarrow H} \triangleright r . \quad (13.48)$$

Let us check whether this θ is a monadic morphism $R \rightsquigarrow \text{Id}$. We need to verify the two laws of monadic morphisms,

$$\text{pu}_R \circ \theta = \text{id} , \quad \text{ftn}_R \circ \theta = \theta^{\uparrow R} \circ \theta = \theta \circ \theta .$$

The identity law, applied to an arbitrary $x : A$, is

$$\begin{aligned} x \triangleright \text{pu}_R \circ \theta &= (x \triangleright \text{pu}_R) \triangleright \theta \\ \text{definition of } r \triangleright \theta : &= h_1 \triangleright (_ \Rightarrow 1)^{\downarrow H} \triangleright (x \triangleright \text{pu}_R) \\ \text{definition of } x \triangleright \text{pu}_R : &= (h_1 \triangleright (_ \Rightarrow 1)^{\downarrow H}) \triangleright (_ \Rightarrow x) \\ \text{function composition :} &= x . \end{aligned}$$

This verifies the identity law. The composition law, applied to an arbitrary r :

⁴This is where the argument is not rigorous: I have not proved rigorously that the type $\forall A. (H^A \Rightarrow A) \Rightarrow A$ is *equivalent* to $\forall A. H^A$.

13 Computations in functor blocks. III. Monad transformers

R^{R^A} , expands to

$$\begin{aligned}
 r \triangleright \text{ftn}_R \circ \theta &= \underline{r \triangleright \text{ftn}_R} \triangleright \theta \\
 \text{definition of } \text{ftn}_R : &= \underline{r \triangleright (t \Rightarrow q \Rightarrow q \triangleright (q \triangleright (x \Rightarrow q \triangleright x)^{\downarrow H} \circ t))} \triangleright \theta \\
 \text{apply to } r : &= \underline{(q \Rightarrow q \triangleright (q \triangleright (x \Rightarrow q \triangleright x)^{\downarrow H} \circ r))} \triangleright \theta \\
 \text{definition (13.48) of } \theta : &= h_1 \triangleright (_ \Rightarrow 1)^{\downarrow H} \triangleright \underline{(q \Rightarrow q \triangleright (q \triangleright (x \Rightarrow q \triangleright x)^{\downarrow H} \circ r))} \\
 \text{apply to } h_1 \triangleright (_ \Rightarrow 1)^{\downarrow H} : &= h_1 \triangleright (_ \Rightarrow 1)^{\downarrow H} \triangleright \underline{(h_1 \triangleright (_ \Rightarrow 1)^{\downarrow H} \circ (x \Rightarrow \dots)^{\downarrow H} \circ r)} \\
 \text{composition under } H : &= h_1 \triangleright (_ \Rightarrow 1)^{\downarrow H} \triangleright (h_1 \triangleright (_ \Rightarrow 1)^{\downarrow H} \triangleright r) \\
 \text{definition (13.48) of } \theta : &= r \triangleright \theta \triangleright \theta = r \triangleright \theta \circ \theta \quad .
 \end{aligned}$$

This verifies the composition law; so θ is indeed a monadic morphism.

Using the code of θ defined in Eq. (13.48), we can now verify the monadic naturality law of $\text{sw}_{R,M}$ (with respect to the runners θ of that form). The left-hand side of the law is given by Eq. (13.47) and is rewritten as

$$\begin{aligned}
 &q \triangleright ((q_1 \Rightarrow m \triangleright (r \Rightarrow q_1 \triangleright \text{pu}_M^{\downarrow H} \circ r)^{\uparrow M}) \triangleright \theta) \\
 \text{definition of } \theta : &= q \triangleright (h_1 \triangleright (_ \Rightarrow 1)^{\downarrow H} \triangleright \underline{(q_1 \Rightarrow m \triangleright (r \Rightarrow q_1 \triangleright \text{pu}_M^{\downarrow H} \circ r)^{\uparrow M})}) \\
 \text{apply to argument :} &= q \triangleright (m \triangleright (r \Rightarrow h_1 \triangleright (_ \Rightarrow 1)^{\downarrow H} \triangleright \underline{\text{pu}_M^{\downarrow H} \circ r})^{\uparrow M}) \\
 H\text{'s composition law :} &= q \triangleright (m \triangleright (r \Rightarrow h_1 \triangleright (\underline{\text{pu}_M \circ (_ \Rightarrow 1)})^{\downarrow H} \triangleright r)^{\uparrow M}) \\
 \text{compose functions :} &= q \triangleright (m \triangleright (r \Rightarrow h_1 \triangleright (_ \Rightarrow 1)^{\downarrow H} \triangleright r)^{\uparrow M}) \quad .
 \end{aligned}$$

The right-hand side is

$$\begin{aligned}
 &q \triangleright (m \triangleright \theta^{\uparrow M}) \\
 \text{function expansion :} &= q \triangleright (m \triangleright (r \Rightarrow r \triangleright \theta)^{\uparrow M}) \\
 \text{definition (13.48) of } \theta : &= q \triangleright (m \triangleright (r \Rightarrow h_1 \triangleright (_ \Rightarrow 1)^{\downarrow H} \triangleright r)^{\uparrow M}) \quad .
 \end{aligned}$$

This expression is now the same as the left-hand side.

13.5.2 Rigid monad construction 2: composition

Functor composition is the second construction that produces rigid monads. This is a consequence of the properties of monad transformer stacks.

13.5 Composed-outside transformers: Rigid monads

Statement 13.5.2.1 The composition $R_1^{R_2^*}$ of two rigid monads R_1 and R_2 is also a rigid monad.

Proof Since R_1 is rigid, its outside-composition $R_1 \circ M$ with any other monad M is a monad. So $R_1 \circ R_2$ is a monad. To show that $R_1 \circ R_2$ is a rigid monad, we need to show that its monad transformer is of the composed-outside kind. By Theorem 13.2.6, the stacking of monad transformers T_{R_1} and T_{R_2} is a lawful monad transformer. Since the transformers for R_1 and R_2 are of the composed-outside kind, $T_{R_1}^M = R_1 \circ M$ and $T_{R_2}^M = R_2 \circ M$, the stack of transformers is

$$T_{R_1}^{T_{R_2}^M} = R_1 \circ T_{R_2}^M = R_1 \circ (R_2 \circ M) = R_1 \circ R_2 \circ M \quad .$$

Therefore $T^M \triangleq R_1 \circ R_2 \circ M$ is a monad transformer applied to the foreign monad M . This shows, by definition of a rigid monad, that $R_1 \circ R_2$ is a rigid monad.

Example 13.5.2.2 Consider the functor composition of the `Search` monad $R_1^A \triangleq (A \Rightarrow Q) \Rightarrow A$ and the `Reader` monad $R_2^A \triangleq Z \Rightarrow A$:

$$P^A \triangleq ((Z \Rightarrow A) \Rightarrow Q) \Rightarrow Z \Rightarrow A \quad .$$

It follows from Statement 13.5.2.1 that the functor P^\bullet is a rigid monad; so P' 's transformer is of the composed-outside kind. The transformed monad for any foreign monad M is

$$T^A \triangleq ((Z \Rightarrow M^A) \Rightarrow Q) \Rightarrow Z \Rightarrow M^A \quad .$$

To define the monad methods for T , we need to use the definitions of the transformers $T_{R_1}^M$ and $T_{R_2}^M$. Since both the `Search` and the `Reader` monads are special cases of the `Choice` monad construction (Section 13.5.1) where the contrafunctor H is chosen to be $H^A \triangleq A \Rightarrow Q$ and $H^A \triangleq Z$ respectively, we can use Eq. (13.26) to define the `flatMap` methods for the transformers $T_{R_1}^M$ and $T_{R_2}^M$:

```
type R1[A] = (A => Q) => A
def flatMap1[A, B, M[_]: Monad](r1: R1[M[A]])(f: A => R1[M[B]]): R1[M[B]] =

type R2[A] = Z => A
def flatMap2[A, B, M[_]: Monad](r2: R2[M[A]])(f: A => R2[M[B]]): R2[M[B]] =
```

Finally, we define the `flatMap` method for T as

```
type T[A] = R1[R2[A]]
def flatMap1[A, B, M[_]: Monad](t: T[M[A]])(f: A => T[M[B]]): T[M[B]] =
```

13 Computations in functor blocks. III. Monad transformers

In Section 13.5.1, we have proved the monad transformer laws for $T_R^M \triangleq R \circ M$ by defining a `swap` function with type signature

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

and proving its laws. Does the composed monad $R_1 \circ R_2$ have a suitable `swap` function,

$$\text{sw}_{R_1 \circ R_2, M} : M \circ R_1 \circ R_2 \rightsquigarrow R_1 \circ R_2 \circ M \quad ,$$

satisfying the required laws?

13.5.3 Rigid monad construction 3: product

Statement 13.5.3.1 The product of rigid monads, $R_1^A \times R_2^A$, is a rigid monad.

13.5.4 Rigid monad construction 4: selector

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.

13.5.5 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 “`fiR`”) 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 `foR`), 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

$$\text{fi}_R \circ \text{fo}_R = \text{id} \quad .$$

***Not sure that it is really easier to reason about `fo`.

13.5 Composed-outside transformers: Rigid monads

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

$$\begin{aligned} \text{fi} &: (A \Rightarrow H^B \Rightarrow B) \Rightarrow H^{A \Rightarrow B} \Rightarrow A \Rightarrow B \quad , \\ \text{fi} &= f \cdot A \Rightarrow H^B \Rightarrow B \Rightarrow h \cdot H^{A \Rightarrow B} \Rightarrow a \Rightarrow f(a) \circ ((b \Rightarrow _ \Rightarrow b) \downarrow^H h) \quad . \end{aligned}$$

It is easier to calculate with 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

$$\begin{aligned} \tilde{\text{fi}} &: (H^B \Rightarrow A \Rightarrow B) \Rightarrow H^{A \Rightarrow B} \Rightarrow A \Rightarrow B \quad , \\ \tilde{\text{fi}} &= f \cdot H^B \Rightarrow A \Rightarrow B \Rightarrow h \cdot H^{A \Rightarrow B} \Rightarrow f \circ ((b \Rightarrow _ \Rightarrow b) \downarrow^H h) \quad , \\ \tilde{\text{fo}} &: (H^{A \Rightarrow B} \Rightarrow A \Rightarrow B) \Rightarrow H^B \Rightarrow A \Rightarrow B \quad , \\ \tilde{\text{fo}} &= g \cdot H^{A \Rightarrow B} \Rightarrow A \Rightarrow B \Rightarrow h \cdot H^B \Rightarrow a \cdot A \Rightarrow g \circ ((p \cdot A \Rightarrow B \Rightarrow p a) \downarrow^H h) a \quad . \end{aligned}$$

To show the non-degeneracy law for R , compute

$$\begin{aligned} &\tilde{\text{fo}} (\tilde{\text{fi}} f) h \cdot H^B a \cdot A \\ \text{definition of } \tilde{\text{fo}}: &= (\tilde{\text{fi}} f) \circ ((p \Rightarrow p a) \downarrow^H h) a \\ \text{definition of } \tilde{\text{fi}}: &= f \circ ((b \Rightarrow _ \Rightarrow b) \downarrow^H (p \Rightarrow p a) \downarrow^H h) a \\ \text{composition law for } H: &= f \circ ((b \Rightarrow _ \Rightarrow b \circ p \Rightarrow p a) \downarrow^H h) a \\ \text{simplify } b \Rightarrow _ \Rightarrow b \circ p \Rightarrow p a \text{ to:} &= f \circ ((b \Rightarrow b) \downarrow^H h) a \\ \text{identity law for } H: &= f h a \quad . \end{aligned}$$

We obtained $\tilde{\text{fo}} (\tilde{\text{fi}} f) h a = f h a$, so the non-degeneracy law $\tilde{\text{fi}} \circ \tilde{\text{fo}} = \text{id}$ holds.

Statement 1. The composition $R_1^{R_2^*}$ 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

13 Computations in functor blocks. III. Monad transformers

fi_M . So we can implement $\tilde{\text{fi}}_T$ using fi_M :

$$\begin{aligned}\tilde{\text{fi}}_T &: (H^{M^B} \Rightarrow A \Rightarrow M^B) \Rightarrow H^{M^{A \Rightarrow B}} \Rightarrow M^{A \Rightarrow B} \\ \tilde{\text{fi}}_T &= f \Rightarrow h \Rightarrow \text{fi}_M \left(f \left((b \Rightarrow _ \Rightarrow b)^{\uparrow M \downarrow H} h \right) \right) \\ \tilde{\text{fo}}_T &: (H^{M^{A \Rightarrow B}} \Rightarrow M^{A \Rightarrow B}) \Rightarrow H^{M^B} \Rightarrow A \Rightarrow M^B \\ \tilde{\text{fo}}_T &= g \Rightarrow h \Rightarrow a \Rightarrow \text{fo}_M \left(g \left((p^{A \Rightarrow B} \Rightarrow p a)^{\uparrow M \downarrow H} h \right) \right) a\end{aligned}$$

To show the non-degeneracy law for T , compute

$$\begin{aligned}\tilde{\text{fo}}_T \left(\tilde{\text{fi}}_T f \right) h &: H^{M^B} a : A \\ \text{insert the definition of } \tilde{\text{fo}}_T &: = \text{fo}_M \left(\left(\tilde{\text{fi}}_T f \right) \left((p \Rightarrow p a)^{\uparrow M \downarrow H} h \right) \right) a \\ \text{insert the definition of } \tilde{\text{fi}}_T &: = \text{fo}_M \left(\text{fi}_M \left(f \left((b \Rightarrow _ \Rightarrow b)^{\uparrow M \downarrow H} (p \Rightarrow p a)^{\uparrow M \downarrow H} h \right) \right) \right) a \\ \text{nondegeneracy law for } \text{fi}_M &: = f \left((b \Rightarrow _ \Rightarrow b)^{\uparrow M \downarrow H} (p \Rightarrow p a)^{\uparrow M \downarrow H} h \right) a \\ \text{composition laws for } M, H &: = f \left((b \Rightarrow _ \Rightarrow b \circ p \Rightarrow p a)^{\uparrow M \downarrow H} h \right) a \\ \text{simplify} &: = f \left((b \Rightarrow b)^{\uparrow M \downarrow H} h \right) a \\ \text{identity laws for } M, H &: = f h a \quad .\end{aligned}$$

We obtained $\tilde{\text{fo}}_T \left(\tilde{\text{fi}}_T f \right) h a = f h a$. Therefore the non-degeneracy law $\tilde{\text{fi}}_T ; \tilde{\text{fo}}_T = \text{id}$ holds.

13.6 Recursive monad transformers

13.6.1 Transformer for the free monad FreeT

13.6.2 Transformer for the list monad ListT

13.7 Monad transformers for monad constructions

13.7.1 Product of monad transformers

13.7.2 Free pointed monad transformer

13.8 Irregular and incomplete monad transformers

13.8.1 The state monad transformer StateT

13.8.2 The continuation monad transformer ContT

13.8.3 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

$$\text{blift} : (\forall B. (A \Rightarrow F^B) \Rightarrow F^B) \Rightarrow \forall C. (A \Rightarrow M^{F^C}) \Rightarrow M^{F^C}$$

13 Computations in functor blocks. III. Monad transformers

since this type signature cannot be implemented. The codensity transformer also does not have any of the required “runner” transformations `mr` and `br`,

$$\begin{aligned} \text{mr} &: (M^\bullet \rightsquigarrow N^\bullet) \Rightarrow (\forall B. (A \Rightarrow M^{F^B}) \Rightarrow M^{F^B}) \Rightarrow \forall C. (A \Rightarrow N^{F^C}) \Rightarrow N^{F^C} \quad , \\ \text{br} &: ((\forall B. (A \Rightarrow F^B) \Rightarrow F^B) \Rightarrow A) \Rightarrow (\forall C. (A \Rightarrow M^{F^C}) \Rightarrow M^{F^C}) \Rightarrow M^A \quad . \end{aligned}$$

13.9 Summary and discussion

14 Recursive types

14.1 Fixpoints and type recursion schemes

14.2 Row polymorphism and OO programming

14.3 Column polymorphism

14.4 Discussion

15 Co-inductive typeclasses. Comonads

15.1 Practical use

15.2 Laws and structure

15.3 Co-free constructions

15.4 Co-free comonads

15.5 Comonad transformers

15.6 Discussion

16 Irregular typeclasses*

16.1 Distributive functors

16.2 Monoidal monads

16.3 Lenses and prisms

16.4 Discussion

17 Summary and discussion

18 Essay: Software engineers and software artisans

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.

18.1 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.

18.2 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.

18.3 Programmers today are artisans, not engineers

Now I will argue that programmers are *not engineers* in the sense we normally see the engineering professions.

18.3.1 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

18.3 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 computers, without ever having studied quantum physics or calculus.

18.3.2 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 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 theoretical 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

18.3 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 “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. choosing 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.

18.3.3 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.

18.4 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 work using artisanal methods, and their culture and processes 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, etc.) 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 is devoted not to showing code but to a peculiar kind of mathematical reasoning.

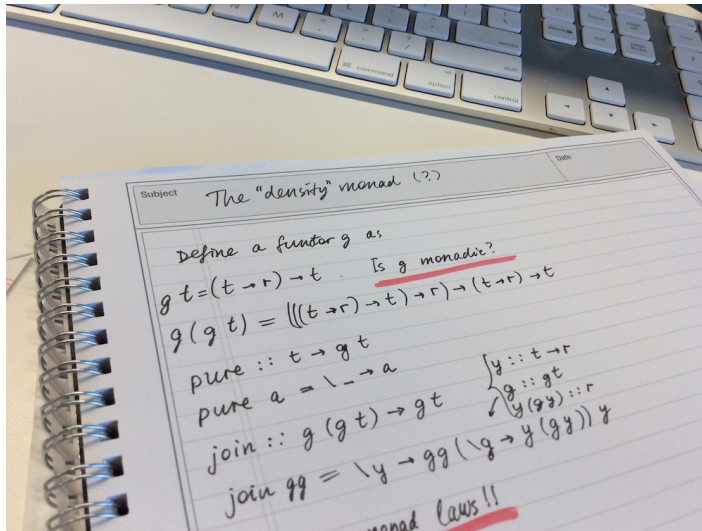


Figure 18.1: Example calculation in functional type theory.

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 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. 18.1). 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 “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.

18.5 Does software need engineers, 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 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

18.5 Does software need engineers, or are artisans good enough?

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.

19 Essay: Towards functional data engineering with Scala

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.

19.1 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. 19.1).

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. 19.2).



Figure 19.1: A nonsensical calculation arises when mixing incompatible data types.

19.2 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⇒B`. All modern functional languages such as OCaml, Haskell, Scala, F#, Swift, Elm, and 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.

19.3 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 19.2: 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 typeclasses and higher-order generic collections.

It remains to see why Scala became the *lingua franca* of big data and not, say, Haskell.

19.4 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 MacOS/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 requirements 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, only Scala and Haskell directly support typeclasses and higher-order generics, which are necessary for expressing abstractions such as automatically parallelized data sets or asynchronous data streams.

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 built-in 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.

19.5 Conclusion

The decisive advantages of Scala over other contenders (such as OCaml, Haskell, F#, or Swift) are

19 Essay: Towards functional data engineering with Scala

1. functional collections in the standard library;
2. a highly sophisticated type system, with support for typeclasses 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.

20 “Applied functional type theory”: A proposal

What exactly is the extent of “theory” that a practicing functional programmer should know in order to be effective at writing code in the functional paradigm? In my view, this question is not yet resolved. In this book, I present a coherent body of theoretical knowledge that I believe fits the description of “practicable functional programming theory”. This body of knowledge is, or should be, understood as a branch of computer science, and I propose to call it **applied functional type theory** (AFTT). This is the area of theoretical computer science that should serve the needs of functional programmers 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 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. This book is a first step towards formulating AFTT.

In this book, I show code in Scala because I am familiar with that language. However, most of this material will work equally well in Haskell, OCaml, and other FP languages. This is so because the science of functional programming, which I call AFTT, is not a set of tricks specific to Scala or Haskell. An advanced user of any other functional programming language will have to face the same questions and struggle with the same practical issues.

20.1 AFTT is not found in computer science curricula

Traditional courses of theoretical computer science (algorithms and data structures, complexity theory, distributed systems, databases, network systems, com-

pilars, operating systems) 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, although typically only at the graduate level.

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 properties 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. Functional programming does not require graduate-level theoretical studies.

As an example of theoretical material that *is* within the scope of AFTT, consider the equational laws imposed on applicative functors (see Chapter 10).

It is essential for a practicing functional programmer to be able to recognize and use applicative functors. An applicative functor is a data structure specifying declaratively a set of operations that run independently of each other. Programs can then easily combine these operations, 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 satisfies the laws of applicative functors. In a given application, a data structure may be dictated in part by the business logic rather than by a programmer’s choice. 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 mathematical 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 and will still work correctly. 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. (Books on category theory and type theory also do not mention applicative functors.)

20.2 AFTT is not category theory, type theory, or formal logic

So far it appears that AFTT includes a selection 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 presentation is too abstract and at the same time lacks specific results necessary for practical programming. 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 do not explain certain concepts relevant to 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 are: “Is X an introduction rule or an elimination rule” and “Does the property Y hold in non-small categories, or only in the category of sets”. Typical questions a Scala programmer might have are: “Can we compute a value of type `Either[Z, R => A]` from a value of type `R => Either[Z, A]`” and “Is the type constructor `F[A] = Option[(A, A, A)]` a monad or only an applicative functor”. 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 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 imple-

mentation of such foundationally important operations as `filter`, `zip`, and `traverse` – operations that functional programmers often use 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 hands-on (“run this code now and see what happens”) and narrow in scope, limited to one or two specific questions and specific applications. Such tutorials usually do not present sufficient mathematical insights to help programmers develop the necessary mathematical intuition.

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 programmer 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, for instance, cannot be used to check 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. Also, the Haskell Wikibooks tutorials rarely show any proofs or derivations of equational laws.

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 cer-

¹“What’s your problem?” as the joke would go.

tain 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 “reasonable”, 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` method 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

I mention these practical tasks as examples because they are actual real-world-coding applications of domain theory and the Curry-Howard correspondence theory. 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 equational 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 equational 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.

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 (tuple) 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 }`

$a \times b$ value of a tuple type; in Scala, `(a, b)`

id the identity function

$\mathbb{1}, 1$ the unit type and its value; in Scala, `Unit` and `()`

$\mathbb{0}$ 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

A Notations

pu_F the standard method `pure` of a monad F

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

$\mathbin{\circledast}$ the forward composition of functions: $f \mathbin{\circledast} g$ 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}

\triangleright use a value as the argument of a function: $x \triangleright f$ is $f(x)$

$f^{\uparrow G}$ a function f raised to a functor G ; same as `fmapG f`

$f^{\uparrow G \uparrow H}$ a function raised first to G and then to H ; in Scala, `h.map(_ .map(f))`

$f^{\downarrow H}$ a function f raised to a contrafunctor

\diamond_M the Kleisli product operation for the monad M

\oplus the binary operation of a monoid; in Scala, `x |+| y`

Δ the “diagonal” function of type $\forall A. A \Rightarrow A \times A$

$\nabla_1, \nabla_2, \dots$ the projections from a tuple to its first, second, ..., part

\boxtimes Cartesian product of functions, $(f \boxtimes g)(x) = f(x) \times g(x)$

$[a, b, c]$ a sequence of values; in Scala, `Seq(a, b, c)`

$\left\| \begin{array}{cc} \text{id} & \mathbb{0} \\ \mathbb{0} & a \Rightarrow a \times a \end{array} \right\|$ a function that works with disjunctive types

A.2 Detailed explanations

F^A means a type constructor F with a type parameter A . In Scala, this is `F[A]`. Type constructors with multiple type parameters are denoted by $F^{A,B,C}$.

$x^{:A}$ means a value x that has type A ; this is a **type annotation**. In Scala, a type annotation is `x:A`. The colon symbol, `:`, in the superscript shows that A is not a type argument (as it would be in a type constructor, F^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 from 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 from 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`.

$x^{:A} \Rightarrow y$ means a nameless function with argument x of type A and function body y . (Usually, the body y will be an expression that uses x .)

`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.

$\mathbb{1}$ means the unit type, and 1 means 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 $\mathbb{1} + A$, which in Scala corresponds to `Option[A]`.

$\mathbb{0}$ 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 $\mathbb{1} + A$ the non-empty part is $\mathbb{0} + A$, which in Scala corresponds to `Some[A]`. The empty part $\mathbb{1} + \mathbb{0}$ corresponds to `None`. Similarly, $A + \mathbb{0}$ denotes the left part of the type $A + B$ (in Scala, `Left[A]`), while $\mathbb{0} + B$ denotes its right part (in Scala, `Right[B]`). Values of disjunctive types are denoted similarly. For instance, $x^{:A} + \mathbb{0}^{:B}$ denotes a value of the left part of the type $A + B$; in Scala, this value is written as `Left[A,B](x)`.

\triangleq means “equal by definition”. Examples:

- $f \triangleq (x^{:Int} \Rightarrow x + 10)$ is a definition of the function f . In Scala, this is written as `val f = { x: Int => x + 10 }`
- $F^A \triangleq \mathbb{1} + A$ is a definition of a type constructor F . In Scala, this is written as `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

A Notations

allows nested tuples to be reordered whenever needed, we can write $(a \times b) \times c \cong a \times (b \times c)$.

A^{F^\bullet} 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 I denote as $F^A \triangleq \mathbb{1}^{F^{\text{Int}}} + A^{F^{A \times \text{String}}}$.

fmap_F means the standard method `fmap` of the `Functor` typeclass, implemented for the functor F . In Scala, this may be written as `Functor[F].fmap`. Since each functor F has its own specific implementation of fmap_F , the subscript “ F ” is not a type parameter of fmap_F . The method fmap_F actually has *two* type parameters, which can be written out as $\text{fmap}_F^{A,B}$. Then the type signature of fmap is written in full as $\text{fmap}_F^{A,B} : (A \Rightarrow B) \Rightarrow F^A \Rightarrow F^B$. For clarity, we 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, the presence of the type parameter A can be denoted by pu_F^A . If we are using the `pure` method with a complicated type, e.g. $\mathbb{1} + P^A$, instead of the type parameter A , we might want to write this type parameter for clarity and write $\text{pu}_F^{1+P^A}$. The type signature of that function is then

$$\text{pu}_F^{1+P^A} : \mathbb{1} + P^A \Rightarrow F^{1+P^A}.$$

But in most cases we will not need to write the type parameter.

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, $T_L^{M,\bullet}$

¹<https://github.com/typelevel/kind-projector>

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 $F \rightsquigarrow G$.

$\forall A. P^A$ 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.

\circ means the forward composition of functions: $f \circ g$ (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]]`.

$x \triangleright f$ means that x is inserted as the argument into the function f . So $x \triangleright f$ is the same as $f(x)$ or $f x$. In Scala, the expression $x \triangleright f$ is written as `x.f` and is the syntax used with many standard methods such as `.size` or `.toSeq`. Because the function f is to the *right* of x in this notation, we can chain forward compositions of functions as $x \triangleright f \triangleright g$ in a left-associative manner, similarly to how this is done in Scala code, for example `x.toSeq.sorted`. The operation \triangleright binds weaker than the forward composition \circ and so $x \triangleright f \circ g = x \triangleright f \triangleright g$ in this notation.

$f^{\uparrow G}$ means a function f raised 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 $f^{\uparrow G} g$ or $g \triangleright f^{\uparrow G}$. In Scala, this is `g.map(f)`. Nested raising (i.e. raising to the functor composition $H \circ G$) 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 raising to a value h of type H^{G^A} is written as $f^{\uparrow G \uparrow H} h$ or $h \triangleright f^{\uparrow G \uparrow H}$. 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 $f^{\downarrow H} h$ or $h \triangleright f^{\downarrow H}$, and yields a value of type 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 $A \Rightarrow M^C$.

\oplus means the binary operation of a monoid, for example $x \oplus y$. The specific

A Notations

monoid type should be defined for this expression to make sense. For example, in Scala the monoidal operation is usually denoted by $x \mid + \mid y$.

\boxtimes means the component-wise Cartesian product of functions, where the result is a pair of the values of the two functions: $(f \boxtimes g)(x) = f(x) \times g(x)$. In Scala, this operation can be defined by

```
def boxtimes[A,P,Q](f: A => P, g: A => Q): A => (P, Q) = x => (f(x), g(x))
```

$[a, b, c]$ means an ordered sequence of values, such as a list or an array. In Scala, this can be `List(a, b, c)`, `Vector(a, b, c)`, `Array(a, b, c)`, or another collection type.

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: Int) => x * 2`. Also known as function expression, function literal, anonymous function, closure, lambda-function, lambda-expression, or simply a “lambda”.

Contrafunctor A type constructor having the properties of a contravariant functor with respect to a type parameter. Instead of saying “contravariant functor”, I use the shorter name “contrafunctor”.

Profunctor A type constructor whose type parameter occurs in both covariant and contravariant positions.

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 Object Pascal and in OCaml). The shortest name is “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

B Glossary of terms

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 explicitly indexed type parameters, such as $F^A \triangleq (A \times A)^{:F^{\text{Int}}} + (\text{Int} \times A)^{:F^1}$. The Scala code for this type constructor is

```
sealed trait F[A]
final case class F1[A](x: A, y: A) extends F[Int]
final case class F2[A](s: Int, t: A) extends F[Unit]
```

Also called a **GADT** (generalized algebraic data type).

Functor block A short syntax for composing several `.map`, `.flatMap`, and `.filter` operations applied to a functor-typed value. The type constructor corresponding to that value must therefore be fixed throughout the entire functor block. (The type constructor *must* be a functor and may additionally be filterable and/or monadic.) 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, 6)`. This is a functor block that “raises” computations to the `List` functor. 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. (The type constructor does not have to)

Method This word is used in two ways: 1) A `method1` is a Scala function defined as a member of a typeclass. For example, `flatMap` is a method defined in the `Monad` typeclass. 2) A `method2` is a Scala function defined as a member of a data type declared as a Java-compatible `class` or `trait`. Trait `methods2` are necessary in Scala when implementing functions whose arguments have type parameters (because Scala function values defined via

B.1 On the current misuse of the term “algebra”

`val` cannot have type parameters). So, many typeclasses 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 `pureM` : $A \Rightarrow M^A$ 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 $A \Rightarrow M^C$).

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 “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 Functions of several arguments vs. tuples

In Scala, there is a difference between a function whose argument is a tuple,

```
def f1(p: (Int, String)): Int = ???
```

and a function with several arguments,

```
def f2(x: Int, s: String): Int = ???
```

This difference is “cosmetic” because these functions are equivalent from the computational point of view. However, Scala syntax for these functions is different.

C.0.3 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
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 propositional logic (IPL)

The intuitionistic propositional logic (sometimes also called the “constructive” propositional logic) describes how programs in functional programming languages may be able to compute values of different types.

The main formal difference between IPL and the classical (Boolean) logic is that IPL does not include the axiom of excluded middle (“*tertium non datur*”), which is

$$\forall A : (A \text{ or } (\text{not } (A))) \text{ is true} \quad .$$

However, given just this information, it is not easy to understand the consequences of *not having* this axiom, or to figure out which statements are true in the IPL.

The reason this axiom is not included in IPL is that IPL propositions such as $\mathcal{CH}(A)$ correspond to the *practical 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. But in classical logic, the axiom of excluded middle says that either $\mathcal{CH}(A)$ or $\text{not } (\mathcal{CH}(A))$ is true. So showing that “ $\text{not } (\mathcal{CH}(A))$ ” is contradictory is sufficient for proving $\mathcal{CH}(A)$, without ever computing any values of type A . For this reason, classical (Boolean) logic does not adequately describe the logic of types in functional programming, i.e. it does not correctly predict the types of values that can be computed by functional programs.

D.1 Example: The logic of types is not Boolean

Here is an explicit example of obtaining an incorrect result when using classical logic to reason about values computed by functional programs. Consider the formula

$$(A \Rightarrow B + C) \Rightarrow (A \Rightarrow B) + (A \Rightarrow C) \quad (\text{D.1})$$

D Intuitionistic propositional logic (IPL)

or, putting in all the parentheses for clarity,

$$(A \Rightarrow (B + C)) \Rightarrow ((A \Rightarrow B) + (A \Rightarrow C)) \quad .$$

This formula is a true theorem in classical logic. To prove this, we only need to show that Eq. (D.1) is always equal to *true* (i.e. Boolean value 1) for any Boolean values of the variables A, B, C . Consider that the only way an implication $A \Rightarrow B$ could be *false* (that is, equal to 0) in Boolean logic is when $A = 1$ and $B = 0$. So, Eq. (D.1) can be false only if $(A \Rightarrow B + C) = 1$ and $(A \Rightarrow B) + (A \Rightarrow C) = 0$. The disjunction can be false only when both parts are false; so we must have $(A \Rightarrow B) = 0$ and $(A \Rightarrow C) = 0$. This is only possible if $A = 1$ and $B = C = 0$. But, with these value assignments, we find $(A \Rightarrow B + C) = 0$ rather than 1. So, we cannot ever make Eq. (D.1) equal to 0 as a Boolean formula. This shows Eq. (D.1) to be a “classically valid” formula, i.e. a theorem that holds in classical Boolean logic.

If we use the Curry-Howard correspondence and apply Eq. (D.1) to propositions such as $\mathcal{CH}(A)$, $\mathcal{CH}(B)$, $\mathcal{CH}(C)$, we obtain the statement that a program should be able to compute a value of type $(A \Rightarrow B) + (A \Rightarrow C)$ given a value of type $A \Rightarrow B + C$. In Scala, such a program would be written as a function with the following type signature,

```
def bad[A, B, C](g: A => Either[B, C]): Either[A=>B, A=>C] = ???
```

However, it is impossible to implement this function in Scala as a total function.

To help build an intuition for the impossibility of implementing `bad`, consider that the only available data is a function $g : A \Rightarrow B + C$, which may return values of type B or C depending on the input value of type A . The function `bad` must return either a function of type $A \Rightarrow B$ or a function of type $A \Rightarrow C$. Can we create a function of type $A \Rightarrow B$? Given a value of type A , we need to compute a value of type B . Since the type B is completely arbitrary (it is a type parameter), we cannot produce a value of type B from scratch. The only potential source of values of type B is the input function g . However, g may produce values of type C for some values of type A . So, in general, we cannot build a function of type $A \Rightarrow B$ out of the function g . Similarly, we find that we cannot build a function of type $A \Rightarrow C$ out of g .

The decision about whether to return $A \Rightarrow B$ or $A \Rightarrow C$ must be somehow made in the code of `bad`. The only input data is the function g that takes an argument of type A . We could imagine calling g on various arguments of type A and to see whether g returns a B or a C . However, the type A is unknown, so the function `bad`

cannot produce any values of that type and call g . So the decision about whether to return $A \Rightarrow B$ or $A \Rightarrow C$ must be made regardless of the function g . Whichever we choose to return, $A \Rightarrow B$ or $A \Rightarrow C$, we will not be able to return a result value of the required type.

We could try to switch between $A \Rightarrow B$ and $A \Rightarrow C$ depending on a given value of type A . This, however, corresponds to a different type signature:

$$(A \Rightarrow B + C) \Rightarrow A \Rightarrow (A \Rightarrow B) + (A \Rightarrow C) \quad .$$

This type signature *can* be implemented, for instance, by this Scala code:

```
def q[A, B, C](g: A => Either[B, C]): A => Either[A=>B, A=>C] = { a =>
  g(a) match {
    case Left(b) => Left(_ => b)
    case Right(c) => Right(_ => c)
  }
}
```

But this is not the type signature that describes Eq. (D.1) via the Curry-Howard correspondence.

In the IPL, it turns out that Eq. (D.1) is not a valid theorem, i.e. it is impossible to find a proof of Eq. (D.1) using the axioms and the derivation rules of the IPL. To *prove* that there is no proof, one needs to use methods of proof theory that are beyond the scope of this book. A good introduction to the required technique is the book “*Proof and Disproof in Formal Logic*” by R. Bornat.¹

This example illustrates that it is precisely the valid theorems in the IPL, and not the valid theorems in the Boolean logic, that correspond to implementable functional programs.

D.2 Using truth values in Boolean logic and in IPL

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 more than two truth values.

¹ R. Bornat, “Proof and Disproof in Formal Logic”, Oxford, 2005 - [link to Amazon.com](https://www.amazon.com/dp/0199286410)

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 A humorous disclaimer

The following text is quoted in part from an anonymous source ("Project Guten Tag") dating back at least to 1997. The original text is no longer available on the Internet.

WARRANTO LIMITENSIS; DISCLAMATANTUS DAMAGENSIS

Solus exceptus "Rectum Replacator Refundiens" describitus ecci,

1. Projectus (etque nunquam partum quis hic etext remitibus cum PROJECT GUTEN TAG-tm identifier) 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 segundibus.

HIC ETEXT VENID "COMO-ASI". NIHIL WARRANTI NUNQUAM CLASSUM, EXPRESSITO NI IMPLICATO, LE MACCHEN COMO SI ETEXTO BENE SIT O IL MEDIO BENE SIT, INCLUTAT ET NON LIMITAT WARRANTI MERCATENSIS, APPROPRIATENSIS PURPOSEM.

Statuen varias non permitatent disclamabaris ni warranti implicatoren ni exclusioni limitatio damagaren consequentialis, ecco lo qua disclamatori exclusato-rique non vobis applicant, et potat optia alia legali.

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.