

J.N. OLIVEIRA

University of Minho

PROGRAM DESIGN BY CALCULATION

(DRAFT of textbook in preparation)

Last update: February 2019

CONTENTS

Preamble	1
1 INTRODUCTION	3
I CALCULATING WITH FUNCTIONS	4
2 AN INTRODUCTION TO POINTFREE PROGRAMMING	5
2.1 Introducing functions and types	6
2.2 Functional application	7
2.3 Functional equality and composition	7
2.4 Identity functions	10
2.5 Constant functions	10
2.6 Monics and epics	11
2.7 Isos	13
2.8 Gluing functions which do not compose — products	14
2.9 Gluing functions which do not compose — coproducts	20
2.10 Mixing products and coproducts	23
2.11 Elementary datatypes	25
2.12 Natural properties	27
2.13 Universal properties	29
2.14 Guards and McCarthy's conditional	32
2.15 Gluing functions which do not compose — exponentials	35
2.16 Finitary products and coproducts	42
2.17 Initial and terminal datatypes	43
2.18 Sums and products in HASKELL	45
2.19 Exercises	48
2.20 Bibliography notes	51
3 RECURSION IN THE POINTFREE STYLE	53
3.1 Motivation	53
3.2 From natural numbers to finite sequences	58
3.3 Introducing inductive datatypes	63
3.4 Observing an inductive datatype	68
3.5 Synthesizing an inductive datatype	71
3.6 Introducing (list) catas, anas and hylos	72
3.7 Inductive types more generally	77
3.8 Functors	78
3.9 Polynomial functors	79
3.10 Polynomial inductive types	81
3.11 F-algebras and F-homomorphisms	82
3.12 F-catamorphisms	83
3.13 Parameterization and type functors	85
3.14 A catalogue of standard polynomial inductive types	90
3.15 Hylo-factorization	92

3.16	Functors and type functors in HASKELL	96
3.17	The mutual-recursion law	97
3.18	“Banana-split”: a corollary of the mutual-recursion law	105
3.19	Inductive datatype isomorphism	106
3.20	Bibliography notes	106
4	WHY MONADS MATTER	108
4.1	Partial functions	108
4.2	Putting partial functions together	109
4.3	Lists	111
4.4	Monads	112
4.5	Monadic application (binding)	115
4.6	Sequencing and the do -notation	115
4.7	Generators and comprehensions	116
4.8	Monads in HASKELL	118
4.9	The state monad	121
4.10	‘Monadification’ of Haskell code made easy	127
4.11	Monadic recursion	131
4.12	Where do monads come from?	132
4.13	Bibliography notes	135
II	CALCULATING WITH RELATIONS	159
5	WHEN EVERYTHING BECOMES A RELATION	160
5.1	Functions are not enough	160
5.2	From functions to relations	162
5.3	Pre/post conditions	164
5.4	Relational composition and converse	166
5.5	Relational equality	169
5.6	Diagrams	171
5.7	Taxonomy of binary relations	173
5.8	Functions, relationally	178
5.9	Meet and join	181
5.10	Relational thinking	184
5.11	Monotonicity	186
5.12	Rules of thumb	188
5.13	Endo-relations	190
5.14	Relational pairing	193
5.15	Relational coproducts	196
5.16	On key-value data models	199
5.17	What about relational “currying”?	200
5.18	Galois connections	202
5.19	Relation division	208
5.20	Predicates also become relations	216
5.21	Guards, coreflexives and the McCarthy conditional	219
5.22	Difunctionality	223
5.23	Other orderings on relations	224
5.24	Back to functions	230

5.25 Bibliography notes	231
6 THEOREMS FOR FREE — BY CALCULATION	233
6.1 Introduction	233
6.2 Polymorphic type signatures	234
6.3 Relators	235
6.4 A relation on functions	236
6.5 Free theorem of type t	238
6.6 Examples	238
6.7 Catamorphism laws as free theorems	242
6.8 Bibliography notes	244
7 CONTRACT-ORIENTED PROGRAMMING	245
7.1 Contracts	246
7.2 Library loan example	249
7.3 Mobile phone example	253
7.4 A calculus of functional contracts	255
7.5 Abstract interpretation	258
7.6 Safety and liveness properties	261
7.7 Examples	262
7.8 “Free contracts”	269
7.9 Reasoning by approximation	269
7.10 Bibliography notes	270
8 PROGRAMS AS RELATIONAL HYLOMORPHISMS	271
9 CALCULATIONAL PROGRAM REFINEMENT	272
III CALCULATING WITH MATRICES	273
10 TOWARDS A LINEAR ALGEBRA OF PROGRAMMING	274
A BACKGROUND — EINDHOVEN QUANTIFIER CALCULUS	275
A.1 Notation	275
A.2 Rules	275
B HASKELL SUPPORT LIBRARY	277

LIST OF EXERCISES

Exercise 2.1	11
Exercise 2.2	11
Exercise 2.3	13
Exercise 2.4	20
Exercise 2.5	23
Exercise 2.6	23
Exercise 2.7	24
Exercise 2.8	25
Exercise 2.9	25
Exercise 2.10	25
Exercise 2.11	27
Exercise 2.12	27
Exercise 2.13	28
Exercise 2.14	28
Exercise 2.15	29
Exercise 2.16	31
Exercise 2.17	31
Exercise 2.18	31
Exercise 2.19	31
Exercise 2.20	31
Exercise 2.21	34
Exercise 2.22	34
Exercise 2.23	35
Exercise 2.24	35
Exercise 2.25	41
Exercise 2.26	41
Exercise 2.27	41
Exercise 2.28	42
Exercise 2.29	42
Exercise 2.30	43
Exercise 2.31	44
Exercise 2.32	48
Exercise 2.33	48
Exercise 2.34	48
Exercise 2.35	48
Exercise 2.36	49
Exercise 2.37	49
Exercise 2.38	49
Exercise 2.39	49
Exercise 2.40	50
Exercise 2.41	50

Exercise 2.42	50
Exercise 2.43	50
Exercise 2.44	50
Exercise 3.1	57
Exercise 3.2	57
Exercise 3.3	57
Exercise 3.4	58
Exercise 3.5	58
Exercise 3.6	68
Exercise 3.7	76
Exercise 3.8	76
Exercise 3.9	76
Exercise 3.10	81
Exercise 3.11	89
Exercise 3.12	89
Exercise 3.13	90
Exercise 3.14	91
Exercise 3.15	91
Exercise 3.16	91
Exercise 3.17	92
Exercise 3.18	94
Exercise 3.19	94
Exercise 3.20	94
Exercise 3.21	95
Exercise 3.22	95
Exercise 3.23	95
Exercise 3.24	97
Exercise 3.25	97
Exercise 3.26	101
Exercise 3.27	101
Exercise 3.28	102
Exercise 3.29	102
Exercise 3.30	102
Exercise 3.31	103
Exercise 3.32	104
Exercise 3.33	106
Exercise 3.34	106
Exercise 3.35	106
Exercise 4.1	111
Exercise 4.2	111
Exercise 4.3	112
Exercise 4.4	114
Exercise 4.5	117
Exercise 4.6	118
Exercise 4.7	119
Exercise 4.8	120

Exercise 4.9	120
Exercise 4.10	131
Exercise 4.11	131
Exercise 4.12	135
Exercise 5.1	168
Exercise 5.2	168
Exercise 5.3	173
Exercise 5.4	176
Exercise 5.5	177
Exercise 5.6	177
Exercise 5.7	177
Exercise 5.8	177
Exercise 5.9	178
Exercise 5.10	179
Exercise 5.11	180
Exercise 5.12	181
Exercise 5.13	181
Exercise 5.14	183
Exercise 5.15	183
Exercise 5.16	183
Exercise 5.17	183
Exercise 5.18	183
Exercise 5.19	186
Exercise 5.20	186
Exercise 5.21	187
Exercise 5.22	187
Exercise 5.23	189
Exercise 5.24	189
Exercise 5.25	191
Exercise 5.26	192
Exercise 5.27	192
Exercise 5.28	192
Exercise 5.29	192
Exercise 5.30	192
Exercise 5.31	192
Exercise 5.32	193
Exercise 5.33	195
Exercise 5.34	195
Exercise 5.35	196
Exercise 5.36	198
Exercise 5.37	198
Exercise 5.38	199
Exercise 5.39	199
Exercise 5.40	202
Exercise 5.41	206
Exercise 5.42	206

Exercise 5.43	207
Exercise 5.44	207
Exercise 5.45	208
Exercise 5.46	211
Exercise 5.47	215
Exercise 5.48	215
Exercise 5.49	215
Exercise 5.50	219
Exercise 5.51	220
Exercise 5.52	222
Exercise 5.53	226
Exercise 5.54	227
Exercise 5.55	227
Exercise 5.56	227
Exercise 5.57	228
Exercise 5.58	228
Exercise 5.59	230
Exercise 6.1	235
Exercise 6.2	240
Exercise 6.3	240
Exercise 6.4	241
Exercise 6.5	241
Exercise 6.6	241
Exercise 6.7	241
Exercise 6.8	242
Exercise 6.9	242
Exercise 6.10	243
Exercise 6.11	243
Exercise 6.12	244
Exercise 7.1	247
Exercise 7.2	248
Exercise 7.3	254
Exercise 7.4	257
Exercise 7.5	259
Exercise 7.6	259
Exercise 7.7	259
Exercise 7.8	260
Exercise 7.9	268
Exercise 7.10	269
Exercise 7.11	269
Exercise 7.12	269

PREAMBLE

This textbook, which has arisen from the author's research and teaching experience, has been in preparation for many years. Its main aim is to draw the attention of software practitioners to a calculational approach to the design of software artifacts ranging from simple algorithms and functions to the specification and realization of information systems.

Put in other words, the book invites software designers to raise standards and adopt mature development techniques found in other engineering disciplines, which (as a rule) are rooted on a sound mathematical basis. *Compositionality* and *parametricity* are central to the whole discipline, granting scalability from school desk exercises to large problems in an industry setting.

It is interesting to note that while coining the phrase *software engineering* in the 1960s, our colleagues of the time were already promising such high quality standards. In March, 1967, ACM President Anthony Oettinger delivered an address in which he said [48]:

*"(...) the scientific, rigorous component of computing, is more like **mathematics** than it is like **physics**" (...) Whatever it is, on the one hand it has components of the purest of mathematics and on the other hand of the dirtiest of engineering.*

As a discipline, software engineering was announced at the Garmisch NATO conference in 1968, from whose report [46] the following excerpt is quoted:

In late 1967 the Study Group recommended the holding of a working conference on Software Engineering. The phrase 'software engineering' was deliberately chosen as being provocative, in implying the need for software manufacture to be based on the types of theoretical foundations and practical disciplines, that are traditional in the established branches of engineering.

Provocative or not, the need for sound theoretical foundations has clearly been under concern since the very beginning of the discipline — exactly fifty years ago, at the time of writing. However, how "scientific" do such foundations turn out to be, now that five decades have since elapsed?¹

Thirty years later (1997), Richard Bird and Oege de Moore published a textbook [10] in the preface of which C.A.R. Hoare writes:

*Programming notation can be expressed by "**formulae** and **equations** (...) which share the **elegance** of those which underlie **physics** and **chemistry** or any other branch of basic science".*

¹ The title of a communication of another ACM President, Vinton Cerf (2012), does not sound particularly optimistic [12].

The formulæ and equations mentioned in this quotation are those of a discipline known as the *Algebra of Programming*. Many others have contributed to this body of knowledge, notably Roland Backhouse and his colleagues at Eindhoven and Nottingham, see eg. [1, 4], Jeremy Gibbons and Ralf Hinze at Oxford see e.g. [23], among many others. Unfortunately, references [1, 4] are still unpublished.

When the author of this draft textbook decided to teach *Algebra of Programming* to 2nd year students of the Minho degrees in computer science, back to 1998, he found textbook [10] too difficult for the students to follow, mainly because of its too explicit categorial (allegorical) flavour. So he decided to start writing slides and notes helping the students to read the book. Eventually, such notes became chapters 2 to 4 of the current version of the monograph. The same procedure was taken when teaching the relational approach of [10] to 4th year students (master level), see chapters 5 to 7.

This draft book is incomplete, all subsequent chapters being still in *slide form*². Such half-finished chapters are omitted from the current print-out. Altogether, the idea is to show that software engineering and, in particular, computer programming can adopt the *scientific method* as other branches of engineering do. Somehow, it's like following in the footsteps of those who marveled at the power of algebraic reasoning in the remote past, in different contexts and disciplines [47]:

“(...) *De manera, que quien sabe por Algebra, sabe científicamente* [In this way, who knows by Algebra knows scientifically].

University of Minho, Braga, February 2019



José N. Oliveira

² See e.g. see technical report [51]. The third part will address a linear algebra of programming intended for quantitative reasoning about software. This is even less stable, but a number of papers exist already about the topic, starting from [50].

INTRODUCTION

Not given in the current version of this textbook.

Part I

CALCULATING WITH FUNCTIONS

AN INTRODUCTION TO POINTFREE PROGRAMMING

Everybody is familiar with the concept of a *function* since the school desk. The functional intuition traverses mathematics from end to end because it has a solid semantics rooted on a well-known mathematical system — the class of “all” sets and set-theoretical functions.

Functional programming literally means “programming with functions”. Programming languages such as LISP or HASKELL allow us to program with functions. However, the functional intuition is far more reaching than producing code which runs on a computer. Since the pioneering work of John McCarthy — the inventor of LISP — in the early 1960s, one knows that other branches of programming can be structured, or expressed functionally. The idea of producing programs by *calculation*, that is to say, that of calculating efficient programs out of abstract, inefficient ones has a long tradition in functional programming.

This book is structured around the idea that functional programming can be used as a basis for teaching programming as a whole, from the successor function $n \mapsto n + 1$ to large information system design.¹

This chapter provides a light-weight introduction to the theory of functional programming. The main emphasis is on *compositionality* — one of the main advantages of “thinking functionally” — by explaining how to construct new functions out of other functions using a minimal set of predefined functional *combinators*. This leads to a programming style that is *point free* in the sense that function descriptions dispense with variables (also known as *points*).

Several technical issues are deliberately ignored and deferred to later chapters. Most programming examples will be provided in the HASKELL functional programming language. Appendix B includes the listings of some HASKELL modules that complement the HASKELL *Standard Prelude* and help to “animate” the main concepts introduced in this chapter.

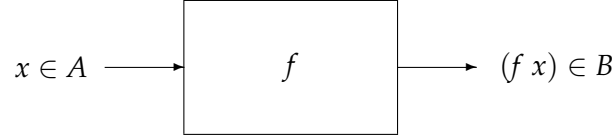
¹ This idea addresses programming in a broad sense, including for instance *reversible* and *quantum programming*, where functional programming already plays leading roles [44, 42, 22].

2.1 INTRODUCING FUNCTIONS AND TYPES

The definition of a function

$$f : A \rightarrow B \quad (2.1)$$

can be regarded as a kind of “process” abstraction: it is a “black box” which produces an output once it is supplied with an input:



The box isn’t really necessary to convey the abstraction, a single labelled arrow sufficing:

$$A \xrightarrow{f} B$$

This simplified notation focusses on what is indeed relevant about f — that it can be regarded as a kind of “contract”:

f commits itself to producing a B-value provided it is supplied with an A-value.

How is such a value produced? In many situations one wishes to ignore it because one is just *using* function f . In others, however, one may want to inspect the internals of the “black box” in order to know the function’s *computation rule*. For instance,

$$\begin{aligned} \text{succ} & : \mathbb{N} \rightarrow \mathbb{N} \\ \text{succ } n & \stackrel{\text{def}}{=} n + 1 \end{aligned}$$

expresses the computation rule of the *successor* function — the function succ which finds “the next natural number” — in terms of natural number addition and of natural number 1. What we above meant by a “contract” corresponds to the *signature* of the function, which is expressed by arrow $\mathbb{N} \rightarrow \mathbb{N}$ in the case of succ and which, by the way, can be shared by other functions, *e.g.* $\text{sq } n \stackrel{\text{def}}{=} n^2$.

In programming terminology one says that succ and sq have the same “type”. Types play a prominent rôle in functional programming (as they do in other programming paradigms). Informally, they provide the “glue”, or interfacing material, for putting functions together to obtain more complex functions. Formally, a “type checking” discipline can be expressed in terms of compositional rules which check for functional expression wellformedness.

It has become standard to use arrows to denote function signatures or function types, recall (2.1). To denote the fact that function f accepts arguments of type A and produces results of type B , we will use the following interchangeable notations: $f : B \leftarrow A$, $f : A \rightarrow B$, $B \xleftarrow{f} A$

or $A \xrightarrow{f} B$. This corresponds to writing $f :: a \rightarrow b$ in the HASKELL functional programming language, where type variables are denoted by lowercase letters. A will be referred to as the *domain* of f and B will be referred to as the *codomain* of f . Both A and B are symbols or expressions which denote sets of values, most often called *types*.

2.2 FUNCTIONAL APPLICATION

What do we want functions for? If we ask this question to a physician or engineer the answer is very likely to be: one wants functions for modelling and reasoning about the behaviour of real things.

For instance, function $distance\ t = 60 \times t$ could be written by a school physics student to model the distance (in, say, kilometers) a car will drive (per hour) at average speed $60km/hour$. When questioned about how far the car has gone in 2.5 hours, such a model provides an immediate answer: just evaluate $distance\ 2.5$ to obtain $150km$.

So we get a naïve purpose of functions: we want them to be *applied* to arguments in order to obtain results. Functional *application* is denoted by juxtaposition, e.g. $f\ a$ for $B \xleftarrow{f} A$ and $a \in A$, and associates to the left: $f\ x\ y$ denotes $(f\ x)\ y$ rather than $f\ (x\ y)$.

2.3 FUNCTIONAL EQUALITY AND COMPOSITION

Application is not everything we want to do with functions. Very soon our physics student will be able to talk about properties of the *distance* model, for instance that property

$$distance\ (2 \times t) = 2 \times (distance\ t) \quad (2.2)$$

holds. Later on, we could learn from her or him that the same property can be restated as $distance\ (twice\ t) = twice\ (distance\ t)$, by introducing function $twice\ x \stackrel{\text{def}}{=} 2 \times x$. Or even simply as

$$distance \cdot twice = twice \cdot distance \quad (2.3)$$

where “.” denotes function-arrow chaining, as suggested by drawing

$$\begin{array}{ccc} \mathbb{R} & \xleftarrow{twice} & \mathbb{R} \\ distance \downarrow & & \downarrow distance \\ \mathbb{R} & \xleftarrow{twice} & \mathbb{R} \end{array} \quad (2.4)$$

where both space and time are modelled by real numbers in \mathbb{R} .

This trivial example illustrates some relevant facets of the functional programming paradigm. Which version of the property presented above is “better”? the version explicitly mentioning variable t and requiring parentheses (2.2)? the version hiding variable t but resorting to function $twice$ (2.3)? or even diagram (2.4) alone?

Expression (2.3) is clearly more compact than (2.2). The trend for notation economy and compactness is well-known throughout the history of mathematics. In the 16th century, for instance, algebrists would write $12.cu.\tilde{p}.18.ce.\tilde{p}.27.co.\tilde{p}.17$ for what is nowadays written as $12x^3 + 18x^2 + 27x + 17$. We may find such *syncopated* notation odd, but we should not forget that at its time it was replacing even more obscure and lengthy expression denotations.

Why do people look for compact notations? A compact notation leads to shorter documents (less lines of code in programming) in which patterns are easier to identify and to reason about. Properties can be stated in clear-cut, one-line long equations which are easy to memorize. And diagrams such as (2.4) can be easily drawn which enable us to visualize maths in a graphical format.

Some people will argue that such compact “pointfree” notation (that is, the notation which hides variables, or function “definition points”) is too cryptic to be useful as a practical programming medium. In fact, pointfree programming languages such as Iverson’s APL or Backus’ FP have been more respected than loved by the programmers community. Virtually all commercial programming languages require variables and so implement the more traditional “pointwise” notation.

Throughout this book we will adopt both, depending upon the context. Our chosen programming medium — HASKELL — blends the pointwise and pointfree programming styles in a quite successful way. In order to switch from one to the other, we need two “bridges”: one lifting equality to the functional level and the other lifting function application.

Concerning equality, note that the “=” sign in (2.2) differs from that in (2.3): while the former states that two real numbers are the same number, the latter states that two $\mathbb{R} \leftarrow \mathbb{R}$ functions are the same function. Formally, we will say that two functions $f, g : B \leftarrow A$ are equal if they agree at pointwise-level, that is²

$$f = g \text{ iff } \langle \forall a : a \in A : f a =_B g a \rangle \quad (2.5)$$

where $=_B$ denotes equality at B -level. Rule (2.5) is known as *extensional equality*.

Concerning application, the pointfree style replaces it by the more generic concept of functional *composition* suggested by function-arrow chaining: wherever two functions are such that the target type of one of them, say $B \xleftarrow{g} A$ is the same as the source type of the other, say $C \xleftarrow{f} B$, then another function can be defined, $C \xleftarrow{f \cdot g} A$ — called the *composition* of f and g , or “ f after g ” — which “glues” f and g together:

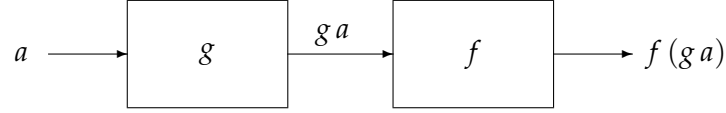
$$(f \cdot g) a \stackrel{\text{def}}{=} f (g a) \quad (2.6)$$

² Quantified notation $\langle \forall x : P : Q \rangle$ means: “for all x in the range P , Q holds”, where P and Q are logical expressions involving x . (See appendix A.) This notation will be used sporadically in the first part of this book.

This situation is pictured by the following arrow-diagram

$$\begin{array}{ccc}
 B & \xleftarrow{g} & A \\
 f \downarrow & \swarrow f \cdot g & \\
 C & &
 \end{array}
 \quad (2.7)$$

or by block-diagram



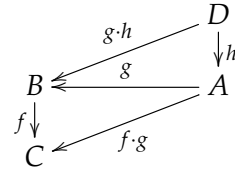
Therefore, the type-rule associated to functional composition can be expressed as follows:³

$$\frac{
 \begin{array}{c}
 C \xleftarrow{f} B \\
 B \xleftarrow{g} A
 \end{array}
 }{
 C \xleftarrow{f \cdot g} A
 }$$

Composition is certainly the most basic of all functional combinators. It is the first kind of “glue” which comes to mind when programmers need to combine, or chain functions (or processes) to obtain more elaborate functions (or processes).⁴ This is because of one of its most relevant properties,

$$(f \cdot g) \cdot h = f \cdot (g \cdot h) \quad (2.8)$$

depicted by diagram



which shares the pattern of, for instance

$$(a + b) + c = a + (b + c)$$

and so is called the *associative* property of composition. This enables us to move parentheses around in pointfree expressions involving functional compositions, or even to omit them altogether, for instance by writing $f \cdot g \cdot h \cdot i$ as an abbreviation of $((f \cdot g) \cdot h) \cdot i$, or of $(f \cdot (g \cdot h)) \cdot i$, or of $f \cdot ((g \cdot h) \cdot i)$, etc. For a chain of n -many function compositions the notation $\bigcirc_{i=1}^n f_i$ will be acceptable as abbreviation of $f_1 \cdot \dots \cdot f_n$.

³ This and other type-rules to come adopt the usual “fractional” layout, reminiscent of that used in school arithmetics for addition, subtraction, etc.

⁴ It even has a place in scripting languages such as UNIX’s shell, where $f \mid g$ is the shell counterpart of $g \cdot f$, for appropriate “processes” f and g .

2.4 IDENTITY FUNCTIONS

How free are we to fulfill the “give me an A and I will give you a B ” contract of equation (2.1)? In general, the choice of f is not unique. Some f s will do as little as possible while others will laboriously compute non-trivial outputs. At one of the extremes, we find functions which “do nothing” for us, that is, the added-value of their output when compared to their input amounts to nothing: $f a = a$. In this case $B = A$, of course, and f is said to be the *identity* function on A :

$$\begin{aligned} id_A &: A \leftarrow A \\ id_A a &\stackrel{\text{def}}{=} a \end{aligned} \quad (2.9)$$

Note that every type X “has” its identity id_X . Subscripts will be omitted wherever implicit in the context. For instance, the arrow notation $\mathbb{N} \xleftarrow{id} \mathbb{N}$ saves us from writing $id_{\mathbb{N}}$. So, we will often refer to “the” identity function rather than to “an” identity function.

How useful are identity functions? At first sight, they look fairly uninteresting. But the interplay between composition and identity, captured by the following equation,

$$f \cdot id = id \cdot f = f \quad (2.10)$$

will be appreciated later on. This property shares the pattern of, for instance,

$$a + 0 = 0 + a = a$$

This is why we say that id is the *unit* (*identity*) of composition. In a diagram, (2.10) looks like this:

$$\begin{array}{ccc} A & \xleftarrow{id} & A \\ f \downarrow & & \downarrow f \\ B & \xleftarrow{id} & B \end{array} \quad (2.11)$$

Note the graphical analogy of diagrams (2.4) and (2.11). The latter is interesting in the sense that it is *generic*, holding for every f . Diagrams of this kind are very common and express important (and rather ‘natural’) properties of functions, as we shall see further on.

2.5 CONSTANT FUNCTIONS

Opposite to the identity functions, which do not lose any information, we find functions which lose all (or almost all) information. Regardless of their input, the output of these functions is always the same value.

Let C be a nonempty data domain and let $c \in C$. Then we define the *everywhere* c function as follows, for arbitrary A :

$$\begin{array}{lcl} \underline{c} & : & A \rightarrow C \\ \underline{c} a & \stackrel{\text{def}}{=} & c \end{array} \quad (2.12)$$

The following property defines constant functions at pointfree level,

$$\underline{c} \cdot f = \underline{c} \quad (2.13)$$

and is depicted by a diagram similar to (2.11):

$$\begin{array}{ccc} C & \xleftarrow{\underline{c}} & A \\ \text{id} \downarrow & & \downarrow f \\ C & \xleftarrow{\underline{c}} & B \end{array} \quad (2.14)$$

Clearly, $\underline{c} \cdot f = \underline{c} \cdot g$, for any f, g , meaning that any difference that may exist in behaviour between such functions is lost.

Note that, strictly speaking, symbol \underline{c} denotes two different functions in diagram (2.14): one, which we should have written \underline{c}_A , accepts inputs from A while the other, which we should have written \underline{c}_B , accepts inputs from B :

$$\underline{c}_B \cdot f = \underline{c}_A \quad (2.15)$$

This property will be referred to as the *constant-fusion* property.

As with identity functions, subscripts will be omitted wherever implicit in the context.

Exercise 2.1. Use (2.5) to show that $f \cdot h = h \cdot f = f$ has the unique solution $h = \text{id}$, cf. (2.10).

□

Exercise 2.2. The HASKELL Prelude provides for constant functions: you write `const c` for \underline{c} . Check that HASKELL assigns the same type to expressions $f \cdot (\text{const } c)$ and `const (f c)`, for every f and c . What else can you say about these functional expressions? Justify.

□

2.6 MONICS AND EPICS

Identity functions and constant functions are limit points of the functional spectrum with respect to information preservation. All the other functions are in between: they lose “some” information, which is regarded as uninteresting for some reason. This remark supports the

following aphorism about a facet of functional programming: it is the *art* of transforming or losing information in a controlled and precise way. That is to say, the art of constructing the exact observation of data which fits in a particular context or requirement.

How do functions lose information? Basically in two different ways: they may be “blind” enough to confuse different inputs, by mapping them onto the same output, or they may ignore values of their codomain. For instance, \underline{c} confuses *all* inputs by mapping them all onto c . Moreover, it ignores all values of its codomain apart from c .

Functions which do not confuse inputs are called *monics* (or *injective* functions) and obey the following property: $B \xleftarrow{f} A$ is *monic* if, for every pair of functions $A \xleftarrow{h,k} C$, if $f \cdot h = f \cdot k$ then $h = k$, cf. diagram

$$B \xleftarrow{f} A \xleftarrow[h]{h} C$$

(we say that f is “post-cancellable”). It is easy to check that “the” identity function is monic,

$$\begin{aligned} & id \cdot h = id \cdot k \Rightarrow h = k \\ \equiv & \quad \{ \text{by (2.10)} \} \\ & h = k \Rightarrow h = k \\ \equiv & \quad \{ \text{predicate logic} \} \\ & \text{TRUE} \end{aligned}$$

and that any constant function \underline{c} is not monic:

$$\begin{aligned} & \underline{c} \cdot h = \underline{c} \cdot k \Rightarrow h = k \\ \equiv & \quad \{ \text{by (2.15)} \} \\ & \underline{c} = \underline{c} \Rightarrow h = k \\ \equiv & \quad \{ \text{function equality is reflexive} \} \\ & \text{TRUE} \Rightarrow h = k \\ \equiv & \quad \{ \text{predicate logic} \} \\ & h = k \end{aligned}$$

So the implication does not hold in general (only if $h = k$).

Functions which do not ignore values of their codomain are called *epics* (or *surjective* functions) and obey the following property: $A \xleftarrow{f} B$ is *epic* if, for every pair of functions $C \xleftarrow{h,k} A$, if $h \cdot f = k \cdot f$ then $h = k$, cf. diagram

$$C \xleftarrow[h]{h} A \xleftarrow{f} B$$

(we say that f is “pre-cancellable”). As expected, identity functions are epic:

$$\begin{aligned}
 & h \cdot id = k \cdot id \Rightarrow h = k \\
 \equiv & \quad \{ \text{by (2.10)} \} \\
 & h = k \Rightarrow h = k \\
 \equiv & \quad \{ \text{predicate logic} \} \\
 & \text{TRUE}
 \end{aligned}$$

Exercise 2.3. Under what circumstances is a constant function epic? Justify.

□

2.7 ISOS

A function $B \xleftarrow{f} A$ which is both monic and epic is said to be *iso* (an isomorphism, or a bijective function). In this situation, f always has a *converse* (or *inverse*) $B \xrightarrow{f^\circ} A$, which is such that

$$f \cdot f^\circ = id_B \quad \wedge \quad f^\circ \cdot f = id_A \quad (2.16)$$

(i.e. f is *invertible*).

Isomorphisms are very important functions because they convert data from one “format”, say A , to another format, say B , without losing information. So f and f° are faithful protocols between the two formats A and B . Of course, these formats contain the same “amount” of information, although the same data adopts a different “shape” in each of them. In mathematics, one says that A is *isomorphic* to B and one writes $A \cong B$ to express this fact.

Isomorphic data domains are regarded as “abstractly” the same. Note that, in general, there is a wide range of isos between two isomorphic data domains. For instance, let *Weekday* be the set of weekdays,

$$\begin{aligned}
 \text{Weekday} = & \\
 & \{ \text{Sunday, Monday, Tuesday, Wednesday, Thursday, Friday, Saturday} \}
 \end{aligned}$$

and let symbol 7 denote the set $\{1, 2, 3, 4, 5, 6, 7\}$, which is the *initial segment* of \mathbb{N} containing exactly seven elements. The following function f , which associates each weekday with its “ordinal” number,

$$\begin{aligned}
 f : \text{Weekday} & \rightarrow 7 \\
 f \text{ Monday} & = 1 \\
 f \text{ Tuesday} & = 2 \\
 f \text{ Wednesday} & = 3 \\
 f \text{ Thursday} & = 4
 \end{aligned}$$

$$\begin{aligned} f \text{ Friday} &= 5 \\ f \text{ Saturday} &= 6 \\ f \text{ Sunday} &= 7 \end{aligned}$$

is iso (guess f°). Clearly, $f d = i$ means “ d is the i -th day of the week”. But note that function $g d \stackrel{\text{def}}{=} \text{rem}(f d, 7) + 1$ is also an iso between Weekday and 7. While f regards *Monday* the first day of the week, g places *Sunday* in that position. Both f and g are witnesses of isomorphism

$$\text{Weekday} \cong 7 \quad (2.17)$$

Isomorphisms are quite flexible in pointwise reasoning. If, for some reason, f° is found handier than isomorphism f in the reasoning, then the shunting rules

$$f \cdot g = h \equiv g = f^\circ \cdot h \quad (2.18)$$

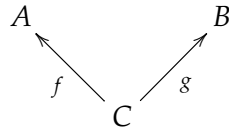
$$g \cdot f = h \equiv g = h \cdot f^\circ \quad (2.19)$$

can be of help.

Finally, note that all classes of functions referred to so far — constants, identities, epics, monics and isos — are closed under composition, that is, the composition of two constants is a constant, the composition of two epics is epic, *etc.*

2.8 GLUING FUNCTIONS WHICH DO NOT COMPOSE — PRODUCTS

Function composition has been presented above as a basis for gluing functions together in order to build more complex functions. However, not every two functions can be glued together by composition. For instance, functions $f : A \leftarrow C$ and $g : B \leftarrow C$ do not compose with each other because the domain of one of them is not the codomain of the other. However, both f and g share the same domain C . So, something we can do about gluing f and g together is to draw a diagram expressing this fact, something like



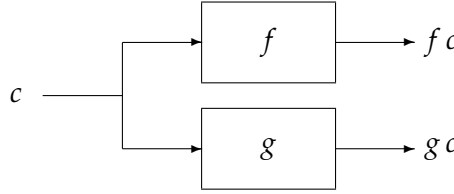
Because f and g share the same domain, their outputs can be paired, that is, we may write ordered pair $(f c, g c)$ for each $c \in C$. Such pairs belong to the Cartesian product of A and B , that is, to the set

$$A \times B \stackrel{\text{def}}{=} \{(a, b) \mid a \in A \wedge b \in B\}$$

So we may think of the operation which pairs the outputs of f and g as a new function combinator $\langle f, g \rangle$ defined as follows:

$$\begin{aligned} \langle f, g \rangle &: C \rightarrow A \times B \\ \langle f, g \rangle c &\stackrel{\text{def}}{=} (f c, g c) \end{aligned} \quad (2.20)$$

Traditionally, the pairing combinator $\langle f, g \rangle$ is pronounced “*f split g*” (or “pair *f* and *g*”) and can be depicted by the following “block”, or “data flow” diagram:



Function $\langle f, g \rangle$ keeps the information of both f and g in the same way Cartesian product $A \times B$ keeps the information of A and B . So, in the same way A data or B data can be retrieved from $A \times B$ data via the implicit *projections* π_1 or π_2 ,

$$A \xleftarrow{\pi_1} A \times B \xrightarrow{\pi_2} B \quad (2.21)$$

defined by

$$\pi_1(a, b) = a \quad \text{and} \quad \pi_2(a, b) = b$$

f and g can be retrieved from $\langle f, g \rangle$ via the same projections:

$$\pi_1 \cdot \langle f, g \rangle = f \quad \text{and} \quad \pi_2 \cdot \langle f, g \rangle = g \quad (2.22)$$

This fact (or pair of facts) will be referred to as the \times -*cancellation* property and is illustrated in the following diagram which puts everything together:

$$\begin{array}{ccccc}
 A & \xleftarrow{\pi_1} & A \times B & \xrightarrow{\pi_2} & B \\
 & \swarrow f & \uparrow \langle f, g \rangle & \searrow g & \\
 & & C & &
 \end{array} \quad (2.23)$$

In summary, the type-rule associated to the “split” combinator is expressed by

$$\frac{
 \begin{array}{c}
 A \xleftarrow{f} C \\
 B \xleftarrow{g} C
 \end{array}
 }{
 A \times B \xleftarrow{\langle f, g \rangle} C
 }$$

A *split* arises wherever two functions do not compose but share the same domain. What about gluing two functions which fail such a requisite, e.g.

$$\frac{
 \begin{array}{c}
 A \xleftarrow{f} C \\
 B \xleftarrow{g} D
 \end{array}
 }{
 \dots?
 }$$

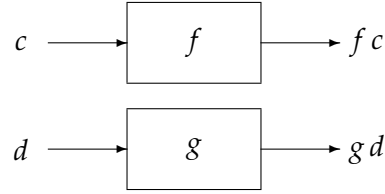
The $\langle f, g \rangle$ *split* combination does not work any more. Nevertheless, a way to “bridge” the domains of f and g , C and D respectively, is to regard them as targets of the projections π_1 and π_2 of $C \times D$:

$$\begin{array}{ccccc} A & \xleftarrow{\pi_1} & A \times B & \xrightarrow{\pi_2} & B \\ f \uparrow & & & & \uparrow g \\ C & \xleftarrow{\pi_1} & C \times D & \xrightarrow{\pi_2} & D \end{array}$$

From this diagram $\langle f \cdot \pi_1, g \cdot \pi_2 \rangle$ arises

$$\begin{array}{ccccc} A & \xleftarrow{\pi_1} & A \times B & \xrightarrow{\pi_2} & B \\ & \searrow f \cdot \pi_1 & \uparrow \langle f \cdot \pi_1, g \cdot \pi_2 \rangle & \nearrow g \cdot \pi_2 & \\ & & C \times D & & \end{array}$$

mapping $C \times D$ to $A \times B$. It corresponds to the “parallel” application of f and g which is suggested by the following data-flow diagram:



Functional combination $\langle f \cdot \pi_1, g \cdot \pi_2 \rangle$ appears so often that it deserves special notation — it will be expressed by $f \times g$. So, by definition, we have

$$f \times g \stackrel{\text{def}}{=} \langle f \cdot \pi_1, g \cdot \pi_2 \rangle \quad (2.24)$$

which is pronounced “product of f and g ” and has typing-rule

$$\frac{\begin{array}{c} A \xleftarrow{f} C \\ B \xleftarrow{g} D \end{array}}{A \times B \xleftarrow{f \times g} C \times D} \quad (2.25)$$

Note the overloading of symbol “ \times ”, which is used to denote both Cartesian product and functional product. This choice of notation will be fully justified later on.

What is the interplay among functional combinators $f \cdot g$ (composition), $\langle f, g \rangle$ (*split*) and $f \times g$ (product)? Composition and *split* relate to each other via the following property, known as \times -fusion:⁵

$$\begin{array}{c}
 A \xleftarrow{\pi_1} A \times B \xrightarrow{\pi_2} B \\
 \swarrow g \quad \searrow h \\
 \quad C \\
 \swarrow g \cdot f \quad \searrow h \cdot f \\
 \quad D
 \end{array}
 \quad \langle g, h \rangle \cdot f = \langle g \cdot f, h \cdot f \rangle \quad (2.26)$$

This shows that *split* is right-distributive with respect to composition. Left-distributivity does not hold but there is something we can say about $f \cdot \langle g, h \rangle$ in case $f = i \times j$:

$$\begin{aligned}
 & (i \times j) \cdot \langle g, h \rangle \\
 = & \quad \{ \text{by (2.24)} \} \\
 & \langle i \cdot \pi_1, j \cdot \pi_2 \rangle \cdot \langle g, h \rangle \\
 = & \quad \{ \text{by } \times\text{-fusion (2.26)} \} \\
 & \langle (i \cdot \pi_1) \cdot \langle g, h \rangle, (j \cdot \pi_2) \cdot \langle g, h \rangle \rangle \\
 = & \quad \{ \text{by (2.8)} \} \\
 & \langle i \cdot (\pi_1 \cdot \langle g, h \rangle), j \cdot (\pi_2 \cdot \langle g, h \rangle) \rangle \\
 = & \quad \{ \text{by } \times\text{-cancellation (2.22)} \} \\
 & \langle i \cdot g, j \cdot h \rangle
 \end{aligned}$$

The law we have just derived is known as \times -absorption. (The intuition behind this terminology is that “*split* absorbs \times ”, as a special kind of fusion.) It is a consequence of \times -fusion and \times -cancellation and is depicted as follows:

$$\begin{array}{c}
 A \xleftarrow{\pi_1} A \times B \xrightarrow{\pi_2} B \\
 \uparrow i \quad \uparrow i \times j \quad \uparrow j \\
 D \xleftarrow{\pi_1} D \times E \xrightarrow{\pi_2} E \\
 \swarrow g \quad \searrow h \\
 \quad C
 \end{array}
 \quad (i \times j) \cdot \langle g, h \rangle = \langle i \cdot g, j \cdot h \rangle \quad (2.27)$$

This diagram provides us with two further results about products and projections which can be easily justified:

$$i \cdot \pi_1 = \pi_1 \cdot (i \times j) \quad (2.28)$$

$$j \cdot \pi_2 = \pi_2 \cdot (i \times j) \quad (2.29)$$

Two special properties of $f \times g$ are presented next. The first one expresses a kind of “bi-distribution” of \times with respect to composition:

$$(g \cdot h) \times (i \cdot j) = (g \times i) \cdot (h \times j) \quad (2.30)$$

⁵ Note how this law can be regarded as a pointfree rendering of (2.20).

We will refer to this property as the \times -*functor property*. The other property, which we will refer to as the \times -*functor-id property*, has to do with identity functions:

$$id_A \times id_B = id_{A \times B} \quad (2.31)$$

These two properties will be identified as the *functorial properties* of product. Once again, this choice of terminology will be explained later on.

Let us finally analyse the particular situation in which a *split* is built involving projections π_1 and π_2 only. These exhibit interesting properties, for instance $\langle \pi_1, \pi_2 \rangle = id$. This property is known as \times -*reflexion* and is depicted as follows:⁶

$$\begin{array}{ccccc} A & \xleftarrow{\pi_1} & A \times B & \xrightarrow{\pi_2} & B \\ & \searrow \pi_1 & \uparrow id_{A \times B} & \nearrow \pi_2 & \\ & & A \times B & & \end{array} \quad \langle \pi_1, \pi_2 \rangle = id_{A \times B} \quad (2.32)$$

What about $\langle \pi_2, \pi_1 \rangle$? This corresponds to a diagram

$$\begin{array}{ccccc} B & \xleftarrow{\pi_1} & B \times A & \xrightarrow{\pi_2} & A \\ & \searrow \pi_2 & \uparrow \langle \pi_2, \pi_1 \rangle & \nearrow \pi_1 & \\ & & A \times B & & \end{array}$$

which looks very much the same if submitted to a 180° clockwise rotation (thus A and B swap with each other). This suggests that swap — the name we adopt for $\langle \pi_2, \pi_1 \rangle$ — is its own inverse; this is checked easily as follows:

$$\begin{aligned} & \text{swap} \cdot \text{swap} \\ = & \quad \{ \text{by definition } \text{swap} \stackrel{\text{def}}{=} \langle \pi_2, \pi_1 \rangle \} \\ & \langle \pi_2, \pi_1 \rangle \cdot \text{swap} \\ = & \quad \{ \text{by } \times\text{-fusion (2.26)} \} \\ & \langle \pi_2 \cdot \text{swap}, \pi_1 \cdot \text{swap} \rangle \\ = & \quad \{ \text{definition of swap twice} \} \\ & \langle \pi_2 \cdot \langle \pi_2, \pi_1 \rangle, \pi_1 \cdot \langle \pi_2, \pi_1 \rangle \rangle \\ = & \quad \{ \text{by } \times\text{-cancellation (2.22)} \} \\ & \langle \pi_1, \pi_2 \rangle \\ = & \quad \{ \text{by } \times\text{-reflexion (2.32)} \} \\ & id \end{aligned}$$

Therefore, swap is iso and establishes the following isomorphism

$$A \times B \cong B \times A \quad (2.33)$$

⁶ For an explanation of the word “*reflexion*” in the name chosen for this law (and for others to come) see section 2.13 later on.

which is known as the *commutative property* of product.

The “product datatype” $A \times B$ is essential to information processing and is available in virtually every programming language. In HASKELL one writes (A,B) to denote $A \times B$, for A and B two pre-defined datatypes, `fst` to denote π_1 and `snd` to denote π_2 . In the C programming language this datatype is called the “struct datatype”,

```
struct {
  A first;
  B second;
};
```

while in PASCAL it is called the “record datatype”:

```
record
  first : A;
  second : B
end
```

Isomorphism (2.33) can be re-interpreted in this context as a guarantee that *one does not lose (or gain) anything in swapping fields in record datatypes*. C or PASCAL programmers know also that record-field nesting has the same status, that is to say that, for instance, datatype

```
record
  f : A;
  s : record
    f : B;
    s : C;
  end
end;
```

is abstractly the same as

```
record
  f : record
    f : A;
    s : B
  end;
  s : C;
end;
```

In fact, this is another well-known isomorphism, known as the *associative property* of product:

$$A \times (B \times C) \cong (A \times B) \times C \quad (2.34)$$

This is established by $A \times (B \times C) \xleftarrow{\text{assocr}} (A \times B) \times C$, which is pronounced “associate to the right” and is defined by

$$\text{assocr} \stackrel{\text{def}}{=} \langle \pi_1 \cdot \pi_1, \pi_2 \times id \rangle \quad (2.35)$$

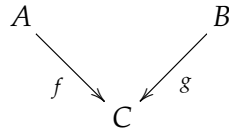
Appendix B lists an extension to the *HASKELL Standard Prelude* that makes isomorphisms such as `swap` and `assoc` available. In this module, the concrete syntax chosen for $\langle f, g \rangle$ is `split f g` and the one chosen for $f \times g$ is `f >< g`.

Exercise 2.4. Rely on (2.24) to prove properties (2.30) and (2.31).

□

2.9 GLUING FUNCTIONS WHICH DO NOT COMPOSE — COPRODUCTS

The *split* functional combinator arose in the previous section as a kind of glue for combining two functions which do not compose but share the same domain. The “dual” situation of two non-composable functions $f : C \leftarrow A$ and $g : C \leftarrow B$ which however share the same codomain is depicted in



It is clear that the kind of glue we need in this case should make it possible to apply f in case we are on the “A-side” or to apply g in case we are on the “B-side” of the diagram. Let us write $[f, g]$ to denote the new kind of combinator. Its codomain will be C . What about its domain?

We need to describe the datatype which is “either an A or a B ”. Since A and B are sets, we may think of $A \cup B$ as such a datatype. This works in case A and B are disjoint sets, but wherever the intersection $A \cap B$ is non-empty it is undecidable whether a value $x \in A \cap B$ is an “A-value” or a “B-value”. In the limit, if $A = B$ then $A \cup B = A = B$, that is to say, we have not invented a new datatype at all. These difficulties can be circumvented by resorting to *disjoint union*,

$$A + B \stackrel{\text{def}}{=} \{i_1 a \mid a \in A\} \cup \{i_2 b \mid b \in B\}$$

assuming the “tagging” functions

$$i_1 a = (t_1, a) \quad , \quad i_2 b = (t_2, b) \tag{2.36}$$

with types⁷ $A \xrightarrow{i_1} A + B \xleftarrow{i_2} B$. Knowing the exact values of tags t_1 and t_2 is not essential to understanding the concept of a disjoint union. It suffices to know that i_1 and i_2 tag differently ($t_1 \neq t_2$) and consistently.

⁷ The tagging functions i_1 and i_2 are usually referred to as the *injections* of the disjoint union.

The values of $A + B$ can be thought of as “copies” of A or B values which are “stamped” with different tags in order to guarantee that values which are simultaneously in A and B do not get mixed up. For instance, the following realizations of $A + B$ in the C programming language,

```
struct {
    int tag; /*1,2*/
    union {
        A ifA;
        B ifB;
    } data;
};
```

or in PASCAL,

```
record
    case tag : integer
    of x =
        1 : (P : A);
        2 : (S : B)
    end;
```

adopt integer tags. In the HASKELL *Standard Prelude*, the $A + B$ datatype is realized by

```
data Either a b = Left a | Right b
```

So, `Left` and `Right` can be thought of as the injections i_1 and i_2 in this realization.

At this level of abstraction, disjoint union $A + B$ is called the *coproduct* of A and B , on top of which we define the new combinator $[f, g]$ (pronounced “either f or g ”) as follows:

$$\begin{aligned}
 [f, g] & : A + B \longrightarrow C \\
 [f, g] x & \stackrel{\text{def}}{=} \begin{cases} x = i_1 a \Rightarrow f a \\ x = i_2 b \Rightarrow g b \end{cases}
 \end{aligned} \tag{2.37}$$

As we did for products, we can express all this in a diagram:

$$\begin{array}{ccccc}
 A & \xrightarrow{i_1} & A + B & \xleftarrow{i_2} & B \\
 & \searrow f & \downarrow [f, g] & \swarrow g & \\
 & & C & &
 \end{array} \tag{2.38}$$

It is interesting to note how similar this diagram is to the one drawn for products — one just has to reverse the arrows, replace projections by injections and the *split* arrow by the *either* one. This expresses the fact that *product* and *coproduct* are *dual* mathematical constructs (compare with *sine* and *cosine* in trigonometry). This duality is of great conceptual economy because everything we can say about product $A \times B$

can be rephrased to coproduct $A + B$. For instance, we may introduce the sum of two functions $f + g$ as the notion dual to product $f \times g$:

$$f + g \stackrel{\text{def}}{=} [i_1 \cdot f, i_2 \cdot g] \quad (2.39)$$

The following list of $+$ -laws provides eloquent evidence of this duality:

$+$ -cancellation :

$$\begin{array}{ccc} A & \xrightarrow{i_1} & A + B \xleftarrow{i_2} B \\ & \searrow g & \downarrow [g,h] \swarrow h \\ & & C \end{array} \quad [g,h] \cdot i_1 = g, [g,h] \cdot i_2 = h \quad (2.40)$$

$+$ -reflexion :

$$\begin{array}{ccc} A & \xrightarrow{i_1} & A + B \xleftarrow{i_2} B \\ & \searrow i_1 & \downarrow id_{A+B} \swarrow i_2 \\ & & A + B \end{array} \quad [i_1, i_2] = id_{A+B} \quad (2.41)$$

$+$ -fusion :

$$\begin{array}{ccc} A & \xrightarrow{i_1} & A + B \xleftarrow{i_2} B \\ & \searrow g & \downarrow [g,h] \swarrow h \\ & f \cdot g & C \swarrow f \cdot h \\ & & D \end{array} \quad f \cdot [g,h] = [f \cdot g, f \cdot h] \quad (2.42)$$

$+$ -absorption :

$$\begin{array}{ccccc} A & \xrightarrow{i_1} & A + B & \xleftarrow{i_2} & B \\ \downarrow i & & \downarrow i+j & & \downarrow j \\ D & \xrightarrow{i_1} & D + E & \xleftarrow{i_2} & E \\ & \searrow g & \downarrow [g,h] \swarrow h & & \\ & & C & & \end{array} \quad [g,h] \cdot (i + j) = [g \cdot i, h \cdot j] \quad (2.43)$$

$+$ -functor :

$$(g \cdot h) + (i \cdot j) = (g + i) \cdot (h + j) \quad (2.44)$$

$+$ -functor-id :

$$id_A + id_B = id_{A+B} \quad (2.45)$$

In summary, the typing-rules of the *either* and *sum* combinators are as follows:

$$\frac{\frac{C \xleftarrow{f} A \quad C \xleftarrow{g} B}{C \xleftarrow{[f,g]} A + B}}{\frac{C \xleftarrow{f} A \quad D \xleftarrow{g} B}{C + D \xleftarrow{f+g} A + B}} \quad (2.46)$$

Exercise 2.5. By analogy (duality) with *swap*, show that $[i_2, i_1]$ is its own inverse and so that fact

$$A + B \cong B + A \quad (2.47)$$

holds.

□

Exercise 2.6. Dualize (2.35), that is, write the iso which witnesses fact

$$A + (B + C) \cong (A + B) + C \quad (2.48)$$

from right to left. Use the *either* syntax available from the HASKELL Standard Prelude to encode this iso in HASKELL.

□

2.10 MIXING PRODUCTS AND COPRODUCTS

Datatype constructions $A \times B$ and $A + B$ have been introduced above as devices required for expressing the codomain of *splits* ($A \times B$) or the domain of *eithers* ($A + B$). Therefore, a function mapping values of a coproduct (say $A + B$) to values of a product (say $A' \times B'$) can be expressed alternatively as an *either* or as a *split*. In the first case, both components of the *either* combinator are *splits*. In the latter, both components of the *split* combinator are *eithers*.

This exchange of format in defining such functions is known as the *exchange law*. It states the functional equality which follows:

$$[\langle f, g \rangle, \langle h, k \rangle] = \langle [f, h], [g, k] \rangle \quad (2.49)$$

It can be checked by type-inference that both the left-hand side and the right-hand side expressions of this equality have type $B \times D \leftarrow A + C$, for $B \xleftarrow{f} A$, $D \xleftarrow{g} A$, $B \xleftarrow{h} C$ and $D \xleftarrow{k} C$.

An example of a function which is in the exchange-law format is isomorphism

$$A \times (B + C) \xleftarrow{\text{undistr}} (A \times B) + (A \times C) \quad (2.50)$$

(pronounce undistr as “un-distribute-right”) which is defined by

$$\text{undistr} \stackrel{\text{def}}{=} [id \times i_1, id \times i_2] \quad (2.51)$$

and witnesses the fact that product distributes through coproduct:

$$A \times (B + C) \cong (A \times B) + (A \times C) \quad (2.52)$$

In this context, suppose that we know of three functions $D \xleftarrow{f} A$, $E \xleftarrow{g} B$ and $F \xleftarrow{h} C$. By (2.46) we infer $E + F \xleftarrow{g+h} B + C$. Then, by (2.25) we infer

$$D \times (E + F) \xleftarrow{f \times (g+h)} A \times (B + C) \quad (2.53)$$

So, it makes sense to combine products and sums of functions and the expressions which denote such combinations have the same “shape” (or symbolic pattern) as the expressions which denote their domain and range — the $\dots \times (\dots + \dots)$ “shape” in this example. In fact, if we *abstract* such a pattern via some symbol, say F — that is, if we define

$$F(\alpha, \beta, \gamma) \stackrel{\text{def}}{=} \alpha \times (\beta + \gamma)$$

— then we can write $F(D, E, F) \xleftarrow{F(f,g,h)} F(A, B, C)$ for (2.53).

This kind of abstraction works for every combination of products and coproducts. For instance, if we now abstract the right-hand side of (2.50) via pattern

$$G(\alpha, \beta, \gamma) \stackrel{\text{def}}{=} (\alpha \times \beta) + (\alpha \times \gamma)$$

we have $G(f, g, h) = (f \times g) + (f \times h)$, a function which maps $G(A, B, C) = (A \times B) + (A \times C)$ onto $G(D, E, F) = (D \times E) + (D \times F)$. All this can be put in a diagram

$$\begin{array}{ccc} F(A, B, C) & \xleftarrow{\text{undistr}} & G(A, B, C) \\ \downarrow F(f,g,h) & & \downarrow G(f,g,h) \\ F(D, E, F) & & G(D, E, F) \end{array}$$

which unfolds to

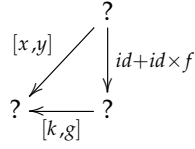
$$\begin{array}{ccc} A \times (B + C) & \xleftarrow{\text{undistr}} & (A \times B) + (A \times C) \\ \downarrow f \times (g+h) & & \downarrow (f \times g) + (f \times h) \\ D \times (E + F) & & (D \times E) + (D \times F) \end{array} \quad (2.54)$$

once the F and G patterns are instantiated. An interesting topic which stems from (completing) this diagram will be discussed in the next section.

Exercise 2.7. Apply the exchange law to undistr.

□

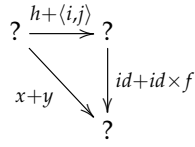
Exercise 2.8. Complete the “?”s in diagram



and then solve the implicit equation for x and y .

□

Exercise 2.9. Repeat exercise 2.8 with respect to diagram



□

Exercise 2.10. Show that $\langle [f, h] \cdot (\pi_1 + \pi_1), [g, k] \cdot (\pi_2 + \pi_2) \rangle$ reduces to $[f \times g, h \times k]$.

□

2.11 ELEMENTARY DATATYPES

So far we have talked mostly about arbitrary datatypes represented by capital letters A , B , etc. (lowercase a , b , etc. in the HASKELL illustrations). We also mentioned \mathbb{R} , Bool and \mathbb{N} and, in particular, the fact that we can associate to each natural number n its *initial segment* $n = \{1, 2, \dots, n\}$. We extend this to \mathbb{N}_0 by stating $0 = \{\}$ and, for $n > 0$, $n + 1 = \{n + 1\} \cup n$.

Initial segments can be identified with enumerated types and are regarded as primitive datatypes in our notation. We adopt the convention that primitive datatypes are written in the *sans serif* font and so, strictly speaking, n is distinct from n : the latter denotes a natural number while the former denotes a datatype.

Datatype 0

Among such enumerated types, 0 is the smallest because it is empty. This is the `Void` datatype in HASKELL, which has no constructor at all.

Datatype 0 (which we tend to write simply as 0) may not seem very “useful” in practice but it is of theoretical interest. For instance, it is easy to check that the following “obvious” properties hold,

$$A + 0 \cong A \quad (2.55)$$

$$A \times 0 \cong 0 \quad (2.56)$$

where the second is actually an equality: $A \times 0 = 0$.

Datatype 1

Next in the sequence of initial segments we find 1, which is singleton set $\{1\}$. How useful is this datatype? Note that every datatype A containing exactly one element is isomorphic to $\{1\}$, e.g. $A = \{\text{NIL}\}$, $A = \{0\}$, $A = \{1\}$, $A = \{\text{FALSE}\}$, etc.. We represent this class of singleton types by 1.

Recall that isomorphic datatypes have the same expressive power and so are “abstractly identical”. So, the actual choice of inhabitant for datatype 1 is irrelevant, and we can replace any particular singleton set by another without losing information. This is evident from the following, observing isomorphism,

$$A \times 1 \cong A \quad (2.57)$$

which can be read informally as follows: if the second component of a record (“struct”) cannot change, then it is useless and can be ignored. Selector π_1 is, in this context, an iso mapping the left-hand side of (2.57) to its right-hand side. Its inverse is $\langle id, \underline{c} \rangle$ where c is a particular choice of inhabitant for datatype 1.

In summary, when referring to datatype 1 we will mean an arbitrary singleton type, and there is a unique iso (and its inverse) between two such singleton types. The HASKELL representative of 1 is datatype $()$, called the *unit type*, which contains exactly constructor $()$. It may seem confusing to denote the type and its unique inhabitant by the same symbol but it is not, since HASKELL keeps track of types and constructors in separate symbol sets.

Note that any function of type $A \rightarrow 1$ is bound to be a constant function. This function, usually called the “bang”, or “sink” function, is denoted by an exclamation mark:

$$! : A \rightarrow 1 \quad (2.58)$$

Clearly, it is *the unique* function of its type. (Can you write a different one, of the same type?)

Finally, what can we say about $1 + A$? Every function $B \xleftarrow{f} 1 + A$ observing this type is bound to be an *either* $[b_0, g]$ for $b_0 \in B$ and $B \xleftarrow{g} A$. This is very similar to the handling of a pointer in C or PASCAL: we “pull a rope” and either we get nothing (1) or we get

something useful of type B . In such a programming context “nothing” above means a predefined value `NIL`. This analogy supports our preference in the sequel for `NIL` as canonical inhabitant of datatype 1. In fact, we will refer to $1 + A$ (or $A + 1$) as the “pointer to A ” datatype. This corresponds to the `Maybe` type constructor of the HASKELL *Standard Prelude*.

Datatype 2

Let us inspect the $1 + 1$ instance of the “pointer” construction just mentioned above. Any observation $B \xleftarrow{f} 1 + 1$ can be decomposed in two constant functions: $f = [b_1, b_2]$. Now suppose that $B = \{b_1, b_2\}$ (for $b_1 \neq b_2$). Then $1 + 1 \cong B$ will hold, for whatever choice of inhabitants b_1 and b_2 . So we are in a situation similar to 1: we will use symbol 2 to represent the abstract class of all such B s containing exactly two elements. Therefore, we can write:

$$1 + 1 \cong 2$$

Of course, $\text{Bool} = \{\text{TRUE}, \text{FALSE}\}$ and initial segment $2 = \{1, 2\}$ are in this abstract class. In the sequel we will show some preference for the particular choice of inhabitants $b_1 = \text{TRUE}$ and $b_2 = \text{FALSE}$, which enables us to use symbol 2 in places where `Bool` is expected. Clearly,

$$2 \times A \cong A + A \quad (2.59)$$

Exercise 2.11. Derive the isomorphism

$$(B + C) \times A \xleftarrow{\text{undistr}} (B \times A) + (C \times A) \quad (2.60)$$

from *undistr* (2.50) and other isomorphisms studied thus far.

□

Exercise 2.12. Furthermore, show that (2.59) follows from (2.60) and, on the practical side, relate HASKELL expression

`either (split (const True) id) (split (const False) id)`

to the same isomorphism (2.59).

□

2.12 NATURAL PROPERTIES

Let us resume discussion about *undistr* and the two other functions in diagram (2.54). What about using *undistr* itself to close this diagram,

at the bottom? Note that definition (2.51) works for D, E and F in the same way it does for A, B and C . (Indeed, the particular choice of symbols A, B and C in (2.50) was rather arbitrary.) Therefore, we get:

$$\begin{array}{ccc} A \times (B + C) & \xleftarrow{\text{undistr}} & (A \times B) + (A \times C) \\ \downarrow f \times (g+h) & & \downarrow (f \times g) + (f \times h) \\ D \times (E + F) & \xleftarrow{\text{undistr}} & (D \times E) + (D \times F) \end{array}$$

which expresses a very important property of undistr :

$$(f \times (g + h)) \cdot \text{undistr} = \text{undistr} \cdot ((f \times g) + (f \times h)) \quad (2.61)$$

This is called the *natural* property of undistr . This kind of property (often called “*free*” instead of “*natural*”) is not a privilege of undistr . As a matter of fact, every function interfacing patterns such as F or G above will exhibit its own *natural* property. Furthermore, we have already quoted *natural* properties without mentioning it. Recall (2.10), for instance. This property (establishing id as the *unit* of composition) is, after all, the *natural* property of id . In this case we have $F \alpha = G \alpha = \alpha$, as can be easily observed in diagram (2.11).

In general, *natural* properties are described by diagrams in which two “copies” of the operator of interest are drawn as horizontal arrows:

$$\begin{array}{ccc} A & F A \xleftarrow{\phi} G A & (F f) \cdot \phi = \phi \cdot (G f) \\ f \downarrow & \downarrow F f \quad \downarrow G f & \\ B & F B \xleftarrow{\phi} G B & \end{array} \quad (2.62)$$

Note that f is universally quantified, that is to say, the *natural* property holds for every $f : B \leftarrow A$.

Diagram (2.62) corresponds to unary patterns F and G . As we have seen with undistr , other functions (g, h etc.) come into play for multiary patterns. A very important rôle will be assigned throughout this book to these F, G , etc. “shapes” or patterns which are shared by pointfree functional expressions and by their domain and codomain expressions. From chapter 3 onwards we will refer to them by their proper name — “functor” — which is standard in mathematics and computer science. Then we will also explain the names assigned to properties such as, for instance, (2.30) or (2.44).

Exercise 2.13. Show that (2.28) and (2.29) are natural properties. Dualize these properties. **Hint:** recall diagram (2.43).

□

Exercise 2.14. Establish the natural properties of the swap (2.33) and assocr (2.35) isomorphisms.

□

Exercise 2.15. Draw the natural property of function $\phi = \text{swap} \cdot (\text{id} \times \text{swap})$ as a diagram, that is, identify F and G in (2.62) for this case.

Do the same for $\phi = \text{coswap} \cdot (\text{swap} + \text{swap})$ where $\text{coswap} = [i_2, i_1]$.

□

2.13 UNIVERSAL PROPERTIES

Functional constructs $\langle f, g \rangle$ and $[f, g]$ — and their derivatives $f \times g$ and $f + g$ — provide good illustration about what is meant by a *program combinator* in a compositional approach to programming: the combinator is put forward equipped with a concise *set of properties* which enable programmers to transform programs, reason about them and perform useful calculations. This raises a *programming methodology* which is scientific and stable.

Such properties bear standard names such as *cancellation*, *reflexion*, *fusion*, *absorption* etc.. Where do these names come from? As a rule, for each combinator to be defined one has to define suitable constructions at “interface”-level⁸, e.g. $A \times B$ and $A + B$. These are not chosen or invented at random: each is defined in a way such that the associated combinator is uniquely defined. This is assured by a so-called *universal property* from which the others can be derived.

Take product $A \times B$, for instance. Its universal property states that, for each pair of arrows $A \xleftarrow{f} C$ and $B \xleftarrow{g} C$, there exists an arrow $A \times B \xleftarrow{\langle f, g \rangle} C$ such that

$$k = \langle f, g \rangle \Leftrightarrow \begin{cases} \pi_1 \cdot k = f \\ \pi_2 \cdot k = g \end{cases} \quad (2.63)$$

holds — recall diagram (2.23) — for all $A \times B \xleftarrow{k} C$.

Note that (2.63) is an *equivalence*, implicitly stating that $\langle f, g \rangle$ is the *unique* arrow satisfying the property on the right. In fact, read (2.63) in the \Rightarrow direction and let k be $\langle f, g \rangle$. Then $\pi_1 \cdot \langle f, g \rangle = f$ and $\pi_2 \cdot \langle f, g \rangle = g$ will hold, meaning that $\langle f, g \rangle$ effectively obeys the property on the right. In other words, we have derived \times -cancellation (2.22). Reading (2.63) in the \Leftarrow direction we understand that, if some k satisfies such properties, then it “has to be” the same arrow as $\langle f, g \rangle$.

The relevance of universal property (2.63) is that it offers a way of *solving equations* of the form $k = \langle f, g \rangle$. Take for instance the follow-

⁸ In the current context, *programs* “are” functions and *program-interfaces* “are” the datatypes involved in functional signatures.

ing exercise: can the identity be expressed, or “reflected”, using this combinator? We just solve the equation $id = \langle f, g \rangle$ for f and g :

$$\begin{aligned}
 id &= \langle f, g \rangle \\
 &\equiv \{ \text{universal property (2.63)} \} \\
 &\quad \left\{ \begin{array}{l} \pi_1 \cdot id = f \\ \pi_2 \cdot id = g \end{array} \right. \\
 &\equiv \{ \text{by (2.10)} \} \\
 &\quad \left\{ \begin{array}{l} \pi_1 = f \\ \pi_2 = g \end{array} \right.
 \end{aligned}$$

The equation has the unique solutions $f = \pi_1$ and $g = \pi_2$ which, once substituted in the equation itself, yield

$$id = \langle \pi_1, \pi_2 \rangle$$

i.e., nothing but the \times -reflexion law (2.32).

All other laws can be calculated from the universal property in a similar way. For instance, the \times -fusion (2.26) law is obtained by solving the equation $k = \langle i, j \rangle$ again for f and g , but this time fixing $k = \langle i, j \rangle \cdot h$, assuming i, j and h given:⁹

$$\begin{aligned}
 \langle i, j \rangle \cdot h &= \langle f, g \rangle \\
 &\equiv \{ \text{universal property (2.63)} \} \\
 &\quad \left\{ \begin{array}{l} \pi_1 \cdot (\langle i, j \rangle \cdot h) = f \\ \pi_2 \cdot (\langle i, j \rangle \cdot h) = g \end{array} \right. \\
 &\equiv \{ \text{composition is associative (2.8)} \} \\
 &\quad \left\{ \begin{array}{l} (\pi_1 \cdot \langle i, j \rangle) \cdot h = f \\ (\pi_2 \cdot \langle i, j \rangle) \cdot h = g \end{array} \right. \\
 &\equiv \{ \text{by } \times\text{-cancellation (derived above)} \} \\
 &\quad \left\{ \begin{array}{l} i \cdot h = f \\ j \cdot h = g \end{array} \right.
 \end{aligned}$$

Substituting the solutions $f = i \cdot h$ and $g = j \cdot h$ in the equation, we get the \times -fusion law: $\langle i, j \rangle \cdot h = \langle i \cdot h, j \cdot h \rangle$.

It will take about the same effort to derive *split* structural equality

$$\langle i, j \rangle = \langle f, g \rangle \Leftrightarrow \left\{ \begin{array}{l} i = f \\ j = g \end{array} \right. \quad (2.64)$$

from universal property (2.63) — just let $k = \langle i, j \rangle$ and solve.

⁹ Solving equations of this kind is reminiscent of many similar calculations carried out in school maths and physics courses. The spirit is the same. The difference is that this time one is not calculating water pump debits, accelerations, velocities, or other physical entities: the solutions of our equations are (just) functional *programs*.

Similar arguments can be built around coproduct's universal property,

$$k = [f, g] \Leftrightarrow \begin{cases} k \cdot i_1 = f \\ k \cdot i_2 = g \end{cases} \quad (2.65)$$

from which structural equality of *either*s can be inferred,

$$[i, j] = [f, g] \Leftrightarrow \begin{cases} i = f \\ j = g \end{cases} \quad (2.66)$$

as well as the other properties we know about this combinator.

Exercise 2.16. Show that *assocr* (2.35) is iso by solving the equation $\text{assocr} \cdot \text{assocl} = \text{id}$ for *assocl*. **Hint:** don't ignore the role of universal property (2.63) in the calculation.

□

Exercise 2.17. Prove the equality: $[\langle f, \underline{k} \rangle, \langle g, \underline{k} \rangle] = \langle [f, g], \underline{k} \rangle$

□

Exercise 2.18. Derive $+$ -cancellation (2.40), $+$ -reflexion (2.41) and $+$ -fusion (2.42) from universal property (2.65). Then derive the exchange law (2.49) from the universal property of product (2.63) or coproduct (2.65).

□

Exercise 2.19. Function $\text{coassocr} = [id + i_1, i_2 \cdot i_2]$ is a witness of isomorphism $(A + B) + C \cong A + (B + C)$, from left to right. Calculate its converse *coassocl* by solving the equation

$$\underbrace{[x, [y, z]]}_{\text{coassocl}} \cdot \text{coassocr} = \text{id} \quad (2.67)$$

for x, y and z .

□

Exercise 2.20. Let δ be a function of which you know that $\pi_1 \cdot \delta = \text{id}$ e $\pi_2 \cdot \delta = \text{id}$ hold. Show that necessarily δ satisfies the natural property $(f \times f) \cdot \delta = \delta \cdot f$.

□

2.14 GUARDS AND MCCARTHY'S CONDITIONAL

Most functional programming languages and notations cater for pointwise conditional expressions of the form

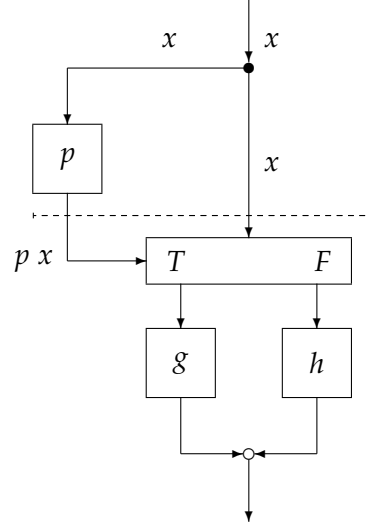
$$\text{if } p \ x \text{ then } g \ x \text{ else } h \ x \quad (2.68)$$

which evaluates to $g \ x$ in case $p \ x$ holds and to $h \ x$ otherwise, that is

$$\begin{cases} p \ x & \Rightarrow g \ x \\ \neg(p \ x) & \Rightarrow h \ x \end{cases}$$

given some predicate $\text{Bool} \xleftarrow{p} A$, some “then”-function $B \xleftarrow{g} A$ and some “else”-function $B \xleftarrow{h} A$.

Can (2.68) be written in the pointfree style?



The drawing above is an attempt to express such a conditional expression as a “block”-diagram. Firstly, the input x is copied, the left copy being passed to predicate p yielding the Boolean $p \ x$. One can easily define this part using $\text{copy} = \langle \text{id}, \text{id} \rangle$.

The informal part of the diagram is the T - F “switch”: it should channel x to g in case $p \ x$ switches the T -output, or channel x to h otherwise.

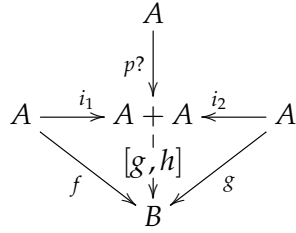
At first sight, this T - F gate should be of type $\mathbb{B} \times A \rightarrow A \times A$. But the output cannot be $A \times A$, as f or g act in *alternation*, not in *parallel* — it should rather be $A + A$, in which case the last step is achieved just by running $[g, h]$. How does one switch from our starting product-based model of conditionals to a coproduct-based one?

The key observation is that the type $\mathbb{B} \times A$ market by the dashed line in the block-diagram is isomorphic to $A + A$, recall (2.59). That is, the information captured by the pair $(p \ x, x) \in \mathbb{B} \times A$ can be converted into a unique $y \in A + A$ with no loss of information. Let us define a new combinator for this, denoted $p?$:

$$(p?)a = \begin{cases} p \ a & \Rightarrow i_1 \ a \\ \neg(p \ a) & \Rightarrow i_2 \ a \end{cases} \quad (2.69)$$

We call $A + A \xleftarrow{p?} A$ a *guard*, or better, the guard associated to a given predicate $\text{Bool} \xleftarrow{p} A$. In a sense, guard $p?$ is more “informative” than p alone: it provides information about the outcome of testing p on some input a , encoded in terms of the coproduct injections (i_1 for a *true* outcome and i_2 for a *false* outcome, respectively) without losing the input a itself.

Finally, the composition $[g, h] \cdot p?$, depicted in the following diagram



has (2.68) as pointwise meaning. It is a well-known functional combinator termed “McCarthy conditional”¹⁰ and usually denoted by the expression $p \rightarrow g, h$. Altogether, we have the definition

$$p \rightarrow g, h \stackrel{\text{def}}{=} [g, h] \cdot p? \quad (2.70)$$

which suggests that, to reason about conditionals, one may seek help in the algebra of coproducts. Indeed, the following fact,

$$f \cdot (p \rightarrow g, h) = p \rightarrow f \cdot g, f \cdot h \quad (2.71)$$

which we shall refer to as the *first McCarthy’s conditional fusion law*¹¹, is nothing but an immediate consequence of $+$ -fusion (2.42).

We shall introduce and define instances of predicate p as long as they are needed. A particularly important assumption of our notation should, however, be mentioned at this point: we assume that, for every datatype A , the equality predicate $\text{Bool} \xleftarrow{=}_A A \times A$ is defined in a way which guarantees three basic properties: reflexivity ($a =_A a$ for every a), transitivity ($a =_A b$ and $b =_A c$ implies $a =_A c$) and symmetry ($a =_A b$ iff $b =_A a$). Subscript A in $=_A$ will be dropped wherever implicit in the context.

In HASKELL programming, the equality predicate for a type becomes available by declaring the type as an instance of class `Eq`, which exports equality predicate `(==)`. This does not, however, guarantee the reflexive, transitive and symmetry properties, which need to be proved by dedicated mathematical arguments.

We close this section with an illustration of how *smart* pointfree algebra can be in reasoning about functions that *one does not actually define explicitly*. It also shows how relevant the *natural properties* studied in section 2.12 are. The issue is that our definition of a guard (2.69) is pointwise and most likely unsuitable to prove facts such as, for instance,

$$p? \cdot f = (f + f) \cdot (p \cdot f)? \quad (2.72)$$

Thinking better, instead of “inventing” (2.69), we might (and perhaps should!) have defined

$$A \xrightarrow{\langle p, id \rangle} 2 \times A \xrightarrow{\alpha} A + A \quad (2.73)$$

$\searrow p?$

¹⁰ After John McCarthy, the computer scientist who first defined it.

¹¹ For the second one go to exercise 2.22.

which actually expresses rather closely our strategy of switching from products to coproducts in the definition of $(p?)$. Isomorphism α (2.59) is the subject of exercise 2.12. Do we need to define it explicitly? Perhaps not: from its type, $2 \times A \rightarrow A + A$, we immediately infer its natural (or “free”) property:

$$\alpha \cdot (id \times f) = (f + f) \cdot \alpha \quad (2.74)$$

It turns out that this is the *knowledge* we need about α in order to prove (2.72), as the following calculation shows:

$$\begin{aligned} & p? \cdot f \\ = & \{ (2.73); \langle p, id \rangle \cdot f = \langle p \cdot f, f \rangle \} \\ & \alpha \cdot \langle p \cdot f, f \rangle \\ = & \{ \times\text{-absorption (2.27)} \} \\ & \alpha \cdot (id \times f) \cdot \langle p \cdot f, id \rangle \\ = & \{ \text{free theorem of } \alpha \text{ (2.74)} \} \\ & (f + f) \cdot \alpha \cdot \langle p \cdot f, id \rangle \\ = & \{ \text{again (2.73); } \langle p, id \rangle \cdot f = \langle p \cdot f, f \rangle \} \\ & (f + f) \cdot (p \cdot f)? \\ & \square \end{aligned}$$

Other examples of this kind of reasoning, based on natural (free) properties of isomorphisms — and often on “shunting” them around using laws (2.18, 2.19) — will be given later in this book.

The less one has to write to solve a problem, the better. One saves time and one’s brain, adding to productivity. This is often called *elegance* when applying a scientific method. (Unfortunately, be prepared for much lack of it in the software engineering field!)

Exercise 2.21. Prove that the following equality between two conditional expressions

$$\begin{aligned} & k \text{ (if } p \text{ x then } f \text{ x else } h \text{ x, if } p \text{ x then } g \text{ x else } i \text{ x)} \\ = & \text{if } p \text{ x then } k \text{ (}\lambda ap \text{ } f \text{ x, } \lambda ap \text{ } g \text{ x)} \text{ else } k \text{ (} h \text{ x, } i \text{ x)} \end{aligned}$$

holds by rewriting it in the pointfree style (using the McCarthy’s conditional combinator) and applying the exchange law (2.49), among others.

□

Exercise 2.22. Prove the first McCarthy’s conditional fusion law (2.71). Then, from (2.70) and property (2.72), infer the second such law:

$$(p \rightarrow f, g) \cdot h = (p \cdot h) \rightarrow (f \cdot h), (g \cdot h) \quad (2.75)$$

□

Exercise 2.23. Prove that property

$$\langle f, \langle p \rightarrow q, h \rangle \rangle = p \rightarrow \langle f, q \rangle, \langle f, h \rangle \quad (2.76)$$

and its corollary

$$(p \rightarrow g, h) \times f = p \cdot \pi_1 \rightarrow g \times f, h \times f \quad (2.77)$$

hold, assuming the basic fact:

$$p \rightarrow f, f = f \quad (2.78)$$

□

Exercise 2.24. Define $p(x, y) = x > y$ and the maximum of two integers, $m(x, y)$, by:

$$m = p \rightarrow \pi_1, \pi_2$$

Then show that

$$\text{succ} \cdot m = m \cdot (\text{succ} \times \text{succ})$$

holds, by using the McCarthy conditional fusion-laws and basic arithmetics.

□

2.15 GLUING FUNCTIONS WHICH DO NOT COMPOSE — EXPONENTIALS

Now that we have made the distinction between the pointfree and pointwise functional notations reasonably clear, it is instructive to revisit section 2.2 and identify *functional application* as the “bridge” between the pointfree and pointwise worlds. However, we should say “a bridge” rather than “the bridge”, for in this section we enrich such an interface with another “bridge” which is very relevant to programming.

Suppose we are given the task to combine two functions, one binary $B \xleftarrow{f} C \times A$ and the other unary: $D \xleftarrow{g} A$. It is clear that none of the combinations $f \cdot g$, $\langle f, g \rangle$ or $[f, g]$ is well-typed. So, f and g cannot be put together directly — they require some extra interfacing.

Note that $\langle f, g \rangle$ would be well-defined in case the C component of f 's domain could be somehow “ignored”. Suppose, in fact, that in

some particular context the first argument of f happens to be “irrelevant”, or to be frozen to some $c \in C$. It is easy to derive a new function

$$\begin{aligned} f_c & : A \rightarrow B \\ f_c a & \stackrel{\text{def}}{=} f(c, a) \end{aligned}$$

from f which combines nicely with g via the *split* combinator: $\langle f_c, g \rangle$ is well-defined and bears type $B \times D \leftarrow A$. For instance, suppose that $C = A$ and f is the equality predicate $=$ on A . Then $\text{Bool} \xleftarrow{=} A$ is the “equal to c ” predicate on A values:

$$=_c a \stackrel{\text{def}}{=} a = c \quad (2.79)$$

As another example, recall function *twice* (2.3) which could be defined as \times_2 using the new notation.

However, we need to be more careful about what is meant by f_c . Such as functional application, expression f_c interfaces the pointfree and the pointwise levels — it involves a function (f) and a value (c).

But, for $B \xleftarrow{f} C \times A$, there is a major distinction between $f c$ and f_c — while the former denotes a value of type B , i.e. $f c \in B$, f_c denotes a function of type $B \leftarrow A$. We will say that $f_c \in B^A$ by introducing a new datatype construct which we will refer to as the *exponential*:

$$B^A \stackrel{\text{def}}{=} \{g \mid g : B \leftarrow A\} \quad (2.80)$$

There are strong reasons to adopt the B^A notation to the detriment of the more obvious $B \leftarrow A$ or $A \rightarrow B$ alternatives, as we shall see shortly.

The B^A exponential datatype is therefore inhabited by functions from A to B , that is to say, functional declaration $g : B \leftarrow A$ means the same as $g \in B^A$. And what do we want functions for? We want to apply them. So it is natural to introduce the *apply* operator

$$\begin{aligned} ap & : B \xleftarrow{ap} B^A \times A \\ ap(f, a) & \stackrel{\text{def}}{=} f a \end{aligned} \quad (2.81)$$

which applies a function f to an argument a .

Back to generic binary function $B \xleftarrow{f} C \times A$, let us now think of the operation which, for every $c \in C$, produces $f_c \in B^A$. This can be regarded as a function of signature $B^A \leftarrow C$ which expresses f as a kind of C -indexed family of functions of signature $B \leftarrow A$. We will denote such a function by \bar{f} (read \bar{f} as “ f transposed”). Intuitively, we want f and \bar{f} to be related to each other by the following property:

$$f(c, a) = (\bar{f} c) a \quad (2.82)$$

Given c and a , both expressions denote the same value. But, in a sense, \bar{f} is more tolerant than f : while the latter is binary and requires *both*

arguments (c, a) to become available before application, the former is happy to be provided with c first and with a later on, if actually required by the evaluation process.

Similarly to $A \times B$ and $A + B$, exponential B^A involves a universal property,

$$k = \bar{f} \Leftrightarrow f = ap \cdot (k \times id) \quad (2.83)$$

from which laws for cancellation, reflexion and fusion can be derived:

Exponentials cancellation :

$$\begin{array}{ccc} B^A & B^A \times A \xrightarrow{ap} B & f = ap \cdot (\bar{f} \times id) \\ \bar{f} \uparrow & \bar{f} \times id \uparrow \nearrow f & \\ C & C \times A & \end{array} \quad (2.84)$$

Exponentials reflexion :

$$\begin{array}{ccc} B^A & B^A \times A \xrightarrow{ap} B & \overline{ap} = id_{B^A} \\ id_{B^A} \uparrow & id_{B^A} \times id_A \uparrow \nearrow ap & \\ B^A & B^A \times A & \end{array} \quad (2.85)$$

Exponentials fusion :

$$\begin{array}{ccc} B^A & B^A \times A \xrightarrow{ap} B & \overline{g \cdot (f \times id)} = \bar{g} \cdot f \\ \bar{g} \uparrow & \bar{g} \times id \uparrow \nearrow g & \\ C & C \times A & \\ f \uparrow & f \times id \uparrow \nearrow g \cdot (f \times id) & \\ D & D \times A & \end{array} \quad (2.86)$$

Note that the cancellation law is nothing but fact (2.82) written in the pointfree style.

Is there an absorption law for exponentials? The answer is affirmative but first we need to introduce a new functional combinator which arises as the transpose of $f \cdot ap$ in the following diagram:

$$\begin{array}{ccc} D^A \times A & \xrightarrow{ap} & D \\ \bar{f \cdot ap} \times id \uparrow & & \uparrow f \\ B^A \times A & \xrightarrow{ap} & B \end{array}$$

We shall denote this by f^A and its type-rule is as follows:

$$\frac{C \xleftarrow{f} B}{C^A \xleftarrow{f^A} B^A}$$

It can be shown that, once A and $C \xleftarrow{f} B$ are fixed, f^A is the function which accepts some input function $B \xleftarrow{g} A$ as argument and produces function $f \cdot g$ as result (see exercise 2.40). So f^A is the “compose with f ” functional combinator:

$$(f^A)g \stackrel{\text{def}}{=} f \cdot g \quad (2.87)$$

Now we are ready to understand the laws which follow:

Exponentials absorption :

$$\begin{array}{ccc} D^A & D^A \times A \xrightarrow{ap} D & \overline{f \cdot g} = f^A \cdot \bar{g} \\ \uparrow f^A & \uparrow f^A \times id & \uparrow f \\ B^A & B^A \times A \xrightarrow{ap} B & \\ \uparrow \bar{g} & \uparrow \bar{g} \times id & \nearrow g \\ C & C \times A & \end{array} \quad (2.88)$$

(Note how, from this, we also get $f^A = \overline{f \cdot ap}$.)

Exponentials-functor :

$$(g \cdot h)^A = g^A \cdot h^A \quad (2.89)$$

Exponentials-functor-id :

$$id^A = id \quad (2.90)$$

WHY THE EXPONENTIAL NOTATION. To conclude this section we need to explain why we have adopted the apparently esoteric B^A notation for the “function from A to B ” data type. This is the opportunity to relate what we have seen above with two (higher order) functions which are very familiar to functional programmers. In the HASKELL Prelude they are defined thus:

```
curry :: ((a,b) -> c) -> (a -> b -> c)
curry f a b = f (a,b)

uncurry :: (a -> b -> c) -> (a,b) -> c
uncurry f (a,b) = f a b
```

In our notation for types, `curry` maps functions in function space $C^{A \times B}$ to functions in $(C^B)^A$; and `uncurry` maps functions from the latter function space to the former.

Let us calculate the meaning of `curry` by removing variables from its definition:

$$\underbrace{(\text{curry } f \ a)}_{\bar{f}} b = f \ (a, b)$$

$$\begin{aligned}
&\equiv \{ \text{introduce } g \} \\
&\quad g \, b = f(a, b) \\
&\equiv \{ \text{since } g \, b = ap(g, b) \text{ (2.81)} \} \\
&\quad ap(g, b) = f(a, b) \\
&\equiv \{ g = \bar{f} \, a ; \text{natural-id} \} \\
&\quad ap(\bar{f} \, a, id \, b) = f(a, b) \\
&\equiv \{ \text{product of functions: } (f \times g)(x, y) = (f \, x, g \, y) \} \\
&\quad ap((\bar{f} \times id)(a, b)) = f(a, b) \\
&\equiv \{ \text{composition} \} \\
&\quad (ap \cdot (\bar{f} \times id))(a, b) = f(a, b) \\
&\equiv \{ \text{extensionality (2.5), i.e. removing points } a \text{ and } b \} \\
&\quad ap \cdot (\bar{f} \times id) = f
\end{aligned}$$

From the above we infer that the definition of `curry` is a re-statement of the cancellation law (2.84). That is,

$$\text{curry } f \stackrel{\text{def}}{=} \bar{f} \tag{2.91}$$

and `curry` is transposition in HASKELL-speak.¹²

Next we do the same for the definition of `uncurry` :

$$\begin{aligned}
&\underbrace{\text{uncurry } f}_{k} (a, b) = f \, a \, b \\
&\equiv \{ \text{introduce } k ; \text{lefthand side as calculated above} \} \\
&\quad k (a, b) = (ap \cdot (f \times id)) (a, b) \\
&\equiv \{ \text{extensionality (2.5)} \} \\
&\quad k = ap \cdot (f \times id) \\
&\equiv \{ \text{universal property (2.83)} \} \\
&\quad f = \bar{k} \\
&\equiv \{ \text{expand } k \} \\
&\quad f = \overline{\text{uncurry } f}
\end{aligned}$$

We conclude that `uncurry` is the inverse of transposition, that is, of `curry`. We shall use the abbreviation \hat{f} for `uncurry` f , whereby the above equality is written $f = \hat{\hat{f}}$. It can also be checked that $f = \hat{\hat{f}}$ also holds, instantiating k above by \hat{f} :

$$\begin{aligned}
&\hat{\hat{f}} = ap \cdot (\bar{f} \times id) \\
&\equiv \{ \text{cancellation (2.84)} \}
\end{aligned}$$

¹² This terminology widely adopted in other functional languages.

$$\widehat{\widehat{f}} = f$$

□

So $\text{uncurry} \text{ — i.e. } (\widehat{_}) \text{ — and } \text{curry} \text{ — i.e. } (\overline{_}) \text{ — are inverses of each other,}$

$$g = \overline{f} \Leftrightarrow \hat{g} = f \quad (2.92)$$

leading to isomorphism

$$A \rightarrow C^B \cong A \times B \rightarrow C$$

which can also be written as

$$(C^B)^A \begin{array}{c} \xrightarrow{\text{uncurry}} \\ \cong \\ \xleftarrow{\text{curry}} \end{array} C^{A \times B} \quad (2.93)$$

decorated with the corresponding witnesses.¹³

Isomorphism (2.93) is at the core of the theory and practice of functional programming. It clearly resembles a well known equality concerning numeric exponentials, $b^{c \times a} = (b^a)^c$. Moreover, other known facts about numeric exponentials, e.g. $a^{b+c} = a^b \times a^c$ or $(b \times c)^a = b^a \times c^a$ also find their counterpart in functional exponentials. The counterpart of the former,

$$A^{B+C} \cong A^B \times A^C \quad (2.94)$$

arises from the uniqueness of the *either* combination: every pair of functions $(f, g) \in A^B \times A^C$ leads to a unique function $[f, g] \in A^{B+C}$ and vice-versa, every function in A^{B+C} is the *either* of some function in A^B and of another in A^C .

The function exponentials counterpart of the second fact about numeric exponentials above is

$$(B \times C)^A \cong B^A \times C^A \quad (2.95)$$

This can be justified by a similar argument concerning the uniqueness of the *split* combinator $\langle f, g \rangle$.

What about other facts valid for numeric exponentials such as $a^0 = 1$ and $1^a = 1$? The reader is invited to go back to section 2.11 and recall what 0 and 1 mean as datatypes: the empty (void) and singleton datatypes, respectively. Our counterpart to $a^0 = 1$ then is

$$A^0 \cong 1 \quad (2.96)$$

where A^0 denotes the set of all functions from the empty set to some A . What does (2.96) mean? It simply tells us that there is only one

¹³ Writing \overline{f} (resp. \widehat{f}) or $\text{curry } f$ (resp. $\text{uncurry } f$) is a matter of taste: the latter are more in the tradition of functional programming and help when the functions have to be named; the former save ink in algebraic expressions and calculations.

function in such a set — the empty function mapping “no” value at all. This fact confirms our choice of notation once again (compare with $a^0 = 1$ in a numeric context).

Next, we may wonder about facts

$$1^A \cong 1 \quad (2.97)$$

$$A^1 \cong A \quad (2.98)$$

which are the functional exponentiation counterparts of $1^a = 1$ and $a^1 = a$. Fact (2.97) is valid: it means that there is only one function mapping A to some singleton set $\{c\}$ — the constant function \underline{c} . There is no room for another function in 1^A because only c is available as output value. Our standard denotation for such a unique function is given by (2.58).

Fact (2.98) is also valid: all functions in A^1 are (single valued) constant functions and there are as many constant functions in such a set as there are elements in A . These functions are often called (abstract) “points” because of the 1-to-1 mapping between A^1 and the elements (points) in A .

Exercise 2.25. Relate the isomorphism involving generic elementary type 2

$$A \times A \cong A^2 \quad (2.99)$$

to the expression `\f->(f True, f False)` written in HASKELL syntax.

□

Exercise 2.26. Consider the witnesses of isomorphism (2.95)

$$(B \times C)^A \begin{array}{c} \xrightarrow{\text{unpair}} \\ \cong \\ \xleftarrow{\text{pair}} \end{array} B^A \times C^A$$

defined by:

$$\begin{aligned} \text{pair } (f, g) &= \langle f, g \rangle \\ \text{unpair } k &= (\pi_1 \cdot k, \pi_2 \cdot k) \end{aligned}$$

Show that $\text{pair} \cdot \text{unpair} = \text{id}$ and $\text{unpair} \cdot \text{pair} = \text{id}$ hold.

□

Exercise 2.27. Show that the following equality

$$\bar{f} a = f \cdot \langle \underline{a}, \text{id} \rangle \quad (2.100)$$

holds.

□

Exercise 2.28. Considering $\alpha = [\bar{i}_1, \bar{i}_2]$, (a) infer the principal (most general) type of α and depict it in a diagram; (b) $\hat{\alpha}$ is a well-known isomorphism — tell which by inferring its type.

□

Exercise 2.29. Prove the equality $\underline{g} = \overline{g \cdot \pi_2}$ knowing that

$$\overline{\pi_2} = \underline{id} \quad (2.101)$$

holds.

□

2.16 FINITARY PRODUCTS AND COPRODUCTS

In section 2.8 it was suggested that product could be regarded as the abstraction behind data-structuring primitives such as `struct` in C or `record` in PASCAL. Similarly, coproducts were suggested in section 2.9 as abstract counterparts of C unions or PASCAL variant records. For a finite A , exponential B^A could be realized as an *array* in any of these languages. These analogies are captured in table 1.

In the same way C `structs` and `unions` may contain finitely many entries, as may PASCAL (variant) records, product $A \times B$ extends to finitary product $A_1 \times \dots \times A_n$, for $n \in \mathbb{N}$, also denoted by $\prod_{i=1}^n A_i$, to which as many projections π_i are associated as the number n of factors involved. Of course, *splits* become n -ary as well

$$\langle f_1, \dots, f_n \rangle : A_1 \times \dots \times A_n \leftarrow B$$

for $f_i : A_i \leftarrow B, i = 1, n$.

Dually, coproduct $A + B$ is extensible to the finitary sum $A_1 + \dots + A_n$, for $n \in \mathbb{N}$, also denoted by $\sum_{j=1}^n A_j$, to which as many injections i_j are assigned as the number n of terms involved. Similarly, *either*s become n -ary

$$[f_1, \dots, f_n] : A_1 + \dots + A_n \rightarrow B$$

for $f_i : B \leftarrow A_i, i = 1, n$.

Datatype n

Next after 2, we may think of 3 as representing the abstract class of all datatypes containing exactly three elements. Generalizing, we may

Abstract notation	PASCAL	C/C++	Description
$A \times B$	<pre>record P: A; S: B; end;</pre>	<pre>struct { A first; B second; };</pre>	Records
$A + B$	<pre>record case tag: integer of x = 1: (P:A); 2: (S:B); end;</pre>	<pre>struct { int tag; /* 1,2 */ union { A ifA; B ifB; } data; };</pre>	Variant records
B^A	<code>array[A] of B</code>	<code>B ...[A]</code>	Arrays
$1 + A$	<code>^A</code>	<code>A *...</code>	Pointers

Table 1.: Abstract notation versus programming language data-structures.

think of n as representing the abstract class of all datatypes containing exactly n elements. Of course, initial segment n will be in this abstract class. (Recall (2.17), for instance: both `Weekday` and `7` are abstractly represented by `7`.) Therefore,

$$n \cong \underbrace{1 + \dots + 1}_n$$

and

$$\underbrace{A \times \dots \times A}_n \cong A^n \quad (2.102)$$

$$\underbrace{A + \dots + A}_n \cong n \times A \quad (2.103)$$

hold.

Exercise 2.30. On the basis of table 1, encode `undistr` (2.51) in C or PASCAL. Compare your code with the HASKELL *pointfree* and *pointwise* equivalents.

□

2.17 INITIAL AND TERMINAL DATATYPES

All properties studied for binary *splits* and binary *either*s extend to the finitary case. For the particular situation $n = 1$, we will have $\langle f \rangle = [f] = f$ and $\pi_1 = i_1 = id$, of course. For the particular situation $n = 0$, finitary products “degenerate” to 1 and finitary coproducts “degenerate” to 0. So diagrams (2.23) and (2.38) are reduced to

$$\begin{array}{ccc} 1 & & 0 \\ \langle \rangle \uparrow & & \downarrow [] \\ C & & C \end{array}$$

The standard notation for the empty *split* $\langle \rangle$ is $!_C$, where subscript C can be omitted if implicit in the context. By the way, this is precisely the only function in 1^C , recall (2.97) and (2.58). Dually, the standard notation for the empty *either* $[]$ is $?_C$, where subscript C can also be omitted. By the way, this is precisely the only function in C^0 , recall (2.96).

In summary, we may think of 0 and 1 as, in a sense, the “extremes” of the whole datatype spectrum. For this reason they are called *initial* and *terminal*, respectively. We conclude this subject with the presentation of their main properties which, as we have said, are instances of properties we have stated for products and coproducts.

Initial datatype reflexion :

$$\begin{array}{c} ?_0 = id_0 \\ \curvearrowright \\ 0 \end{array} \qquad ?_0 = id_0 \qquad (2.104)$$

Initial datatype fusion :

$$\begin{array}{ccc} 0 & & \\ ?_A \downarrow & \searrow ?_B & \\ A & \xrightarrow{f} & B \end{array} \qquad f \cdot ?_A = ?_B \qquad (2.105)$$

Terminal datatype reflexion :

$$\begin{array}{c} !_1 = id_1 \\ \curvearrowleft \\ 1 \end{array} \qquad !_1 = id_1 \qquad (2.106)$$

Terminal datatype fusion :

$$\begin{array}{ccc} 1 & & \\ !_A \uparrow & \nwarrow !_B & \\ A & \xleftarrow{f} & B \end{array} \qquad !_A \cdot f = !_B \qquad (2.107)$$

Exercise 2.31. Particularize the exchange law (2.49) to empty products and empty coproducts, i.e. 1 and 0.

□

2.18 SUMS AND PRODUCTS IN HASKELL

We conclude this chapter with an analysis of the main primitive available in HASKELL for creating datatypes: the `data` declaration. Suppose we declare

```
data CostumerId = P  $\mathbb{Z}$  | C  $\mathbb{Z}$ 
```

meaning to say that, for some company, a client is identified either by its passport number or by its credit card number, if any. What does this piece of syntax precisely mean?

If we enquire the HASKELL *interpreter* about what it knows about `CostumerId`, the reply will contain the following information:

```
Main> :i CostumerId
-- type constructor
data CostumerId

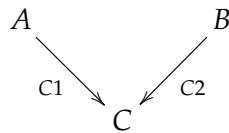
-- constructors:
P :: Int -> CostumerId
C :: Int -> CostumerId
```

In general, let A and B be two known datatypes. Via declaration

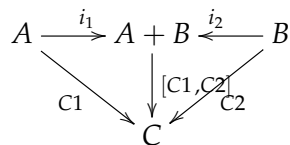
```
data C = C1 A | C2 B
```

(2.108)

one obtains from HASKELL a new datatype C equipped with constructors $C \xleftarrow{C1} A$ and $C \xleftarrow{C2} B$, in fact the only ones available for constructing values of C :



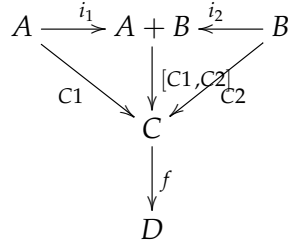
This diagram leads to an obvious instance of coproduct diagram (2.38),



describing that a `data` declaration in HASKELL means the *either* of its constructors.

Because there are no other means to build C data, it follows that C is isomorphic to $A + B$. So $[C1, C2]$ has an inverse, say inv , which is such that $inv \cdot [C1, C2] = id$. How do we calculate inv ? Let us first

think of the generic situation of a function $D \xleftarrow{f} C$ which observes datatype C :



This is an opportunity for $+fusion$ (2.42), whereby we obtain

$$f \cdot [C1, C2] = [f \cdot C1, f \cdot C2]$$

Therefore, the observation will be fully described provided we explain how f behaves with respect to $C1$ — *cf.* $f \cdot C1$ — and with respect to $C2$ — *cf.* $f \cdot C2$. This is what is behind the typical *inductive* structure of pointwise f , which will be made of two and only two clauses:

$$\begin{aligned} f &: C \rightarrow D \\ f(C1\ a) &= \dots \\ f(C2\ b) &= \dots \end{aligned}$$

Let us use this in calculating the inverse inv of $[C1, C2]$:

$$\begin{aligned} inv \cdot [C1, C2] &= id \\ \equiv \quad \{ \text{by } +fusion \text{ (2.42)} \} \\ [inv \cdot C1, inv \cdot C2] &= id \\ \equiv \quad \{ \text{by } +-reflexion \text{ (2.41)} \} \\ [inv \cdot C1, inv \cdot C2] &= [i_1, i_2] \\ \equiv \quad \{ \text{either structural equality (2.66)} \} \\ inv \cdot C1 &= i_1 \wedge inv \cdot C2 = i_2 \end{aligned}$$

Therefore:

$$\begin{aligned} inv &: C \rightarrow A + B \\ inv(C1\ a) &= i_1\ a \\ inv(C2\ b) &= i_2\ b \end{aligned}$$

In summary, $C1$ is a “renaming” of injection i_1 , $C2$ is a “renaming” of injection i_2 and C is a “renamed” replica of $A + B$:

$$C \xleftarrow{[C1, C2]} A + B \tag{2.109}$$

$[C1, C2]$ is called the *algebra* of datatype C and its inverse inv is called the *coalgebra* of C . The algebra contains the constructors $C1$ and $C2$ of

type C , that is, it is used to “build” C -values. In the opposite direction, co-algebra inv enables us to “destroy” or observe values of C :

$$\begin{array}{ccc} & \xrightarrow{inv} & \\ C & \cong & A + B \\ & \xleftarrow{[C1,C2]} & \end{array}$$

Algebra/coalgebras also arise about product datatypes. For instance, suppose that one wishes to describe datatype *Point* inhabited by pairs $(x_0, y_0), (x_1, y_1)$ (etc.) of Cartesian coordinates of a given type, say A . Although $A \times A$ equipped with projections π_1, π_2 “is” such a datatype, one may be interested in a suitably named replica of $A \times A$ in which points are built explicitly by some constructor (say *Point*) and observed by dedicated selectors (say x and y):

$$\begin{array}{ccccc} A & \xleftarrow{\pi_1} & A \times A & \xrightarrow{\pi_2} & A \\ & \searrow x & \downarrow \text{Point} & \nearrow y & \\ & & \text{Point} & & \end{array} \quad (2.110)$$

This gives birth to the algebra *Point* and the coalgebra $\langle x, y \rangle$ of datatype *Point*:

$$\begin{array}{ccc} & \xrightarrow{\langle x, y \rangle} & \\ \text{Point} & \cong & A \times A \\ & \xleftarrow{\text{Point}} & \end{array}$$

In HASKELL one writes

data *Point* $a = \text{Point } \{ x :: a, y :: a \}$

but be warned that HASKELL delivers *Point* in curried form:

Point $:: a \rightarrow a \rightarrow \text{Point } a$

Finally, what is the “HASKELL-equivalent” to handling a *pointer* in (say) C ? This corresponds to $A = 1$ in (2.109),

$$C \xleftarrow{[C1,C2]} 1 + B$$

and to the following HASKELL declaration:

data $C = C1 () \mid C2 B$

Note that HASKELL allows for a more programming-oriented alternative in this case, in which the unit type $()$ is eliminated:

data $C = C1 \mid C2 B$

The difference is that here $C1$ denotes an inhabitant of C (and so a clause $f(C1 a) = \dots$ is rewritten to $f C1 = \dots$) while above $C1$ denotes a (constant) function $C \xleftarrow{C1} 1$. Isomorphism (2.98) helps in comparing these two alternative situations.

2.19 EXERCISES

Exercise 2.32. Let A and B be two disjoint datatypes, that is, $A \cap B = \emptyset$ holds. Show that isomorphism

$$A \cup B \cong A + B \quad (2.111)$$

holds. **Hint:** define $A \cup B \xleftarrow{i} A + B$ as $i = [\text{emb}_A, \text{emb}_B]$ for $\text{emb}_A a = a$ and $\text{emb}_B b = b$, and find its inverse. By the way, why didn't we define i as simply as $i \stackrel{\text{def}}{=} [\text{id}_A, \text{id}_B]$?

□

Exercise 2.33. Knowing that a given function xr satisfies property

$$\text{xr} \cdot \langle \langle f, g \rangle, h \rangle = \langle \langle f, h \rangle, g \rangle \quad (2.112)$$

for all f, g and h , derive from (2.112) the definition of xr :

$$\text{xr} = \langle \pi_1 \times \text{id}, \pi_2 \cdot \pi_1 \rangle \quad (2.113)$$

□

Exercise 2.34. Let distr (read: 'distribute right') be the bijection which witnesses isomorphism $A \times (B + C) \cong A \times B + A \times C$. Fill in the "... " in the diagram which follows so that it describes bijection distl (red: 'distribute left') which witnesses isomorphism $(B + C) \times A \cong B \times A + C \times A$:

$$(B + C) \times A \xrightarrow{\text{swap}} \dots \xrightarrow{\text{distr}} \dots \xrightarrow{\dots} B \times A + C \times A$$

$\xrightarrow{\text{distl}}$

□

Exercise 2.35. In the context of exercise 2.34, prove

$$[g, h] \times f = [g \times f, h \times f] \cdot \text{distl} \quad (2.114)$$

knowing that

$$f \times [g, h] = [f \times g, f \times h] \cdot \text{distr} \quad (2.115)$$

holds.

□

Exercise 2.36. The arithmetic law $(a + b)(c + d) = (ac + ad) + (bc + bd)$ corresponds to the isomorphism

$$(A+B) \times (C+D) \cong (A \times C + A \times D) + (B \times C + B \times D)$$

$$h = [[i_1 \times i_1, i_1 \times i_2], [i_2 \times i_1, i_2 \times i_2]]$$

From universal property (2.65) infer the following definition of function h , written in Haskell syntax:

$$\begin{aligned} h(\text{Left}(\text{Left}(a, c))) &= (\text{Left } a, \text{Left } c) \\ h(\text{Left}(\text{Right}(a, d))) &= (\text{Left } a, \text{Right } d) \\ h(\text{Right}(\text{Left}(b, c))) &= (\text{Right } b, \text{Left } c) \\ h(\text{Right}(\text{Right}(b, d))) &= (\text{Right } b, \text{Right } d) \end{aligned}$$

☐

Exercise 2.37. Every C programmer knows that a struct of pointers

$$(A + 1) \times (B + 1)$$

offers a data type which represents both product $A \times B$ (struct) and coproduct $A + B$ (union), alternatively. Express in pointfree notation the isomorphisms i_1 to i_5 of

$$\begin{array}{ccc}
 (A+1) \times (B+1) & \xleftarrow{i_1} & ((A+1) \times B) + ((A+1) \times 1) \\
 & & \uparrow i_2 \\
 & & (A \times B + 1 \times B) + (A \times 1 + 1 \times 1) \\
 & & \uparrow i_3 \\
 & & (A \times B + B) + (A + 1) \\
 & & \uparrow i_4 \\
 (A \times B + (B + A)) + 1 & \xrightarrow{i_5} & A \times B + (B + (A + 1))
 \end{array}$$

which witness the observation above.

□

Exercise 2.38. Prove the following property of McCarthy conditionals:

$$p \rightarrow f \cdot g, h \cdot k = [f, h] \cdot (p \rightarrow i_1 \cdot g, i_2 \cdot k) \quad (2.116)$$

☐

Exercise 2.39. Assuming the fact

$$(p^? + p^?) \cdot p^? = (i_1 + i_2) \cdot p^? \quad (2.117)$$

show that nested conditionals can be simplified:

$$p \rightarrow (p \rightarrow f, g), (p \rightarrow h, k) = p \rightarrow f, k \quad (2.118)$$

□

Exercise 2.40. Show that $(\overline{f \cdot ap})g = f \cdot g$ holds, cf. (2.87).

□

Exercise 2.41. Consider the higher-order isomorphism *flip* defined as follows:

$$\begin{aligned} (C^B)^A &\cong C^{A \times B} \cong C^{B \times A} \cong (C^A)^B \\ f &\mapsto \widehat{f} \mapsto \widehat{f}.\text{swap} \mapsto \overline{\widehat{f} \cdot \text{swap}} = \text{flip } f \end{aligned}$$

Show that $\text{flip } f \ x \ y = f \ y \ x$.

□

Exercise 2.42. Let $C \xrightarrow{\text{const}} C^A$ be the function of exercise 2.2, that is, $\text{const } c = c_A$. Which fact is expressed by the following diagram featuring *const*?

$$\begin{array}{ccc} C & \xrightarrow{\text{const}} & C^A \\ f \downarrow & & \downarrow f^A \\ B & \xrightarrow{\text{const}} & B^A \end{array} \quad (2.119)$$

Write it at point-level and describe it by your own words.

□

Exercise 2.43. Show that $\overline{\pi_2} \cdot f = \overline{\pi_2}$ holds for every f . Thus $\overline{\pi_2}$ is a constant function — which one?

□

Exercise 2.44. Establish the difference between the following two declarations in HASKELL,

```
data D = D1 A | D2 B C
data E = E1 A | E2 (B,C)
```

for A, B and C any three predefined types. Are D and E isomorphic? If so, can you specify and encode the corresponding isomorphism?

□

2.20 BIBLIOGRAPHY NOTES

A few decades ago John Backus read, in his Turing Award Lecture, a revolutionary paper [7]. This paper proclaimed conventional command-oriented programming languages obsolete because of their inefficiency arising from retaining, at a high-level, the so-called “memory access bottleneck” of the underlying computation model — the well-known *von Neumann* architecture. Alternatively, the (at the time already mature) *functional programming* style was put forward for two main reasons. Firstly, because of its potential for concurrent and parallel computation. Secondly — and Backus emphasis was really put on this —, because of its strong algebraic basis.

Backus *algebra of (functional) programs* was providential in alerting computer programmers that computer languages alone are insufficient, and that only languages which exhibit an *algebra* for reasoning about the objects they purport to describe will be useful in the long run.

The impact of Backus first argument in the computing science and computer architecture communities was considerable, in particular if assessed in quality rather than quantity and in addition to the almost contemporary *structured programming* trend ¹⁴. By contrast, his second argument for changing computer programming was by and large ignored, and only the so-called *algebra of programming* research minorities pursued in this direction. However, the advances in this area throughout the last two decades are impressive and can be fully appreciated by reading a textbook written relatively recently by Bird and de Moor [10]. A comprehensive review of the voluminous literature available in this area can also be found in this book.

Although the need for a pointfree algebra of programming was first identified by Backus, perhaps influenced by Iverson’s APL growing popularity in the USA at that time, the idea of reasoning and using mathematics to transform programs is much older and can be traced to the times of McCarthy’s work on the foundations of computer programming [38], of Floyd’s work on program meaning [15] and of Paterson and Hewitt’s *comparative schematology* [55]. Work of the so-called *program transformation* school was already very expressive in the mid 1970s, see for instance references [11].

The mathematics adequate for the effective integration of these related but independent lines of thought was provided by the categorical approach of Manes and Arbib compiled in a textbook [37] which has very strongly influenced the last decade of 20th century theoretical computer science.

A so-called MPC (“Mathematics of Program Construction”) community has been among the most active in producing an integrated body

¹⁴ Even the C programming language and the UNIX operating system, with their implicit functional flavour, may be regarded as subtle outcomes of the “going functional” trend.

of knowledge on the algebra of programming which has found in functional programming an eloquent and paradigmatic medium. Functional programming has a tradition of absorbing fresh results from theoretical computer science, algebra and category theory. Languages such as HASKELL [9] have been competing to integrate the most recent developments and therefore are excellent *prototyping* vehicles in courses on program calculation, as happens with this book.

For fairly recent work on this topic see e.g. [19, 24, 25, 18].

RECURSION IN THE POINTFREE STYLE

How useful from a programmer's point of view are the abstract concepts presented in the previous chapter? Recall that a table was presented — table 1 — which records an analogy between abstract type notation and the corresponding data-structures available in common, imperative languages.

This analogy will help in showing how to extend the abstract notation studied thus far towards a most important field of programming: *recursion*. This, however, will be preceded by a simpler introduction to the subject rooted on very basic and intuitive notions of mathematics.

3.1 MOTIVATION

Where do algorithms come from? From human imagination only? Surely not — they actually emerge from mathematics. In a sense, in the same way one may say that hardware follows the laws of physics (eg. semiconductor electronics) one might say that software is governed by the laws of mathematics.

This section provides a naive introduction to algorithm analysis and synthesis by showing how a quite elementary class of algorithms — equivalent to for-loops in C or any other imperative language — arise from elementary properties of the underlying maths domain.

We start by showing how the arithmetic operation of multiplying two natural numbers (in \mathbb{N}_0) is a for-loop which emerges solely from the algebraic properties of multiplication:

$$\left\{ \begin{array}{l} a \times 0 = 0 \\ a \times 1 = a \\ a \times (b + c) = a \times b + a \times c \end{array} \right. \quad (3.1)$$

These properties are known as the *absorption*, *unit* and *distributive* properties of multiplication, respectively.

Start by making $c := 1$ in the third (distributive) property, obtaining $a \times (b + 1) = a \times b + a \times 1$, and then simplify. The second clause is

useful in this simplification but it is not required in the final system of two equations,

$$\begin{cases} a \times 0 = 0 \\ a \times (b + 1) = a \times b + a \end{cases} \quad (3.2)$$

since it is derivable from the remaining two, for $b := 0$ and property $0 + a = a$ of addition.

System (3.2) is *already* a runnable program in a functional language such as Haskell (among others). The moral of this trivial exercise is that programs *arise* from the underlying maths, instead of being *invented* or coming out of the blue. Novices in functional programming do this kind of reasoning all the time without even noticing it, when writing their first programs. For instance, the function which computes discrete exponentials will scale up the same procedure, thanks to the properties

$$\begin{cases} a^0 = 1 \\ a^1 = a \\ a^{b+c} = a^b \times a^c \end{cases}$$

where the program just developed for multiplication can be re-used, and so and so on.

Type-wise, the multiplication algorithm just derived for natural numbers is not immediate to generalize. Intuitively, it will diverge for b a negative integer and for b a real number less than 1, at least. Argument a , however, does not seem to be constrained.

Indeed, the two arguments a and b will have different types in general. Let us see why and how. Starting by looking at infix operators (\times) and $(+)$ as *curried* operators — recall section 2.15 — we can resort to the corresponding *sections* and write:

$$\begin{cases} (a \times) 0 = 0 \\ (a \times) (b + 1) = (a +) ((a \times) b) \end{cases} \quad (3.3)$$

It can be easily checked that

$$(a \times) = \text{for } (a +) 0 \quad (3.4)$$

by introducing a **for-loop** combinator given by

$$\begin{cases} \text{for } f \ i \ 0 = i \\ \text{for } f \ i \ (n + 1) = f \ (\text{for } f \ i \ n) \end{cases} \quad (3.5)$$

where f is the loop-body and i is the initialization value. In fact, $(\text{for } f \ i) n = f^n i$, that is, f is iterated n times over the initial value i .

For-loops are a primitive construct available in many programming languages. In C, for instance, one will write something like

```

int mul(int a, int n)
{
  int s=0; int i;
  for (i=1; i<n+1; i++) {s += a;}
  return s;
};

```

for (the uncurried version of) for $(a+)$ 0 loop.

To better understand this construct let us remove variables from both equations in (3.3) by lifting function application to function composition and lifting 0 to the “everywhere 0” (constant) function:

$$\begin{cases} (a \times) \cdot \underline{0} = \underline{0} \\ (a \times) \cdot (+1) = (+a) \cdot (a \times) \end{cases}$$

Using the *junc* (“either”) pointfree combinator we merge the two equations into a single one,

$$[(a \times) \cdot \underline{0}, (a \times) \cdot (+1)] = [\underline{0}, (+a) \cdot (a \times)]$$

— thanks to the Eq-+ rule (2.66) — then single out the common factor $(a \times)$ in the left hand side,

$$(a \times) \cdot [\underline{0}, (+1)] = [\underline{0}, (+a) \cdot (a \times)]$$

— thanks to +-fusion (2.42) — and finally do a similar *fission* operation on the other side,

$$(a \times) \cdot [\underline{0}, (+1)] = [\underline{0}, (+a)] \cdot (id + (a \times)) \quad (3.6)$$

— thanks to +-absorption (2.43).

As we already know, equalities of compositions are nicely drawn as diagrams. That of (3.6) is as follows:

$$\begin{array}{ccc} \mathbb{N}_0 & \xleftarrow{[\underline{0}, (+1)]} & A + \mathbb{N}_0 \\ (a \times) \downarrow & & \downarrow id + (a \times) \\ \mathbb{N}_0 & \xleftarrow{[\underline{0}, (+a)]} & A + \mathbb{N}_0 \end{array}$$

Function $(+1)$ is the successor function *succ* on natural numbers. Type A is any (non-empty) type. For the particular case of $A = 1$, the diagram is more interesting, as $[\underline{0}, \text{succ}]$ becomes an isomorphism, telling a *unique* way of building natural numbers:¹

Every natural number in \mathbb{N}_0 either is 0 or the successor of another natural number.

¹ This is nothing but a re-statement of the well-known *Peano* axioms for the natural numbers. Giuseppe Peano (1858-1932) was a famous Italian mathematician.

We will denote such an isomorphism by *in* and its converse by *out* in the following version of the same diagram

$$\begin{array}{ccc}
 \mathbb{N}_0 & \xrightarrow{\text{out}=\text{in}^\circ} & 1 + \mathbb{N}_0 \\
 \downarrow (a \times) & \begin{array}{c} \cong \\ \text{in}=[\underline{0}, \text{succ}] \end{array} & \downarrow \text{id}+(a \times) \\
 \mathbb{N}_0 & \xleftarrow{[\underline{0}, (+a)]} & 1 + \mathbb{N}_0
 \end{array}$$

capturing both the isomorphism and the $(a \times)$ recursive function. By solving the isomorphism equation $\text{out} \cdot \text{in} = \text{id}$ we easily obtain the definition of *out*, the converse of *in*²:

$$\begin{aligned}
 \text{out } 0 &= i_1() \\
 \text{out}(n+1) &= i_2 n
 \end{aligned}$$

Finally, we generalize the target \mathbb{N}_0 to any non-empty type B , $(+a)$ to any function $B \xrightarrow{g} B$ and 0 to any constant k in B (this is why B has to be non-empty). The corresponding generalization of $(a \times)$ is denoted by f below:

$$\begin{array}{ccc}
 \mathbb{N}_0 & \xrightarrow{\text{out}=\text{in}^\circ} & 1 + \mathbb{N}_0 \\
 \downarrow f & \begin{array}{c} \cong \\ \text{in}=[\underline{0}, \text{succ}] \end{array} & \downarrow \text{id}+f \\
 B & \xleftarrow{[\underline{k}, g]} & 1 + B
 \end{array}$$

It turns out that, given k and g , there is a unique solution to the equation (in f) captured by the diagram: $f \cdot \text{in} = [\underline{k}, g] \cdot (\text{id} + f)$. We know this solution already, recall (3.5):

$$f = \text{for } g \ k$$

As we have seen earlier on, solution uniqueness is captured by universal properties. In this case we have the following property, which we will refer to by writing “for-loop-universal”:

$$f = \text{for } g \ k \quad \equiv \quad f \cdot \text{in} = [\underline{k}, g] \cdot (\text{id} + f) \quad (3.7)$$

From this property it is possible to infer a basic theory of for-loops. For instance, by making $f = \text{id}$ and solving the for-loop-universal equation (3.7) for g and k we obtain the reflexion law:

$$\text{for succ } 0 = \text{id} \quad (3.8)$$

This can be compared with the following (useless) program in C:

² Note how the singularity of type 1 ensures *out* a function: what would the outcome of *out* 0 be should A be arbitrary?

```

int id(int n)
{
  int s=0; int i;
  for (i=1; i<n+1; i++) {s += 1;}
  return s;
};

```

(Clearly, the value returned in s is that of input n .)

More knowledge about for-loops can be extracted from (3.7). Later on we will show that these constructs are special cases of a more general concept termed *catamorphism*.³ In the usual “*banana-bracket*” notation of catamorphisms, to be introduced later, the for-combinator will be written:

$$\text{for } g \ k = ([k, g]) \quad (3.9)$$

In the sequel, we shall study the (more general) theory of catamorphisms and come back to for-loops as an instantiation. Then we will understand how more interesting for-loops can be synthesized, for instance those handling more than one “global variable”, thanks to catamorphism theory (for instance, the mutual recursion laws).

As a generalization to what we’ve just seen happening between for-loops and natural numbers, it will be shown that a catamorphism is intimately connected to the data-structure it processes, for instance a finite list (sequence) or a binary tree. A good understanding of such structures is therefore required. We proceed to studying the list data structure first, wherefrom trees stem as natural extensions.

Exercise 3.1. Addition is known to be associative ($a + (b + c) = (a + b) + c$) and have unit 0 ($a + 0 = a$). Following the same strategy that was adopted above for $(a \times)$, show that

$$(a+) = \text{for succ } a \quad (3.10)$$

□

Exercise 3.2. The following fusion-law

$$h \cdot (\text{for } g \ k) = \text{for } j \ (h \ k) \iff h \cdot g = j \cdot h \quad (3.11)$$

can be derived from universal-property (3.7)⁴. Since $(a+) \cdot \text{id} = (a+)$, provide an alternative derivation of (3.10) using the fusion-law above.

□

Exercise 3.3. From (3.4) and fusion-law (3.11) infer: $(a*) \cdot \text{succ} = \text{for } a \ (a+)$.

□

³ See eg. section 3.6.

⁴ A generalization of this property will be derived in section 3.12.

Exercise 3.4. Show that $f = \text{for } \underline{k} \ k$ and $g = \text{for } id \ k$ are the same program (function).

□

Exercise 3.5. Generic function $k = \text{for } f \ i$ can be encoded in the syntax of C by writing

```
int k(int n) {
    int r=i;
    int x;
    for (x=1; x<n+1; x++) {r=f(x);}
    return r;
};
```

for some predefined f . Encode the functions f and g of exercise 3.4 in C and compare them.

□

3.2 FROM NATURAL NUMBERS TO FINITE SEQUENCES

Let us consider a very common data-structure in programming: “linked-lists”. In PASCAL one will write

```
L = ^N;
N = record
    first: A;
    next: ^N
end;
```

to specify such a data-structure L . This consists of a pointer to a *node* (N), where a node is a record structure which puts some predefined type A together with a pointer to another node, and so on. In the C programming language, every $x \in L$ will be declared as $L \ x$ in the context of datatype definition

```
typedef struct N {
    A first;
    struct N *next;
} *L;
```

and so on.

What interests us in such “first year programming course” datatype declarations? Records and pointers have already been dealt with in table 1. So we can use this table to find the abstract version of datatype

L , by replacing pointers by the “ $1 + \dots$ ” notation and records (*structs*) by the “ $\dots \times \dots$ ” notation:

$$\begin{cases} L &= 1 + N \\ N &= A \times (1 + N) \end{cases} \quad (3.12)$$

We obtain a system of two equations on unknowns L and N , in which L 's dependence on N can be removed by substitution:

$$\begin{aligned} &\begin{cases} L &= 1 + N \\ N &= A \times (1 + N) \end{cases} \\ \equiv &\quad \{ \text{substituting } L \text{ for } 1 + N \text{ in the second equation} \} \\ &\begin{cases} L &= 1 + N \\ N &= A \times L \end{cases} \\ \equiv &\quad \{ \text{substituting } A \times L \text{ for } N \text{ in the first equation} \} \\ &\begin{cases} L &= 1 + A \times L \\ N &= A \times L \end{cases} \end{aligned}$$

System (3.12) is thus equivalent to:

$$\begin{cases} L &= 1 + A \times L \\ N &= A \times (1 + N) \end{cases} \quad (3.13)$$

Intuitively, L abstracts the “possibly empty” linked-list of elements of type A , while N abstracts the “non-empty” linked-list of elements of type A . Note that L and N are independent of each other, but also that each depends on itself. Can we solve these equations in a way such that we obtain “solutions” for L and N , in the same way we do with school equations such as, for instance,

$$x = 1 + \frac{x}{2} \quad ? \quad (3.14)$$

Concerning this equation, let us recall how we would go about it in school mathematics:

$$\begin{aligned} &x = 1 + \frac{x}{2} \\ \equiv &\quad \{ \text{adding } -\frac{x}{2} \text{ to both sides of the equation} \} \\ &x - \frac{x}{2} = 1 + \frac{x}{2} - \frac{x}{2} \\ \equiv &\quad \{ -\frac{x}{2} \text{ cancels } \frac{x}{2} \} \\ &x - \frac{x}{2} = 1 \\ \equiv &\quad \{ \text{multiplying both sides of the equation by 2 etc.} \} \\ &2 \times x - x = 2 \\ \equiv &\quad \{ \text{subtraction} \} \\ &x = 2 \end{aligned}$$

We very quickly get solution $x = 2$. However, many steps were omitted from the actual calculation. This unfolds into the longer sequence of more elementary steps which follows, in which notation $a - b$ abbreviates $a + (-b)$ and $\frac{a}{b}$ abbreviates $a \times \frac{1}{b}$, for $b \neq 0$:

$$\begin{aligned}
 & x = 1 + \frac{x}{2} \\
 \equiv & \quad \{ \text{adding } -\frac{x}{2} \text{ to both sides of the equation} \} \\
 & x - \frac{x}{2} = (1 + \frac{x}{2}) - \frac{x}{2} \\
 \equiv & \quad \{ + \text{ is associative} \} \\
 & x - \frac{x}{2} = 1 + (\frac{x}{2} - \frac{x}{2}) \\
 \equiv & \quad \{ -\frac{x}{2} \text{ is the additive inverse of } \frac{x}{2} \} \\
 & x - \frac{x}{2} = 1 + 0 \\
 \equiv & \quad \{ 0 \text{ is the unit of addition} \} \\
 & x - \frac{x}{2} = 1 \\
 \equiv & \quad \{ \text{multiplying both sides of the equation by 2} \} \\
 & 2 \times (x - \frac{x}{2}) = 2 \times 1 \\
 \equiv & \quad \{ 1 \text{ is the unit of multiplication} \} \\
 & 2 \times (x - \frac{x}{2}) = 2 \\
 \equiv & \quad \{ \text{multiplication distributes over addition} \} \\
 & 2 \times x - 2 \times \frac{x}{2} = 2 \\
 \equiv & \quad \{ 2 \text{ cancels its inverse } \frac{1}{2} \} \\
 & 2 \times x - 1 \times x = 2 \\
 \equiv & \quad \{ \text{multiplication distributes over addition} \} \\
 & (2 - 1) \times x = 2 \\
 \equiv & \quad \{ 2 - 1 = 1 \text{ and } 1 \text{ is the unit of multiplication} \} \\
 & x = 2
 \end{aligned}$$

Back to (3.13), we would like to submit each of the equations, *e.g.*

$$L = 1 + A \times L \quad (3.15)$$

to a similar reasoning. Can we do it? The analogy which can be found between this equation and (3.14) goes beyond pattern similarity. From chapter 2 we know that many properties required in the reasoning above hold in the context of (3.15), provided the “=” sign is replaced by the “ \cong ” sign, that of set-theoretical isomorphism. Recall that, for

instance, $+$ is associative (2.48), 0 is the unit of addition (2.55), 1 is the unit of multiplication (2.57), multiplication distributes over addition (2.52) *etc.* Moreover, the first step above assumed that addition is compatible (monotonic) with respect to equality,

$$\frac{\begin{array}{ccc} a & = & b \\ c & = & d \end{array}}{a + c = b + d}$$

a fact which still holds when numeric equality gives place to isomorphism and numeric addition gives place to coproduct:

$$\frac{\begin{array}{ccc} A & \cong & B \\ C & \cong & D \end{array}}{A + C \cong B + D}$$

— recall (2.46) for isos f and g .

Unfortunately, the main steps in the reasoning above are concerned with two basic *cancellation properties*

$$\begin{aligned} x + b = c &\equiv x = c - b \\ x \times b = c &\equiv x = \frac{c}{b} \quad (b \neq 0) \end{aligned}$$

which hold about numbers but do not hold about datatypes. In fact, neither products nor coproducts have arbitrary inverses⁵, and so we cannot “calculate by cancellation”. How do we circumvent this limitation?

Just think of how we would have gone about (3.14) in case we didn’t know about the *cancellation properties*: we would be bound to the x by $1 + \frac{x}{2}$ substitution plus the other properties. By performing such a substitution over and over again we would obtain...

$$\begin{aligned} x &= 1 + \frac{x}{2} \\ &\equiv \{ x \text{ by } 1 + \frac{x}{2} \text{ substitution followed by simplification} \} \\ x &= 1 + \frac{1 + \frac{x}{2}}{2} = 1 + \frac{1}{2} + \frac{x}{4} \\ &\equiv \{ \text{the same as above} \} \\ x &= 1 + \frac{1}{2} + \frac{1 + \frac{x}{2}}{4} = 1 + \frac{1}{2} + \frac{1}{4} + \frac{x}{8} \\ &\equiv \{ \text{over and over again, } n\text{-times} \} \\ &\dots \\ &\equiv \{ \text{simplification} \} \\ x &= \sum_{i=0}^n \frac{1}{2^i} + \frac{x}{2^{n+1}} \end{aligned}$$

⁵ The initial and terminal datatypes do have inverses — 0 is its own “additive inverse” and 1 is its own “multiplicative inverse” — but not all the others.

$$\begin{aligned}
&\equiv \{ \text{sum of } n \text{ first terms of a geometric progression} \} \\
&x = (2 - \frac{1}{2^n}) + \frac{x}{2^{n+1}} \\
&\equiv \{ \text{let } n \rightarrow \infty \} \\
&x = (2 - 0) + 0 \\
&\equiv \{ \text{simplification} \} \\
&x = 2
\end{aligned}$$

Clearly, this is a much more complicated way of finding solution $x = 2$ for equation (3.14). But we would have loved it in case it were the only known way, and this is precisely what happens with respect to (3.15). In this case we have:

$$\begin{aligned}
&L = 1 + A \times L \\
&\equiv \{ \text{substitution of } 1 + A \times L \text{ for } L \} \\
&L = 1 + A \times (1 + A \times L) \\
&\equiv \{ \text{distributive property (2.52)} \} \\
&L \cong 1 + A \times 1 + A \times (A \times L) \\
&\equiv \{ \text{unit of product (2.57) and associativity of product (2.34)} \} \\
&L \cong 1 + A + (A \times A) \times L \\
&\equiv \{ \text{by (2.96), (2.98) and (2.102)} \} \\
&L \cong A^0 + A^1 + A^2 \times L \\
&\equiv \{ \text{another substitution as above and similar simplifications} \} \\
&L \cong A^0 + A^1 + A^2 + A^3 \times L \\
&\equiv \{ \text{after } (n+1)\text{-many similar steps} \} \\
&L \cong \sum_{i=0}^n A^i + A^{n+1} \times L
\end{aligned}$$

Bearing a large n in mind, let us deliberately (but temporarily) ignore term $A^{n+1} \times L$. Then L will be isomorphic to the sum of n -many contributions A^i ,

$$L \cong \sum_{i=0}^n A^i$$

each of them consisting of i -long tuples, or *sequences*, of values of A . (Number i is said to be the *length* of any sequence in A^i .) Such sequences will be denoted by enumerating their elements between square brackets, for instance the *empty sequence* $[\]$ which is the only inhabitant in A^0 , the two element sequence $[a_1, a_2]$ which belongs to A^2 provided $a_1, a_2 \in A$, and so on. Note that all such contributions are mutually disjoint, that is, $A^i \cap A^j = \emptyset$ wherever $i \neq j$. (In other words, a sequence of length i is never a sequence of length j , for $i \neq j$.)

If we join all contributions A^i into a single set, we obtain the set of all *finite sequences* on A , denoted by A^* and defined as follows:

$$A^* \stackrel{\text{def}}{=} \bigcup_{i \geq 0} A^i \quad (3.16)$$

The intuition behind taking the limit in the numeric calculation above was that term $\frac{x}{2^{n+1}}$ was getting smaller and smaller as n went larger and larger and, “in the limit”, it could be ignored. By analogy, taking a similar limit in the calculation just sketched above will mean that, for a “sufficiently large” n , the sequences in A^n are so long that it is very unlikely that we will ever use them! So, for $n \rightarrow \infty$ we obtain

$$L \cong \sum_{i=0}^{\infty} A^i$$

Because $\sum_{i=0}^{\infty} A^i$ is isomorphic to $\bigcup_{i=0}^{\infty} A^i$ (see exercise 2.32), we finally have:

$$L \cong A^*$$

All in all, we have obtained A^* as a solution to equation (3.15). In other words, datatype L is isomorphic to the datatype which contains all finite sequences of some predefined datatype A . This corresponds to the HASKELL `[a]` datatype, in general. Recall that we started from the “linked-list datatype” expressed in PASCAL or C. In fact, wherever the C programmer thinks of linked-lists, the HASKELL programmer will think of finite sequences.

But, what does equation (3.15) mean in fact? Is A^* the only solution to this equation? Back to the numeric field, we know of equations which have more than one solution — for instance $x = \frac{x^2+3}{4}$, which admits two solutions 1 and 3 —, which have no solution at all — for instance $x = x + 1$ —, or which admit an infinite number of — for instance $x = x$.

We will address these topics in the next section about *inductive* datatypes and — more generally — in chapter 8, where the formal semantics of recursion will be made explicit. This is where the “limit” constructions used informally in this section will be shown to make sense.

3.3 INTRODUCING INDUCTIVE DATATYPES

Datatype L as defined by (3.15) is said to be *recursive* because L “re-curs” in the definition of L itself ⁶. From the discussion above, it is clear that set-theoretical equality “=” in this equation should give place to set-theoretical isomorphism (“ \cong ”):

$$L \cong 1 + A \times L \quad (3.17)$$

⁶ By analogy, we may regard (3.14) as a “recursive definition” of number 2.

Which isomorphism $L \xleftarrow{in} 1 + A \times L$ do we expect to witness (3.15)? This will depend on which particular solution to (3.15) we are thinking of. So far we have seen only one, A^* . By recalling the notion of *algebra* of a datatype (section 2.18), so we may rephrase the question as: which algebra

$$A^* \xleftarrow{in} 1 + A \times A^*$$

do we expect to witness the tautology which arises from (3.15) by replacing unknown L with solution A^* , that is

$$A^* \cong 1 + A \times A^* \quad ?$$

It will have to be of the form $in = [in_1, in_2]$ as depicted by the following diagram:

$$\begin{array}{ccc} 1 & \xrightarrow{i_1} 1 + A \times A^* \xleftarrow{i_2} & A \times A^* \\ & \searrow in_1 \quad \downarrow in \quad \swarrow in_2 & \\ & A^* & \end{array} \quad (3.18)$$

Arrows in_1 and in_2 can be guessed rather intuitively: $in_1 = []$, which will express the “NIL pointer” by the empty sequence, at A^* level, and $in_2 = cons$, where $cons$ is the standard “left append” sequence constructor, which we for the moment introduce rather informally as follows:

$$\begin{aligned} cons &: A \times A^* \rightarrow A^* \\ cons(a, [a_1, \dots, a_n]) &= [a, a_1, \dots, a_n] \end{aligned} \quad (3.19)$$

In a diagram:

$$\begin{array}{ccc} 1 & \xrightarrow{i_1} 1 + A \times A^* \xleftarrow{i_2} & A \times A^* \\ & \searrow [] \quad \downarrow [[], cons] \quad \swarrow cons & \\ & A^* & \end{array} \quad (3.20)$$

Of course, for in to be iso it needs to have an inverse, which is not hard to guess,

$$out \stackrel{\text{def}}{=} (! + \langle hd, tl \rangle) \cdot (=_{[]}?) \quad (3.21)$$

where sequence operators hd (*head of a nonempty sequence*) and tl (*tail of a nonempty sequence*) are (again informally) described as follows:

$$\begin{aligned} hd &: A^* \rightarrow A \\ hd [a_1, a_2, \dots, a_n] &= a_1 \end{aligned} \quad (3.22)$$

$$\begin{aligned} tl &: A^* \rightarrow A^* \\ tl [a_1, a_2, \dots, a_n] &= [a_2, \dots, a_n] \end{aligned} \quad (3.23)$$

Showing that *in* and *out* are each other inverses is not a hard task either:

$$\begin{aligned}
& in \cdot out = id \\
\equiv & \quad \{ \text{definitions of } in \text{ and } out \} \\
& [\underline{\quad}, cons] \cdot (! + \langle hd, tl \rangle) \cdot (=_{\underline{\quad}}?) = id \\
\equiv & \quad \{ +\text{-absorption (2.43) and (2.15)} \} \\
& [\underline{\quad}, cons \cdot \langle hd, tl \rangle] \cdot (=_{\underline{\quad}}?) = id \\
\equiv & \quad \{ \text{property of sequences: } cons(hd\ s, tl\ s) = s \} \\
& [\underline{\quad}, id] \cdot (=_{\underline{\quad}}?) = id \\
\equiv & \quad \{ \text{going pointwise (2.69)} \} \\
& \left\{ \begin{array}{l} =_{\underline{\quad}} a \Rightarrow [\underline{\quad}, id] (i_1 a) \\ \neg(=_{\underline{\quad}} a) \Rightarrow [\underline{\quad}, id] (i_2 a) \end{array} \right. = a \\
\equiv & \quad \{ +\text{-cancellation (2.40)} \} \\
& \left\{ \begin{array}{l} =_{\underline{\quad}} a \Rightarrow \underline{\quad} a \\ \neg(=_{\underline{\quad}} a) \Rightarrow id\ a \end{array} \right. = a \\
\equiv & \quad \{ a = \underline{\quad} \text{ in one case and identity function (2.9) in the other} \} \\
& \left\{ \begin{array}{l} a = \underline{\quad} \Rightarrow a \\ \neg(a = \underline{\quad}) \Rightarrow a \end{array} \right. = a \\
\equiv & \quad \{ \text{property } (p \rightarrow f, f) = f \text{ holds} \} \\
& a = a
\end{aligned}$$

A comment on the particular choice of terminology above: symbol *in* suggests that we are going inside, or constructing (synthesizing) values of A^* ; symbol *out* suggests that we are going out, or destructing (analyzing) values of A^* . We shall often resort to this duality in the sequel.

Are there more solutions to equation (3.17)? In trying to implement this equation, a HASKELL programmer could have written, after the declaration of type A , the following datatype declaration:

```
data L = Nil () | Cons (A, L)
```

which, as we have seen in section 2.18, can be written simply as

```
data L = Nil | Cons (A, L) (3.24)
```

and generates diagram

$$\begin{array}{ccc}
1 & \xrightarrow{i_1} & 1 + A \times L \xleftarrow{i_2} A \times L \\
& \searrow \underline{Nil} & \downarrow in' \swarrow Cons \\
& & L
\end{array} \quad (3.25)$$

leading to algebra $in' = [\underline{Nil}, Cons]$.

HASKELL seems to have generated another solution for the equation, which it calls L . To avoid the inevitable confusion between this symbol denoting the newly created datatype and symbol L in equation (3.17), which denotes a mathematical variable, let us use symbol T to denote the former (T stands for “type”). This can be coped with very simply by writing T instead of L in (3.24):

$$\text{data } T = Nil \mid Cons \ (A, T) \quad (3.26)$$

In order to make T more explicit, we will write in_T instead of in' .

Some questions are on demand at this point. First of all, what is datatype T ? What are its inhabitants? Next, is $T \xleftarrow{in_T} 1 + A \times T$ an iso or not?

HASKELL will help us to answer these questions. Suppose that A is a primitive numeric datatype, and that we add `deriving Show` to (3.26) so that we can “see” the inhabitants of the T datatype. The information associated to T is thus:

```
Main> :i T
-- type constructor
data T

-- constructors:
Nil :: T
Cons :: (A,T) -> T

-- instances:
instance Show T
instance Eval T
```

By typing `Nil`

```
Main> Nil
Nil :: T
```

we confirm that *Nil* is itself an inhabitant of T , and by typing `Cons`

```
Main> Cons
<<function>> :: (A,T) -> T
```

we realize that *Cons* is not so (as expected), but it can be used to build such inhabitants, for instance:

```
Main> Cons (1, Nil)
Cons (1, Nil) :: T
```

or

```
Main> Cons (2, Cons (1, Nil))
Cons (2, Cons (1, Nil)) :: T
```

etc. We conclude that *expressions* involving *Nil* and *Cons* are inhabitants of type T . Are these the *only* ones? The answer is *yes* because, by design of the HASKELL language, the constructors of type T will remain fixed once its declaration is interpreted, that is, no further constructor can be added to T . Does in_T have an inverse? Yes, its inverse is coalgebra

$$\begin{aligned} out_T &: T \rightarrow 1 + A \times T \\ out_T Nil &= i_1 NIL \\ out_T(Cons(a, l)) &= i_2(a, l) \end{aligned} \quad (3.27)$$

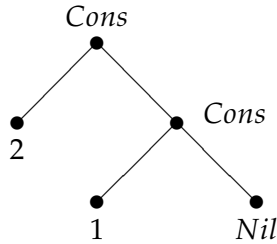
which can be straightforwardly encoded in HASKELL using the `Either` realization of $+$ (recall sections 2.9 and 2.18):

```
outT :: T -> Either () (A, T)
outT Nil = Left ()
outT (Cons(a, l)) = Right(a, l)
```

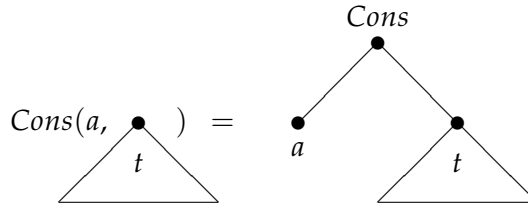
In summary, isomorphism

$$T \begin{array}{c} \xrightarrow{out_T} \\ \cong \\ \xleftarrow{in_T} \end{array} 1 + A \times T \quad (3.28)$$

holds, where datatype T is inhabited by symbolic expressions which we may visualize very conveniently as trees, for instance



picturing expression $Cons(2, Cons(1, Nil))$. *Nil* is the empty tree and *Cons* may be regarded as the operation which adds a new root and a new branch, say a , to a tree t :



The choice of symbols T , *Nil* and *Cons* was rather arbitrary in (3.26). Therefore, an alternative declaration such as, for instance,

$$\text{data } U = \text{Stop} \mid \text{Join } (A, U) \quad (3.29)$$

would have been perfectly acceptable, generating another solution for the equation under algebra $[Stop, Join]$. It is easy to check that (3.29) is but a renaming of Nil to $Stop$ and of $Cons$ to $Join$. Therefore, both datatypes are isomorphic, or “abstractly the same”.

Indeed, any other datatype X *inductively* defined by a constant and a binary constructor accepting A and X as parameters will be a solution to the equation. Because we are just renaming symbols in a consistent way, all such solutions are abstractly the same. All of them capture the abstract notion of a *list* of symbols.

We wrote “inductively” above because the set of all expressions (trees) which inhabit the type is defined by induction. Such types are called *inductive* and we shall have a lot more to say about them in chapter 8.

Exercise 3.6. Obviously,

```
either (const []) (:
```

does not work as a HASKELL realization of the mediating arrow in diagram (3.20). What do you need to write instead?

□

3.4 OBSERVING AN INDUCTIVE DATATYPE

Suppose that one is asked to express a particular *observation* of an inductive such as T (3.26), that is, a function of signature $B \xleftarrow{f} T$ for some target type B . Suppose, for instance, that A is \mathbb{N}_0 (the set of all non-negative integers) and that we want to add all elements which occur in a T -list. Of course, we have to ensure that addition is available in \mathbb{N}_0 ,

$$\begin{aligned} add : \mathbb{N}_0 \times \mathbb{N}_0 &\rightarrow \mathbb{N}_0 \\ add(x, y) &\stackrel{\text{def}}{=} x + y \end{aligned}$$

and that $0 \in \mathbb{N}_0$ is a value denoting “the addition of nothing”. So constant arrow $\mathbb{N}_0 \xleftarrow{0} 1$ is available. Of course, $add(0, x) = add(x, 0) = x$ holds, for all $x \in \mathbb{N}_0$. This property means that \mathbb{N}_0 , together with operator add and constant 0 , forms a *monoid*, a very important algebraic structure in computing which will be exploited intensively later in this book. The following arrow “packaging” \mathbb{N}_0 , add and 0 ,

$$\mathbb{N}_0 \xleftarrow{[0, add]} 1 + \mathbb{N}_0 \times \mathbb{N}_0 \quad (3.30)$$

is a convenient way to express such a structure. Combining this arrow with the algebra

$$T \xleftarrow{in_T} 1 + \mathbb{N}_0 \times T \quad (3.31)$$

which defines T , and the function f we want to define, the target of which is $B = \mathbb{N}_0$, we get the almost closed diagram which follows, in which only the dashed arrow is yet to be filled in:

$$\begin{array}{ccc} T & \xleftarrow{in_T} & 1 + \mathbb{N}_0 \times T \\ f \downarrow & & \vdots \downarrow \\ \mathbb{N}_0 & \xleftarrow{[0, add]} & 1 + \mathbb{N}_0 \times \mathbb{N}_0 \end{array} \quad (3.32)$$

We know that $in_T = [\underline{Nil}, Cons]$. A pattern for the missing arrow is not difficult to guess: in the same way f bridges T and \mathbb{N}_0 on the left-hand side, it will do the same job on the right-hand side. So pattern $\dots + \dots \times f$ comes to mind (recall section 2.10), where the “ \dots ” are very naturally filled in by identity functions. All in all, we obtain diagram

$$\begin{array}{ccc} T & \xleftarrow{[\underline{Nil}, Cons]} & 1 + \mathbb{N}_0 \times T \\ f \downarrow & & \downarrow id + id \times f \\ \mathbb{N}_0 & \xleftarrow{[0, add]} & 1 + \mathbb{N}_0 \times \mathbb{N}_0 \end{array} \quad (3.33)$$

which pictures the following property of f

$$f \cdot [\underline{Nil}, Cons] = [0, add] \cdot (id + id \times f) \quad (3.34)$$

and is easy to convert to pointwise notation:

$$\begin{aligned} & f \cdot [\underline{Nil}, Cons] = [0, add] \cdot (id + id \times f) \\ \equiv & \quad \{ \text{(2.42) on the lefthand side, (2.43) and identity } id \text{ on the righthand side} \} \\ & [f \cdot \underline{Nil}, f \cdot Cons] = [0, add \cdot (id \times f)] \\ \equiv & \quad \{ \text{either structural equality (2.66)} \} \\ & \begin{cases} f \cdot \underline{Nil} = 0 \\ f \cdot Cons = add \cdot (id \times f) \end{cases} \\ \equiv & \quad \{ \text{going pointwise} \} \\ & \begin{cases} (f \cdot \underline{Nil})x = 0x \\ (f \cdot Cons)(a, x) = (add \cdot (id \times f))(a, x) \end{cases} \\ \equiv & \quad \{ \text{composition (2.6), constant (2.12), product (2.24) and definition of } add \} \\ & \begin{cases} f \underline{Nil} = 0 \\ f(Cons(a, x)) = a + f x \end{cases} \end{aligned}$$

Note that we could have used out_T in diagram (3.32),

$$\begin{array}{ccc} T & \xrightarrow{out_T} & 1 + \mathbb{N}_0 \times T \\ f \downarrow & & \downarrow id + id \times f \\ \mathbb{N}_0 & \xleftarrow{[0, add]} & 1 + \mathbb{N}_0 \times \mathbb{N}_0 \end{array} \quad (3.35)$$

obtaining another version of the *definition* of f ,

$$f = [\underline{0}, add] \cdot (id + id \times f) \cdot out_T \quad (3.36)$$

which would lead to exactly the same pointwise recursive definition:

$$\begin{aligned} f &= [\underline{0}, add] \cdot (id + id \times f) \cdot out_T \\ \equiv & \quad \{ \text{(2.43) and identity } id \text{ on the righthand side} \} \\ f &= [\underline{0}, add \cdot (id \times f)] \cdot out_T \\ \equiv & \quad \{ \text{going pointwise on } out_T \text{ (3.27)} \} \\ & \begin{cases} f \text{ Nil} = ([\underline{0}, add \cdot (id \times f)] \cdot out_T) \text{ Nil} \\ f(\text{Cons}(a, x)) = ([\underline{0}, add \cdot (id \times f)] \cdot out_T)(a, x) \end{cases} \\ \equiv & \quad \{ \text{definition of } out_T \text{ (3.27)} \} \\ & \begin{cases} f \text{ Nil} = ([\underline{0}, add \cdot (id \times f)] \cdot i_1) \text{ Nil} \\ f(\text{Cons}(a, x)) = ([\underline{0}, add \cdot (id \times f)] \cdot i_2)(a, x) \end{cases} \\ \equiv & \quad \{ \text{+cancellation (2.40)} \} \\ & \begin{cases} f \text{ Nil} = \underline{0} \text{ Nil} \\ f(\text{Cons}(a, x)) = (add \cdot (id \times f))(a, x) \end{cases} \\ \equiv & \quad \{ \text{simplification} \} \\ & \begin{cases} f \text{ Nil} = 0 \\ f(\text{Cons}(a, x)) = a + f x \end{cases} \end{aligned}$$

Pointwise f mirrors the structure of type T in having as many definition clauses as constructors in T . Such functions are said to be defined *by induction on* the structure of their input type. If we repeat this calculation for \mathbb{N}_0^* instead of T , that is, for

$$out = (! + \langle hd, tl \rangle) \cdot (=_{[]}?)$$

— recall (3.21) — taking place of out_T , we get a “more algorithmic” version of f :

$$\begin{aligned} f &= [\underline{0}, add] \cdot (id + id \times f) \cdot (! + \langle hd, tl \rangle) \cdot (=_{[]}?) \\ \equiv & \quad \{ \text{+functor (2.44), identity and } \times \text{-absorption (2.27)} \} \\ f &= [\underline{0}, add] \cdot (! + \langle hd, f \cdot tl \rangle) \cdot (=_{[]}?) \\ \equiv & \quad \{ \text{+absorption (2.43) and constant } \underline{0} \} \\ f &= [\underline{0}, add \cdot \langle hd, f \cdot tl \rangle] \cdot (=_{[]}?) \\ \equiv & \quad \{ \text{going pointwise on guard } =_{[]}? \text{ (2.69) and simplifying} \} \\ f l &= \begin{cases} l = [] & \Rightarrow \underline{0} l \\ \neg(l = []) & \Rightarrow (add \cdot \langle hd, f \cdot tl \rangle) l \end{cases} \\ \equiv & \quad \{ \text{simplification} \} \\ f l &= \begin{cases} l = [] & \Rightarrow 0 \\ \neg(l = []) & \Rightarrow hd l + f(tl l) \end{cases} \end{aligned}$$

The outcome of this calculation can be encoded in HASKELL syntax as

$$\begin{aligned} f\ l \mid l \equiv [] &= 0 \\ \mid \text{otherwise} &= \text{head } l + f\ (\text{tail } l) \end{aligned}$$

or

$$f\ l = \text{if } l \equiv [] \text{ then } 0 \text{ else head } l + f\ (\text{tail } l)$$

both requiring the equality predicate \equiv and destructors `head` and `tail`.

3.5 SYNTHESIZING AN INDUCTIVE DATATYPE

The issue which concerns us in this section dualizes what we have just dealt with: instead of analyzing or *observing* an inductive type such as T (3.26), we want to be able to synthesize (generate) particular inhabitants of T . In other words, we want to be able to specify functions with signature $B \xrightarrow{f} T$ for some given source type B . Let $B = \mathbb{N}_0$ and suppose we want f to generate, for a given natural number $n > 0$, the list containing all numbers less or equal to n in decreasing order

$$\text{Cons}(n, \text{Cons}(n-1, \text{Cons}(\dots, \text{Nil})))$$

or the empty list Nil , in case $n = 0$.

Let us try and draw a diagram similar to (3.35) applicable to the new situation. In trying to “re-use” this diagram, it is immediate that arrow f should be reversed. Bearing duality in mind, we may feel tempted to reverse all arrows just to see what happens. Identity functions are their own inverses, and in_T takes the place of out_T :

$$\begin{array}{ccc} T & \xleftarrow{\text{in}_T} & 1 + \mathbb{N}_0 \times T \\ \uparrow f & & \uparrow \text{id} + \text{id} \times f \\ \mathbb{N}_0 & \xrightarrow{\quad\quad\quad} & 1 + \mathbb{N}_0 \times \mathbb{N}_0 \end{array}$$

Interestingly enough, the bottom arrow is the one which is not obvious to reverse, meaning that we have to “invent” a particular destructor of \mathbb{N}_0 , say

$$\mathbb{N}_0 \xrightarrow{g} 1 + \mathbb{N}_0 \times \mathbb{N}_0$$

fitting in the diagram and *generating* the particular computational effect we have in mind. Once we do this, a recursive definition for f will pop out immediately,

$$f = \text{in}_T \cdot (\text{id} + \text{id} \times f) \cdot g \quad (3.37)$$

which is equivalent to:

$$f = [\underline{\text{Nil}}, \text{Cons} \cdot (\text{id} \times f)] \cdot g \quad (3.38)$$

Because we want $f\ 0 = Nil$ to hold, g (the actual generator of the computation) should distinguish input 0 from all the others. One thus decomposes g as follows,

$$\mathbb{N}_0 \xrightarrow{=0?} \mathbb{N}_0 + \mathbb{N}_0 \xrightarrow{!+h} 1 + \mathbb{N}_0 \times \mathbb{N}_0$$

g

leaving h to fill in. This will be a *split* providing, on the lefthand side, for the value to be *Cons*'ed to the output and, on the righthand side, for the “seed” to the next recursive call. Since we want the output values to be produced contiguously and in decreasing order, we may define $h = \langle id, pred \rangle$ where, for $n > 0$,

$$pred\ n \stackrel{\text{def}}{=} n - 1 \quad (3.39)$$

computes the *predecessor* of n . Altogether, we have synthesized

$$g = (! + \langle id, pred \rangle) \cdot (=0?) \quad (3.40)$$

Filling this in (3.38) we get

$$\begin{aligned} f &= [\underline{Nil}, Cons \cdot (id \times f)] \cdot (! + \langle id, pred \rangle) \cdot (=0?) \\ &\equiv \{ \text{+ -absorption (2.43) followed by } \times \text{-absorption (2.27) etc.} \} \\ f &= [\underline{Nil}, Cons \cdot \langle id, f \cdot pred \rangle] \cdot (=0?) \\ &\equiv \{ \text{going pointwise on guard } =0? \text{ (2.69) and simplifying} \} \\ f\ n &= \begin{cases} n = 0 & \Rightarrow Nil \\ \neg(n = 0) & \Rightarrow Cons(n, f\ (n - 1)) \end{cases} \end{aligned}$$

which matches the function we had in mind:

$$\begin{aligned} f\ n & \\ &| \ n \equiv 0 = Nil \\ &| \text{otherwise} = Cons\ (n, f\ (n - 1)) \end{aligned}$$

We shall see briefly that the constructions of the f function adding up a list of numbers in the previous section and, in this section, of the f function generating a list of numbers are very standard in algorithm design and can be broadly generalized. Let us first introduce some standard terminology.

3.6 INTRODUCING (LIST) CATAS, ANAS AND HYLOS

Suppose that, back to section 3.4, we want to *multiply*, rather than add, the elements occurring in lists of type T (3.26). How much of the program synthesis effort presented there can be reused in the design of the new function?

It is intuitive that only the bottom arrow $\mathbb{N}_0 \xleftarrow{[0, add]} 1 + \mathbb{N}_0 \times \mathbb{N}_0$ of diagram (3.35) needs to be replaced, because this is the only place

where we can specify that target datatype \mathbb{N}_0 is now regarded as the carrier of another (multiplicative rather than additive) monoidal structure,

$$\mathbb{N}_0 \xleftarrow{[1, mul]} 1 + \mathbb{N}_0 \times \mathbb{N}_0 \quad (3.41)$$

for $mul(x, y) \stackrel{\text{def}}{=} x \cdot y$. We are saying that the argument list is now to be reduced by the multiplication operator and that output value 1 is expected as the result of “nothing left to multiply”.

Moreover, in the previous section we might have wanted our number-list generator to produce the list of even numbers smaller than a given number, in decreasing order (see exercise 3.9). Intuition will once again help us in deciding that only arrow g in (3.37) needs to be updated.

The following diagrams generalize both constructions by leaving such bottom arrows unspecified,

$$\begin{array}{ccc} T & \xrightarrow{out_T} & 1 + \mathbb{N}_0 \times T \\ f \downarrow & & \downarrow id + id \times f \\ B & \xleftarrow{g} & 1 + \mathbb{N}_0 \times B \end{array} \quad \begin{array}{ccc} T & \xleftarrow{in_T} & 1 + \mathbb{N}_0 \times T \\ f \uparrow & & \uparrow id + id \times f \\ B & \xrightarrow{g} & 1 + \mathbb{N}_0 \times B \end{array} \quad (3.42)$$

and express their duality (cf. the directions of the arrows). It so happens that, for each of these diagrams, f is uniquely dependent on the g arrow, that is to say, each particular instantiation of g will determine the corresponding f . So both g s can be regarded as “seeds” or “genetic material” of the f functions they uniquely define⁷.

Following the standard terminology, we express these facts by writing $f = \langle g \rangle$ with respect to the lefthand side diagram and by writing $f = \llbracket g \rrbracket$ with respect to the righthand side diagram. Read $\langle g \rangle$ as “the T -catamorphism induced by g ” and $\llbracket g \rrbracket$ as “the T -anamorphism induced by g ”. This terminology is derived from the Greek words $\kappa\alpha\tau\alpha$ (cata) and $\alpha\nu\alpha$ (ana) meaning, respectively, “downwards” and “upwards” (compare with the direction of the f arrow in each diagram). The exchange of parentheses “()” and “[]” in double parentheses “ $\langle \rangle$ ” and “ $\llbracket \rrbracket$ ” is aimed at expressing the duality of both concepts.

We shall have a lot to say about catamorphisms and anamorphisms of a given type such as T . For the moment, it suffices to say that

- the T -catamorphism induced by $B \xleftarrow{g} 1 + \mathbb{N}_0 \times B$ is the unique function $B \xleftarrow{\langle g \rangle} T$ which obeys to property (or is defined by)

$$\langle g \rangle = g \cdot (id + id \times \langle g \rangle) \cdot out_T \quad (3.43)$$

which is the same as

$$\langle g \rangle \cdot in_T = g \cdot (id + id \times \langle g \rangle) \quad (3.44)$$

⁷ The theory which supports the statements of this paragraph will not be dealt with until chapter 8.

- given $B \xrightarrow{g} 1 + \mathbb{N}_0 \times B$ the T-anamorphism induced by g is the unique function $B \xrightarrow{[g]} T$ which obeys to property (or is defined by)

$$[g] = in_T \cdot (id + id \times [g]) \cdot g \quad (3.45)$$

From (3.42) it can be observed that T can act as a mediator between any T -anamorphism and any T -catamorphism, that is to say, $B \xleftarrow{[g]} T$ composes with $T \xleftarrow{[h]} C$, for some $C \xrightarrow{h} 1 + \mathbb{N}_0 \times C$. In other words, a T -catamorphism call always observe (consume) the output of a T -anamorphism. The latter produces a list of \mathbb{N}_0 s which is consumed by the former. This is depicted in the diagram which follows:

$$\begin{array}{ccc} B & \xleftarrow{g} & 1 + \mathbb{N}_0 \times B \\ \uparrow [g] & & \uparrow id + id \times [g] \\ T & \xleftarrow{in_T} & 1 + \mathbb{N}_0 \times T \\ \uparrow [h] & & \uparrow id + id \times [h] \\ C & \xrightarrow{h} & 1 + \mathbb{N}_0 \times C \end{array} \quad (3.46)$$

What can we say about the $[g] \cdot [h]$ composition? It is a function from C to B which resorts to T as an *intermediate* data-structure and can be subject to the following calculation (cf. outermost rectangle in (3.46)):

$$\begin{aligned} [g] \cdot [h] &= g \cdot (id + id \times [g]) \cdot (id + id \times [h]) \cdot h \\ &\equiv \{ \text{+functor (2.44)} \} \\ [g] \cdot [h] &= g \cdot ((id \cdot id) + (id \times [g]) \cdot (id \times [h])) \cdot h \\ &\equiv \{ \text{identity and } \times\text{-functor (2.30)} \} \\ [g] \cdot [h] &= g \cdot (id + id \times [g] \cdot [h]) \cdot h \end{aligned}$$

This calculation shows how to define $C \xleftarrow{[g] \cdot [h]} B$ in one go, that is to say, doing without any intermediate data-structure:

$$\begin{array}{ccc} B & \xleftarrow{g} & 1 + \mathbb{N}_0 \times B \\ \uparrow [g] \cdot [h] & & \uparrow id + id \times [g] \cdot [h] \\ C & \xrightarrow{h} & 1 + \mathbb{N}_0 \times C \end{array} \quad (3.47)$$

As an example, let us see what comes out of $[g] \cdot [h]$ for h and g respectively given by (3.40) and (3.41):

$$\begin{aligned} [g] \cdot [h] &= g \cdot (id + id \times [g]) \cdot [h] \cdot h \\ &\equiv \{ [g] \cdot [h] \text{ abbreviated to } f \text{ and instantiating } h \text{ and } g \} \end{aligned}$$

$$\begin{aligned}
f &= [\underline{1}, mul] \cdot (id + id \times f) \cdot (! + \langle id, pred \rangle) \cdot (=_0?) \\
&\equiv \{ \text{+-functor (2.44) and identity} \} \\
f &= [\underline{1}, mul] \cdot (! + (id \times f) \cdot \langle id, pred \rangle) \cdot (=_0?) \\
&\equiv \{ \text{\(\times\)-absorption (2.27) and identity} \} \\
f &= [\underline{1}, mul] \cdot (! + \langle id, f \cdot pred \rangle) \cdot (=_0?) \\
&\equiv \{ \text{+-absorption (2.43) and constant } \underline{1} \text{ (2.15)} \} \\
f &= [\underline{1}, mul \cdot \langle id, f \cdot pred \rangle] \cdot (=_0?) \\
&\equiv \{ \text{McCarthy conditional (2.70)} \} \\
f &= (=_0?) \rightarrow \underline{1}, mul \cdot \langle id, f \cdot pred \rangle
\end{aligned}$$

Going pointwise, we get — via (2.70) —

$$\begin{aligned}
f\ 0 &= [\underline{1}, mul \cdot \langle id, f \cdot pred \rangle](i_1\ 0) \\
&= \{ \text{+-cancellation (2.40)} \} \\
&\quad \underline{1}\ 0 \\
&= \{ \text{constant function (2.12)} \} \\
&\quad 1
\end{aligned}$$

and

$$\begin{aligned}
f(n+1) &= [\underline{1}, mul \cdot \langle id, f \cdot pred \rangle](i_2(n+1)) \\
&= \{ \text{+-cancellation (2.40)} \} \\
&\quad mul \cdot \langle id, f \cdot pred \rangle(n+1) \\
&= \{ \text{pointwise definitions of } split, \text{ identity, predecessor and } mul \} \\
&\quad (n+1) \times f\ n
\end{aligned}$$

In summary, f is but the well-known factorial function:

$$\begin{cases} f\ 0 = 1 \\ f(n+1) = (n+1) \times f\ n \end{cases}$$

This result comes to no surprise if we look at diagram (3.46) for the particular g and h we have considered above and recall a popular “definition” of factorial:

$$n! = \underbrace{n \times (n-1) \times \dots \times 1}_{n \text{ times}} \tag{3.48}$$

In fact, $[(h)]\ n$ produces T-list

$$Cons(n, Cons(n-1, \dots Cons(1, Nil)))$$

as an intermediate data-structure which is consumed by $(\llbracket g \rrbracket)$, the effect of which is but the “replacement” of $Cons$ by \times and Nil by 1 , therefore accomplishing (3.48) and realizing the computation of factorial.

The moral of this example is that a function as simple as factorial can be *decomposed* into two components (producer/consumer functions) which share a common intermediate inductive datatype. The producer function is an anamorphism which “represents” or produces a “view” of its input argument as a value of the intermediate datatype. The consumer function is a catamorphism which reduces this intermediate data-structure and produces the final result. Like factorial, many functions can be handsomely expressed by a $\llbracket g \rrbracket \cdot \llbracket h \rrbracket$ composition for a suitable choice of the intermediate type, and of g and h .

The intermediate data-structure is said to be *virtual* in the sense that it only exists as a means to induce the associated pattern of recursion and disappears by calculation. The composition

$$\llbracket g \rrbracket \cdot \llbracket h \rrbracket$$

of a T-catamorphism with a T-anamorphism is called a T-hylomorphism⁸ and is denoted by $\llbracket g, h \rrbracket$. Because g and h fully determine the behaviour of the $\llbracket g, h \rrbracket$ function, they can be regarded as the “genes” of the function they define. As we shall see, this analogy with biology will prove specially useful for algorithm analysis and classification.

Exercise 3.7. A way of computing n^2 , the square of a given natural number n , is to sum up the n first odd numbers. In fact, $1^2 = 1$, $2^2 = 1 + 3$, $3^2 = 1 + 3 + 5$, etc., $n^2 = (2n - 1) + (n - 1)^2$. Following this hint, express function

$$sq\ n \stackrel{\text{def}}{=} n^2 \tag{3.49}$$

as a T-hylomorphism and encode it in HASKELL.

□

Exercise 3.8. Write function x^n as a T-hylomorphism and encode it in HASKELL.

□

Exercise 3.9. The following function in HASKELL computes the T-sequence of all even numbers less or equal to n :

$f\ n = \text{if } n \leq 1 \text{ then Nil else Cons } (m, f\ (m - 2))$
where $m = \text{if even } n \text{ then } n \text{ else } n - 1$

Find its “genetic material”, that is, function g such that $f = \llbracket g \rrbracket$ in

$$\begin{array}{ccc} T & \xleftarrow{\text{in}_T} & 1 + \mathbb{N}_0 \times T \\ \uparrow \llbracket g \rrbracket & & \uparrow \text{id} + \text{id} \times \llbracket g \rrbracket \\ \mathbb{N}_0 & \xrightarrow{g} & 1 + \mathbb{N}_0 \times \mathbb{N}_0 \end{array}$$

8 This terminology is derived from the Greek word *υλος* (hylos) meaning “matter”.

□

3.7 INDUCTIVE TYPES MORE GENERALLY

So far we have focussed our attention exclusively to a particular inductive type T (3.31) — that of finite sequences of non-negative integers. This is, of course, of a very limited scope. First, because one could think of finite sequences of other datatypes, *e.g.* Booleans or many others. Second, because other datatypes such as trees, hash-tables *etc.* exist which our notation and method should be able to take into account.

Although a generic theory of arbitrary datatypes requires a theoretical elaboration which cannot be explained at once, we can move a step further by taking the two observations above as starting points. We shall start from the latter in order to talk generically about inductive types. Then we introduce parameterization and functorial behaviour.

Suppose that, as a mere notational convention, we abbreviate every expression of the form “ $1 + \mathbb{N}_0 \times \dots$ ” occurring in the previous section by “ $F \dots$ ”, *e.g.* $1 + \mathbb{N}_0 \times B$ by $F B$, *e.g.* $1 + \mathbb{N}_0 \times T$ by $F T$

$$\begin{array}{ccc} & \xrightarrow{\text{out}_T} & \\ T & \cong & F T \\ & \xleftarrow{\text{in}_T} & \end{array} \quad (3.50)$$

etc. This is the same as introducing a datatype-level operator

$$F X \stackrel{\text{def}}{=} 1 + \mathbb{N}_0 \times X \quad (3.51)$$

which maps every datatype A into datatype $1 + \mathbb{N}_0 \times A$. Operator F captures the pattern of recursion which is associated to so-called “right” lists (of non-negative integers), that is, lists which grow to the right. The slightly different pattern $G X \stackrel{\text{def}}{=} 1 + X \times \mathbb{N}_0$ will generate a different, although related, inductive type

$$X \cong 1 + X \times \mathbb{N}_0 \quad (3.52)$$

— that of so-called “left” lists (of non-negative integers). And it is not difficult to think of the pattern which merges both right and left lists and gives rise to bi-linear lists, better known as *binary trees*:

$$X \cong 1 + X \times \mathbb{N}_0 \times X \quad (3.53)$$

One may think of many other expressions $F X$ and guess the inductive datatype they generate, for instance $H X \stackrel{\text{def}}{=} \mathbb{N}_0 + \mathbb{N}_0 \times X$ generating

non-empty lists of non-negative integers (\mathbb{N}_0^+). The general rule is that, given an inductive datatype definition of the form

$$X \cong F X \quad (3.54)$$

(also called a domain equation), its pattern of recursion is captured by a so-called *functor* F .

3.8 FUNCTORS

The concept of a functor F , borrowed from category theory, is a most generic and useful device in programming⁹. As we have seen, F can be regarded as a datatype constructor which, given datatype A , builds a more elaborate datatype $F A$; given another datatype B , builds a similarly elaborate datatype $F B$; and so on. But what is more important and has the most beneficial consequences is that, if F is regarded as a functor, then its data-structuring effect extends smoothly to functions

in the following way: suppose that $B \xleftarrow{f} A$ is a function which observes A into B , which are parameters of $F A$ and $F B$, respectively. By definition, if F is a functor then $F B \xleftarrow{Ff} F A$ exists for every such f :

$$\begin{array}{ccc} A & \xrightarrow{\quad} & F A \\ f \downarrow & & \downarrow Ff \\ B & \xrightarrow{\quad} & F B \end{array}$$

Ff extends f to F -structures and will, by definition, obey to two very basic properties: it commutes with identity

$$F id_A = id_{(F A)} \quad (3.55)$$

and with composition

$$F(g \cdot h) = (Fg) \cdot (Fh) \quad (3.56)$$

Two simple examples of a functor follow:

- Identity functor: define $F X = X$, for every datatype X , and $F f = f$. Properties (3.55) and (3.56) hold trivially just by removing symbol F wherever it occurs.
- Constant functors: for a given C , define $F X = C$ (for all datatypes X) and $F f = id_C$, as expressed in the following diagram:

$$\begin{array}{ccc} A & \xrightarrow{\quad} & C \\ f \downarrow & & \downarrow id_C \\ B & \xrightarrow{\quad} & C \end{array}$$

Properties (3.55) and (3.56) hold trivially again.

⁹ The category theory practitioner must be warned of the fact that the word *functor* is used here in a too restrictive way. A proper (generic) definition of a functor will be provided later in this book.

Data construction	Universal construct	Functor	Description
$A \times B$	$\langle f, g \rangle$	$f \times g$	Product
$A + B$	$[f, g]$	$f + g$	Coproduct
B^A	\overline{f}	f^A	Exponential

Table 2.: Datatype constructions and associated operators.

In the same way functions can be unary, binary, *etc.*, we can have functors with more than one argument. So we get binary functors (also called *bifunctors*), ternary functors *etc.*. Of course, properties (3.55) and (3.56) have to hold for every parameter of an n -ary functor. For a binary functor B , for instance, equation (3.55) becomes

$$B(id_A, id_B) = id_{B(A,B)} \quad (3.57)$$

and equation (3.56) becomes

$$B(g \cdot h, i \cdot j) = B(g, i) \cdot B(h, j) \quad (3.58)$$

Product and coproduct are typical examples of bifunctors. In the former case one has $B(A, B) = A \times B$ and $B(f, g) = f \times g$ — recall (2.24). Properties (2.31) and (2.30) instantiate (3.57) and (3.58), respectively, and this explains why we called them the functorial properties of product. In the latter case, one has $B(A, B) = A + B$ and $B(f, g) = f + g$ — recall (2.39) — and functorial properties (2.45) and (2.44). Finally, exponentiation is a functorial construction too: assuming A , one has $F X \stackrel{\text{def}}{=} X^A$ and $F f \stackrel{\text{def}}{=} \overline{f \cdot ap}$ and functorial properties (2.89) and (2.90). All this is summarized in table 2.

Such as functions, functors may compose with each other in the obvious way: the composition of F and G , denoted $F \cdot G$, is defined by

$$(F \cdot G)X \stackrel{\text{def}}{=} F(G X) \quad (3.59)$$

$$(F \cdot G)f \stackrel{\text{def}}{=} F(G f) \quad (3.60)$$

3.9 POLYNOMIAL FUNCTORS

We may put constant, product, coproduct and identity functors together to obtain so-called *polynomial functors*, which are described by polynomial expressions, for instance

$$F X = 1 + A \times X$$

— recall (3.17). A polynomial functor is either

- a constant functor or the identity functor, or
- the (finitary) product or coproduct (sum) of other polynomial functors, or

- the composition of other polynomial functors.

So the effect on arrows of a polynomial functor is computed in an easy and structured way, for instance:

$$\begin{aligned}
 Ff &= (1 + A \times X)f \\
 &= \{ \text{sum of two functors where } A \text{ is a constant and } X \text{ is a variable} \} \\
 &\quad (1)f + (A \times X)f \\
 &= \{ \text{constant functor and product of two functors} \} \\
 &\quad id_1 + (A)f \times (X)f \\
 &= \{ \text{constant functor and identity functor} \} \\
 &\quad id_1 + id_A \times f \\
 &= \{ \text{subscripts dropped for simplicity} \} \\
 &\quad id + id \times f
 \end{aligned}$$

So, $1 + A \times f$ denotes the same as $id_1 + id_A \times f$, or even the same as $id + id \times f$ if one drops the subscripts.

It should be clear at this point that what was referred to in section 2.10 as a “symbolic pattern” applicable to both datatypes and arrows is after all a functor in the mathematical sense. The fact that the same polynomial expression is used to denote both the data and the operators which structurally transform such data is of great conceptual economy and practical application. For instance, once polynomial functor (3.51) is assumed, the diagrams in (3.42) can be written as simply as

$$\begin{array}{ccc}
 T & \xrightarrow{out_T} & FT \\
 f \downarrow & & \downarrow Ff \\
 B & \xleftarrow{g} & FB
 \end{array}
 \quad
 \begin{array}{ccc}
 T & \xleftarrow{in_T} & FT \\
 f \uparrow & & \uparrow Ff \\
 B & \xrightarrow{g} & FB
 \end{array}
 \quad (3.61)$$

It is useful to know that, thanks to the isomorphism laws studied in chapter 2, every polynomial functor F may be put into the canonical form,

$$\begin{aligned}
 FX &\cong C_0 + (C_1 \times X) + (C_2 \times X^2) + \cdots + (C_n \times X^n) \\
 &= \sum_{i=0}^n C_i \times X^i
 \end{aligned} \quad (3.62)$$

and that *Newton’s binomial formula*

$$(A + B)^n \cong \sum_{p=0}^n {}^nC_p \times A^{n-p} \times B^p \quad (3.63)$$

can be used in such conversions. These are performed up to isomorphism, that is to say, after the conversion one gets a different but isomorphic datatype. Consider, for instance, functor

$$FX \stackrel{\text{def}}{=} A \times (1 + X)^2$$

(where A is a constant datatype) and check the following reasoning:

$$\begin{aligned}
FX &= A \times (1 + X)^2 \\
&\cong \{ \text{law (2.102)} \} \\
&\quad A \times ((1 + X) \times (1 + X)) \\
&\cong \{ \text{law (2.52)} \} \\
&\quad A \times ((1 + X) \times 1 + (1 + X) \times X) \\
&\cong \{ \text{laws (2.57), (2.33) and (2.52)} \} \\
&\quad A \times ((1 + X) + (1 \times X + X \times X)) \\
&\cong \{ \text{laws (2.57) and (2.102)} \} \\
&\quad A \times ((1 + X) + (X + X^2)) \\
&\cong \{ \text{law (2.48)} \} \\
&\quad A \times (1 + (X + X) + X^2) \\
&\cong \{ \text{canonical form obtained via laws (2.52) and (2.103)} \} \\
&\quad \underbrace{A}_{C_0} + \underbrace{A \times 2 \times X}_{C_1} + \underbrace{A}_{C_2} \times X^2
\end{aligned}$$

Exercise 3.10. Synthesize the isomorphism

$$A + A \times 2 \times X + A \times X^2 \xleftarrow{\nu} A \times (1 + X^2)$$

implicit in the above reasoning.

□

3.10 POLYNOMIAL INDUCTIVE TYPES

An inductive datatype is said to be *polynomial* wherever its pattern of recursion is described by a polynomial functor, that is to say, wherever F in equation (3.54) is polynomial. For instance, datatype T (3.31) is polynomial ($n = 1$) and its associated polynomial functor is canonically defined with coefficients $C_0 = 1$ and $C_1 = \mathbb{N}_0$. For reasons that will become apparent later on, we shall always impose $C_0 \neq 0$ to hold in a *polynomial* datatype expressed in canonical form.

Polynomial types are easy to encode in HASKELL wherever the associated functor is in canonical polynomial form, that is, wherever one has

$$T \xleftarrow[\text{in}_T]{\cong} \sum_{i=0}^n C_i \times T^i \quad (3.64)$$

Then we have

$$\text{in}_T \stackrel{\text{def}}{=} [f_1, \dots, f_n]$$

where, for $i = 1, n$, f_i is an arrow of type $T \leftarrow C_i \times T^i$. Since n is finite, one may expand exponentials according to (2.102) and encode this in HASKELL as follows:

data $T = C0 \mid C1 (C1, T) \mid C2 (C2, (T, T)) \mid \dots \mid Cn (Cn, (T, \dots, T))$

Of course the choice of symbol C_i to realize each f_i is arbitrary¹⁰. Several instances of polynomial inductive types (in canonical form) will be mentioned in section 3.14. Section 3.19 will address the conversion between inductive datatypes induced by so-called *natural transformations*.

The concepts of catamorphism, anamorphism and hylomorphism introduced in section 3.6 can be extended to arbitrary polynomial types. We devote the following sections to explaining catamorphisms in the polynomial setting. Polynomial anamorphisms and hylomorphisms will not be dealt with until chapter 8.

3.11 F-ALGEBRAS AND F-HOMOMORPHISMS

Our interest in polynomial types is basically due to the fact that, for polynomial F , equation (3.54) always has a particularly interesting solution which corresponds to our notion of a recursive datatype.

In order to explain this, we need two notions which are easy to understand: first, that of an *F-algebra*, which simply is any function α of signature $A \xleftarrow{\alpha} F A$. A is called the *carrier* of F -algebra α and contains the values which α manipulates by computing new A -values out of existing ones, according to the F -pattern (the “type” of the algebra). As examples, consider $[0, add]$ (3.30) and in_T (3.31), which are both algebras of type $F X = 1 + \mathbb{N}_0 \times X$. The type of an algebra clearly determines its form. For instance, any algebra α of type $F X = 1 + X \times X$ will be of form $[\alpha_1, \alpha_2]$, where α_1 is a constant and α_2 is a binary operator. So monoids are algebras of this type¹¹.

Secondly, we introduce the notion of an *F-homomorphism* which is but a function observing a particular F -algebra α into another F -algebra β :

$$\begin{array}{ccc} A & \xleftarrow{\alpha} & F A \\ f \downarrow & & \downarrow F f \\ B & \xleftarrow{\beta} & F B \end{array} \quad f \cdot \alpha = \beta \cdot (F f) \quad (3.65)$$

¹⁰ A more traditional (but less close to (3.64)) encoding will be

data $T = C0 \mid C1 C1 T \mid C2 C2 T T \mid \dots \mid Cn Cn T \dots T$

delivering every constructor in curried form.

¹¹ But not every algebra of this type is a monoid, since the type of an algebra only fixes its syntax and does not impose any properties such as associativity, *etc.*

Clearly, f can be regarded as a structural translation between A and B , that is, A and B have a similar structure¹². Note that — thanks to (3.55) — identity functions are always (trivial) F -homomorphisms and that — thanks to (3.56) — these homomorphisms compose, that is, the composition of two F -homomorphisms is an F -homomorphism.

3.12 F-CATAMORPHISMS

An F -algebra can be epic, monic or both, that is, iso. Iso F -algebras are particularly relevant to our discussion because they describe solutions to the $X \cong F X$ equation (3.54). Moreover, for polynomial F a particular iso F -algebra always exists, which is denoted by $\mu F \xleftarrow{in} F \mu F$ and has special properties. First, its carrier is the smallest among the carriers of other iso F -algebras, and this is why it is denoted by μF — μ for “minimal”¹³. Second, it is the so-called *initial* F -algebra. What does this mean?

It means that, for every F -algebra α there exists one and only one F -homomorphism between in and α . This unique arrow mediating in and α is therefore determined by α itself, and is called the *F -catamorphism* generated by α . This construct, which was introduced in 3.6, is in general denoted by $\langle \alpha \rangle_F$:

$$\begin{array}{ccc} \mu F & \xleftarrow{in} & F \mu F \\ f = \langle \alpha \rangle_F \downarrow & & \downarrow F \langle \alpha \rangle_F \\ A & \xleftarrow{\alpha} & F A \end{array} \quad (3.66)$$

We will drop the F subscript in $\langle \alpha \rangle_F$ wherever deducible from the context, and often call α the “gene” of $\langle \alpha \rangle_F$.

As happens with *splits*, *ethers* and transposes, the uniqueness of the catamorphism construct is captured by a universal property established in the class of all F -homomorphisms:

$$k = \langle \alpha \rangle \Leftrightarrow k \cdot in = \alpha \cdot F k \quad (3.67)$$

According to the experience gathered from section 2.13 onwards, a few properties can be expected as consequences of (3.67). For instance, one may wonder about the “gene” of the identity catamorphism. Just let $k = id$ in (3.67) and see what happens:

$$\begin{aligned} id &= \langle \alpha \rangle \Leftrightarrow id \cdot in = \alpha \cdot F id \\ &= \{ \text{identity (2.10) and } F \text{ is a functor (3.55)} \} \\ id &= \langle \alpha \rangle \Leftrightarrow in = \alpha \cdot id \end{aligned}$$

¹² Cf. *homomorphism* = *homo* (the same) + *morphos* (structure, shape).

¹³ μF means the least fixpoint solution of equation $X \cong F X$, as will be described in chapter 8.

$$\begin{aligned}
&= \{ \text{identity (2.10) once again} \} \\
&\quad id = (\lceil \alpha \rceil) \Leftrightarrow in = \alpha \\
&= \{ \alpha \text{ replaced by } in \text{ and simplifying} \} \\
&\quad id = (\lceil in \rceil)
\end{aligned}$$

Thus one finds out that the genetic material of the identity catamorphism is the initial algebra in . Which is the same as establishing the *reflection property* of catamorphisms:

Cata-reflection :

$$\begin{array}{ccc}
\mu F & \xleftarrow{in} & F \mu F \\
(\lceil in \rceil) \downarrow & & \downarrow F(\lceil in \rceil) \\
\mu F & \xleftarrow{in} & F \mu F
\end{array} \quad (\lceil in \rceil) = id_{\mu F} \quad (3.68)$$

In a more intuitive way, one might have observed that $(\lceil in \rceil)$ is, by definition of in , the unique arrow mediating μF and itself. But another arrow of the same type is already known: the identity $id_{\mu F}$. So these two arrows must be the same.

Another property following immediately from (3.67), for $k = (\lceil \alpha \rceil)$, is

Cata-cancellation :

$$(\lceil \alpha \rceil) \cdot in = \alpha \cdot F(\lceil \alpha \rceil) \quad (3.69)$$

Because in is iso, this law can be rephrased as follows

$$(\lceil \alpha \rceil) = \alpha \cdot F(\lceil \alpha \rceil) \cdot out \quad (3.70)$$

where out denotes the inverse of in :

$$\begin{array}{ccc}
& out & \\
\mu F & \xrightarrow{\quad} & F \mu F \\
& \cong & \\
& in &
\end{array}$$

Now, let f be F -homomorphism (3.65) between F -algebras α and β . How does it relate to $(\lceil \alpha \rceil)$ and $(\lceil \beta \rceil)$? Note that $f \cdot (\lceil \alpha \rceil)$ is an arrow mediating μF and B . But B is the carrier of β and $(\lceil \beta \rceil)$ is the unique arrow mediating μF and B . So the two arrows are the same:

Cata-fusion :

$$\begin{array}{ccc}
\mu F & \xleftarrow{in} & F \mu F \\
(\lceil \alpha \rceil) \downarrow & & \downarrow F(\lceil \alpha \rceil) \\
A & \xleftarrow{\alpha} & F A \\
f \downarrow & & \downarrow F f \\
B & \xleftarrow{\beta} & F B
\end{array} \quad f \cdot (\lceil \alpha \rceil) = (\lceil \beta \rceil) \quad \text{if} \quad f \cdot \alpha = \beta \cdot F f \quad (3.71)$$

Of course, this law is also a consequence of the universal property, for $k = f \cdot \langle \alpha \rangle$:

$$\begin{aligned}
 f \cdot \langle \alpha \rangle = \langle \beta \rangle &\Leftrightarrow (f \cdot \langle \alpha \rangle) \cdot \text{in} = \beta \cdot F(f \cdot \langle \alpha \rangle) \\
 &\Leftrightarrow \{ \text{composition is associative and } F \text{ is a functor (3.56)} \} \\
 &\quad f \cdot (\langle \alpha \rangle \cdot \text{in}) = \beta \cdot (F f) \cdot (F \langle \alpha \rangle) \\
 &\Leftrightarrow \{ \text{cata-cancellation (3.69)} \} \\
 &\quad f \cdot \alpha \cdot F \langle \alpha \rangle = \beta \cdot F f \cdot F \langle \alpha \rangle \\
 &\Leftrightarrow \{ \text{require } f \text{ to be a } F\text{-homomorphism (3.65)} \} \\
 &\quad f \cdot \alpha \cdot F \langle \alpha \rangle = f \cdot \alpha \cdot F \langle \alpha \rangle \wedge f \cdot \alpha = \beta \cdot F f \\
 &\Leftrightarrow \{ \text{simplify} \} \\
 &\quad f \cdot \alpha = \beta \cdot F f
 \end{aligned}$$

The presentation of the *absorption* property of catamorphisms entails the very important issue of parameterization and deserves to be treated in a separate section, as follows.

3.13 PARAMETERIZATION AND TYPE FUNCTORS

By analogy with what we have done about *splits* (product), *eithers* (co-product) and transposes (exponential), we now look forward to identifying F -catamorphisms which exhibit functorial behaviour.

Suppose that one wishes to square all numbers that are members of lists of type T (3.31). It can be checked that

$$(\langle \text{Nil}, \text{Cons} \cdot (sq \times id) \rangle) \quad (3.72)$$

will do this for us, where $\mathbb{N}_0 \xleftarrow{sq} \mathbb{N}_0$ is given by (3.49). This catamorphism, which converted to pointwise notation is nothing but function h which follows

$$\begin{cases} h \text{ Nil} = \text{Nil} \\ h(\text{Cons}(a, l)) = \text{Cons}(sq a, h l) \end{cases}$$

maps type T to itself. This is because sq maps \mathbb{N}_0 to \mathbb{N}_0 . Now suppose that, instead of sq , one would like to apply a given function

$B \xleftarrow{f} \mathbb{N}_0$ (for some B other than \mathbb{N}_0) to all elements of the argument list. It is easy to see that it suffices to replace f for sq in (3.72). However, the output type no longer is T , but rather type $T' \cong 1 + B \times T$.

Types T and T' are very close to each other. They share the same “shape” (recursive pattern) and only differ with respect to the type of elements — \mathbb{N}_0 in T and B in T' . This suggests that these two types can be regarded as instances of a more generic list datatype List

$$\begin{array}{ccc}
 \text{List } X & \cong & 1 + X \times \text{List } X \\
 & \nwarrow \text{in} = [\text{Nil}, \text{Cons}] & \\
 & &
 \end{array} \quad (3.73)$$

in which the type of elements X is allowed to vary. Thus one has $T = \text{List } \mathbb{N}_0$ and $T' = \text{List } B$.

By inspection, it can be checked that, for every $B \xleftarrow{f} A$,

$$(|[\underline{Nil}, \text{Cons} \cdot (f \times id)]|) \quad (3.74)$$

maps $\text{List } A$ to $\text{List } B$. Moreover, for $f = id$ one has:

$$\begin{aligned} & (|[\underline{Nil}, \text{Cons} \cdot (id \times id)]|) \\ = & \quad \{ \text{by the } \times\text{-functor-id property (2.31) and identity} \} \\ & (|[\underline{Nil}, \text{Cons}]|) \\ = & \quad \{ \text{cata-reflection (3.68)} \} \\ & id \end{aligned}$$

Therefore, by defining

$$\text{List } f \stackrel{\text{def}}{=} (|[\underline{Nil}, \text{Cons} \cdot (f \times id)]|)$$

what we have just seen can be written thus:

$$\text{List } id_A = id_{\text{List } A}$$

This is nothing but law (3.55) for F replaced by List . Moreover, it will not be too difficult to check that

$$\text{List } (g \cdot f) = \text{List } g \cdot \text{List } f$$

also holds — cf. (3.56). Altogether, this means that List can be regarded as a functor.

In programming terminology one says that $\text{List } X$ (the “lists of X s datatype”) is *parametric* and that, by instantiating parameter X , one gets ground lists such as lists of integers, booleans, *etc.* The illustration above deepens one’s understanding of parameterization by identifying the functorial behaviour of the parametric datatype along with its parameter instantiations.

All this can be broadly generalized and leads to what is commonly known by a *type functor*. First of all, it should be clear that the generic format

$$T \cong F T$$

adopted so far for the definition of an inductive type is not sufficiently detailed because it does not provide a parametric view of T . For simplicity, let us suppose (for the moment) that only one parameter is identified in T . Then we may factor this out via *type variable* X and write (overloading symbol T)

$$T X \cong B(X, T X)$$

where B is called the type’s *base functor*. Binary functor $B(X, Y)$ is given this name because it is the basis of the whole inductive type

definition. By instantiation of X one obtains F . In the example above, $B(X, Y) = 1 + X \times Y$ and in fact $F Y = B(\mathbb{N}_0, Y) = 1 + \mathbb{N}_0 \times Y$, recall (3.51). Moreover, one has

$$F f = B(id, f) \quad (3.75)$$

and so every F -homomorphism can be written in terms of the base-functor of F , e.g.

$$f \cdot \alpha = \beta \cdot B(id, f)$$

instead of (3.65).

$\mathsf{T} X$ will be referred to as the *type functor* generated by B :

$$\underbrace{\mathsf{T} X}_{\text{type functor}} \cong \underbrace{B(X, \mathsf{T} X)}_{\text{base functor}}$$

We proceed to the description of its functorial behaviour — $\mathsf{T} f$ — for a given $B \xleftarrow{f} A$. As far as typing rules are concerned, we shall have

$$\frac{B \xleftarrow{f} A}{\mathsf{T} B \xleftarrow{\mathsf{T} f} \mathsf{T} A}$$

So we should be able to express $\mathsf{T} f$ as a $B(A, _)$ -catamorphism $\langle\!\langle g \rangle\!\rangle$:

$$\begin{array}{ccc} A & & \mathsf{T} A \xleftarrow{\text{in}_{\mathsf{T} A}} B(A, \mathsf{T} A) \\ f \downarrow & \mathsf{T} f = \langle\!\langle g \rangle\!\rangle \downarrow & \downarrow B(id, \mathsf{T} f) \\ B & \mathsf{T} B \xleftarrow{g} B(A, \mathsf{T} B) \end{array}$$

As we know that $\text{in}_{\mathsf{T} B}$ is the standard constructor of values of type $\mathsf{T} B$, we may put it into the diagram too:

$$\begin{array}{ccc} A & & \mathsf{T} A \xleftarrow{\text{in}_{\mathsf{T} A}} B(A, \mathsf{T} A) \\ f \downarrow & \mathsf{T} f = \langle\!\langle g \rangle\!\rangle \downarrow & \downarrow B(id, \mathsf{T} f) \\ B & \mathsf{T} B \xleftarrow{g} B(A, \mathsf{T} B) \\ & \swarrow \text{in}_{\mathsf{T} B} & \nearrow \text{dashed} \\ & B(B, \mathsf{T} B) \end{array}$$

The catamorphism's gene g will be synthesized by filling the dashed arrow in the diagram with the "obvious" $B(f, id)$, whereby one gets

$$\mathsf{T} f \stackrel{\text{def}}{=} \langle\!\langle \text{in}_{\mathsf{T} B} \cdot B(f, id) \rangle\!\rangle \quad (3.76)$$

and a final diagram, where $\text{in}_{\mathsf{T} A}$ is abbreviated by in_A (ibid. $\text{in}_{\mathsf{T} B}$ by in_B):

$$\begin{array}{ccc} A & & \mathsf{T} A \xleftarrow{\text{in}_A} B(A, \mathsf{T} A) \\ f \downarrow & \mathsf{T} f = \langle\!\langle \text{in}_B \cdot B(f, id) \rangle\!\rangle \downarrow & \downarrow B(id, \mathsf{T} f) \\ B & \mathsf{T} B \xleftarrow{\text{in}_B} B(B, \mathsf{T} B) \xleftarrow{B(f, id)} B(A, \mathsf{T} B) \end{array}$$

Next, we proceed to derive the useful law of *cata-absorption*

$$(\llbracket g \rrbracket) \cdot \mathsf{T} f = (\llbracket g \cdot \mathsf{B}(f, \mathsf{id}) \rrbracket) \quad (3.77)$$

as consequence of the laws studied in section 3.12. Our target is to show that, for $k = (\llbracket g \rrbracket) \cdot \mathsf{T} f$ in (3.67), one gets $\alpha = g \cdot \mathsf{B}(f, \mathsf{id})$:

$$\begin{aligned}
 & (\llbracket g \rrbracket) \cdot \mathsf{T} f = (\llbracket \alpha \rrbracket) \\
 \Leftrightarrow & \quad \{ \text{type-functor definition (3.76)} \} \\
 & (\llbracket g \rrbracket) \cdot (\llbracket \mathsf{in}_B \cdot \mathsf{B}(f, \mathsf{id}) \rrbracket) = (\llbracket \alpha \rrbracket) \\
 \Leftarrow & \quad \{ \text{cata-fusion (3.71)} \} \\
 & (\llbracket g \rrbracket) \cdot \mathsf{in}_B \cdot \mathsf{B}(f, \mathsf{id}) = \alpha \cdot \mathsf{B}(\mathsf{id}, (\llbracket g \rrbracket)) \\
 \Leftrightarrow & \quad \{ \text{cata-cancellation (3.69)} \} \\
 & g \cdot \mathsf{B}(\mathsf{id}, (\llbracket g \rrbracket)) \cdot \mathsf{B}(f, \mathsf{id}) = \alpha \cdot \mathsf{B}(\mathsf{id}, (\llbracket g \rrbracket)) \\
 \Leftrightarrow & \quad \{ \text{B is a bi-functor (3.58)} \} \\
 & g \cdot \mathsf{B}(\mathsf{id} \cdot f, (\llbracket g \rrbracket) \cdot \mathsf{id}) = \alpha \cdot \mathsf{B}(\mathsf{id}, (\llbracket g \rrbracket)) \\
 \Leftrightarrow & \quad \{ \text{id is natural (2.11)} \} \\
 & g \cdot \mathsf{B}(f \cdot \mathsf{id}, \mathsf{id} \cdot (\llbracket g \rrbracket)) = \alpha \cdot \mathsf{B}(\mathsf{id}, (\llbracket g \rrbracket)) \\
 \Leftrightarrow & \quad \{ (3.58) \text{ again, this time from left to right} \} \\
 & g \cdot \mathsf{B}(f, \mathsf{id}) \cdot \mathsf{B}(\mathsf{id}, (\llbracket g \rrbracket)) = \alpha \cdot \mathsf{B}(\mathsf{id}, (\llbracket g \rrbracket)) \\
 \Leftarrow & \quad \{ \text{Leibniz} \} \\
 & g \cdot \mathsf{B}(f, \mathsf{id}) = \alpha
 \end{aligned}$$

The following diagram pictures this property of catamorphisms:

$$\begin{array}{ccccc}
 A & & \mathsf{T} A & \xleftarrow{\mathsf{in}_A} & \mathsf{B}(A, \mathsf{T} A) \\
 f \downarrow & & \mathsf{T} f \downarrow & & \downarrow \mathsf{B}(\mathsf{id}, \mathsf{T} f) \\
 C & & \mathsf{T} C & \xleftarrow{\mathsf{in}_C} \mathsf{B}(C, \mathsf{T} C) & \xleftarrow{\mathsf{B}(f, \mathsf{id})} \mathsf{B}(A, \mathsf{T} C) \\
 & & \downarrow \llbracket g \rrbracket & \downarrow \mathsf{B}(\mathsf{id}, \llbracket g \rrbracket) & \downarrow \mathsf{B}(\mathsf{id}, \llbracket g \rrbracket) \\
 & & D & \xleftarrow{g} \mathsf{B}(C, D) & \xleftarrow{\mathsf{B}(f, \mathsf{id})} \mathsf{B}(A, D)
 \end{array}$$

It remains to show that (3.76) indeed defines a functor. This can be verified by checking properties (3.55) and (3.56) for $F = \mathsf{T}$:

- Property **type-functor-id**, cf. (3.55):

$$\begin{aligned}
 & \mathsf{T} \mathsf{id} \\
 = & \quad \{ \text{by definition (3.76)} \} \\
 & (\llbracket \mathsf{in}_B \cdot \mathsf{B}(\mathsf{id}, \mathsf{id}) \rrbracket) \\
 = & \quad \{ \text{B is a bi-functor (3.57)} \}
 \end{aligned}$$

$$\begin{aligned}
& (\|in_B \cdot id\|) \\
= & \quad \{ \text{identity and cata-reflection (3.68)} \} \\
& id
\end{aligned}$$

- Property **type-functor**, cf. (3.56) :

$$\begin{aligned}
& T(f \cdot g) \\
= & \quad \{ \text{by definition (3.76)} \} \\
& (\|in_B \cdot B(f \cdot g, id)\|) \\
= & \quad \{ id \cdot id = id \text{ and } B \text{ is a bi-functor (3.58)} \} \\
& (\|in_B \cdot B(f, id) \cdot B(g, id)\|) \\
= & \quad \{ \text{cata-absorption (3.77)} \} \\
& (\|in_B \cdot B(f, id)\|) \cdot Tg \\
= & \quad \{ \text{again cata-absorption (3.77)} \} \\
& (\|in_B\|) \cdot Tf \cdot Tg \\
= & \quad \{ \text{cata-reflection (3.68) followed by identity} \} \\
& Tf \cdot Tg
\end{aligned}$$

Exercise 3.11. Function $\text{length} = (\|[zero, succ \cdot \pi_2]\|)$ counts the number of elements of a finite list. If the input list has at least one element it suffices to count the elements of its tail starting with count 1 instead of 0:

$$\text{length} \cdot (a:) = (\|[one, succ \cdot \pi_2]\|) \quad (3.78)$$

Prove (3.78) knowing that

$$\text{length} \cdot (a:) = succ \cdot \text{length}$$

follows from the definition of length . (**NB:** assume $zero _ = 0$ and $one _ = 1$.)

□

Exercise 3.12. Function concat , extracted from Haskell's Prelude, can be defined as list catamorphism,

$$\text{concat} = (\|[nil, conc]\|) \quad (3.79)$$

where $\text{conc } (x, y) = x ++ y$, $nil _ = []$, $B(f, g) = id + f \times g$, $Ff = B(id, f)$, and $Tf = \text{map } f$. Prove property

$$\text{length} \cdot \text{concat} = \text{sum} \cdot \text{map length} \quad (3.80)$$

resorting to cata-fusion and cata-absorption.

□

3.14 A CATALOGUE OF STANDARD POLYNOMIAL INDUCTIVE TYPES

The following table contains a collection of standard polynomial inductive types and associated base type bi-functors, which are in canonical form (3.64). The table contains two extra columns which may be used as bookmarks for equations (3.75) and (3.76), respectively ¹⁴:

Description	$T X$	$B(X, Y)$	$B(id, f)$	$B(f, id)$
"Right" Lists	List X	$1 + X \times Y$	$id + id \times f$	$id + f \times id$
"Left" Lists	LList X	$1 + Y \times X$	$id + f \times id$	$id + id \times f$
Non-empty Lists	NList X	$X + X \times Y$	$id + id \times f$	$f + f \times id$
Binary Trees	BTree X	$1 + X \times Y^2$	$id + id \times f^2$	$id + f \times id$
"Leaf" Trees	LTree X	$X + Y^2$	$id + f^2$	$f + id$

(3.81)

All type functors T in this table are unary. In general, one may think of inductive datatypes which exhibit more than one type parameter. Should n parameters be identified in T , then this will be based on an $n + 1$ -ary base functor B , cf.

$$T(X_1, \dots, X_n) \cong B(X_1, \dots, X_n, T(X_1, \dots, X_n))$$

So, every $n + 1$ -ary polynomial functor $B(X_1, \dots, X_n, X_{n+1})$ can be identified as the basis of an inductive n -ary type functor (the convention is to stick to the canonical form and reserve the last variable X_{n+1} for the "recursive call"). While type bi-functors ($n = 2$) are often found in programming, the situation in which $n > 2$ is relatively rare. For instance, the combination of leaf-trees with binary-trees in (3.81) leads to the so-called "full tree" type bi-functor

Description	$T(X_1, X_2)$	$B(X_1, X_2, Y)$	$B(id, id, f)$	$B(f, g, id)$
"Full" Trees	FTree(X_1, X_2)	$X_1 + X_2 \times Y^2$	$id + id \times f^2$	$f + g \times id$

(3.82)

As we shall see later on, these types are widely used in programming. In the actual encoding of these types in HASKELL, exponentials are normally expanded to products according to (2.102), see for instance

$$\mathbf{data} \text{ BTree } a = \text{Empty} \mid \text{Node } (a, (\text{BTree } a, \text{BTree } a)) \quad (3.83)$$

Moreover, one may choose to curry the type constructors as in, e.g.

$$\mathbf{data} \text{ BTree } a = \text{Empty} \mid \text{Node } a \ (\text{BTree } a) \ (\text{BTree } a)$$

Exercise 3.13. Write as a catamorphisms

- the function which counts the number of elements of a non-empty list (type NList in (3.81)).
- the function which computes the maximum element of a binary-tree of natural numbers.

¹⁴ Since $(id_A)^2 = id_{(A^2)}$ one writes id^2 for id in this table.

□

Exercise 3.14. Let

$$\begin{aligned}\text{nil } _ &= [] \\ \text{singl } a &= [a]\end{aligned}$$

be defined. Characterize the function which is defined by $(\llbracket \text{nil}, h \rrbracket)$ for each of the following definitions of h :

$$h(x, (y_1, y_2)) = y_1 ++ [x] ++ y_2 \quad (3.84)$$

$$h = ++ \cdot (\text{singl} \times ++) \quad (3.85)$$

$$h = ++ \cdot (++ \times \text{singl}) \cdot \text{swap} \quad (3.86)$$

Identify in (3.81) which datatypes are involved as base functors.

□

Exercise 3.15. Write as a catamorphism the function which computes the frontier of a tree of type `LTree` (3.81), listed from left to right.

□

Exercise 3.16. Function

$$\begin{aligned}\text{mirror } (\text{Leaf } a) &= \text{Leaf } a \\ \text{mirror } (\text{Fork } (x, y)) &= \text{Fork } (\text{mirror } y, \text{mirror } x)\end{aligned}$$

which mirrors binary trees of type `LTree` $a = \text{Leaf } a \mid \text{Fork } (\text{LTree } a, \text{LTree } a)$ can be defined both as a catamorphism

$$\text{mirror} = (\llbracket \text{in} \cdot (\text{id} + \text{swap}) \rrbracket) \quad (3.87)$$

and as an anamorphism

$$\text{mirror} = \llbracket ((\text{id} + \text{swap}) \cdot \text{out}) \rrbracket \quad (3.88)$$

where `out` is the converse of

$$\text{in} = [\text{Leaf}, \text{Fork}] \quad (3.89)$$

Show that both definitions are effectively the same, that is, complete the etc steps of the reasoning:

$$\begin{aligned}\text{mirror} &= (\llbracket \text{in} \cdot (\text{id} + \text{swap}) \rrbracket) \\ &\equiv \{ \dots \text{etc} \dots \} \\ \text{mirror} &= \llbracket ((\text{id} + \text{swap}) \cdot \text{out}) \rrbracket\end{aligned}$$

□



Figure 3.1.: Towers of Hanoi.

(*Hint: recall that $Ff = id + f \times f$ for this type and mind the natural property of $id + \text{swap}$.*)

□

Exercise 3.17. Let parametric type T be given with base B , that is, such that $Tf = (\lambda \text{ in} \cdot B(f, id))$. Define the so-called triangular combinator of T , $tri f$, as follows:

$$tri f = (\lambda \text{ in} \cdot B(id, Tf)) \quad (3.90)$$

Show that the instance of this combinator for type $LTree\ a = Leaf\ a \mid Fork\ (LTree\ a, LTree\ a)$ — such that $\text{in} = [Leaf, Fork]$ and $B(f, g) = f + g \times g$ — is the following function

$$\begin{aligned} tri &:: (a \rightarrow a) \rightarrow LTree\ a \rightarrow LTree\ a \\ tri\ f\ (Leaf\ x) &= Leaf\ x \\ tri\ f\ (Fork\ (t, t')) &= Fork\ (fmap\ f\ (tri\ f\ t), fmap\ f\ (tri\ f\ t')) \end{aligned}$$

written in Haskell syntax.

□

3.15 HYLO-FACTORIZATION

A well-known example of a hylomorphism is the algorithm than computes the sequence of disk moves in the Towers of Hanoi puzzle:

$$\begin{aligned} hanoi\ (d, 0) &= [] \\ hanoi\ (d, n + 1) &= hanoi\ (\neg d, n) ++ [(n, d)] ++ hanoi\ (\neg d, n) \end{aligned} \quad (3.91)$$

Here is a nice account of this puzzle, taken from [5]:

The Towers of Hanoi problem comes from a puzzle marketed in 1883 by the French mathematician Édouard Lucas, under the pseudonym Claus. The puzzle is based on a legend according to which there is a temple, apparently in Bramah rather than in Hanoi as one might expect, where there are three giant poles fixed in the ground. On the first of these poles, at the time of the world's creation, God placed sixty

four golden disks, each of different size, in decreasing order of size. The Bramin monks were given the task of moving the disks, one per day, from one pole to another subject to the rule that no disk may ever be above a smaller disk. The monks' task would be complete when they had succeeded in moving all the disks from the first of the poles to the second and, on the day that they completed their task the world would come to an end!

There is a wellknown inductive solution to the problem (...) In this solution we make use of the fact that the given problem is symmetrical with respect to all three poles. Thus it is undesirable to name the individual poles. Instead we visualize the poles as being arranged in a circle [See figure 3.1]; the problem is to move the tower of disks from one pole to the next pole in a specified direction around the circle. The code defines $H \ n \ d \S$ to be a sequence of pairs (k, d') where n is the number of disks, k is a disk number and d and d' are directions. Disks are numbered from 0 onwards, disk 0 being the smallest. (Assigning number 0 to the smallest rather than the largest disk has the advantage that the number of the disk that is moved on any day is independent of the total number of disks to be moved.) Directions are boolean values, true representing a clockwise movement and false an anticlockwise movement. The pair (k, d') means move the disk numbered k from its current position in the direction d' . (...) Taking the pairs in order from left to right, the complete sequence (...) prescribes how to move the n smallest disks one-by-one from one pole to the next pole in the direction d following the rule of never placing a larger disk on top of a smaller disk.

Next, here is how the same function (3.91) can be viewed as a hylo-morphism:¹⁵

$$\begin{aligned} hanoi &= ([inord]) \cdot [(strategy)] \text{ where} \\ strategy \ (d, 0) &= i_1 \ () \\ strategy \ (d, n+1) &= i_2 \ ((d, n), ((\neg d, n), (\neg d, n))) \\ inord &= [nil, f] \\ f \ (x, (l, r)) &= l ++ [x] ++ r \end{aligned}$$

This means that, for some functor F ,

$$hanoi = inord \cdot F \ hanoi \cdot strategy \quad (3.92)$$

holds. The question is: what is the functor F capturing the *recursive pattern* of the algorithm? From $strategy \ (d, 0) = i_1 \ ()$ we infer the type

$$strategy : \mathbb{B} \times \mathbb{N}_0 \rightarrow 1 + \dots$$

and from $strategy \ (d, n+1) = i_2 \ ((d, n), ((\neg d, n), (\neg d, n)))$ we infer

$$strategy : \mathbb{B} \times \mathbb{N}_0 \rightarrow \dots + (\mathbb{B} \times \mathbb{N}_0) \times (\mathbb{B} \times \mathbb{N}_0)^2$$

Altogether:

$$\begin{array}{ccc} (\mathbb{B} \times \mathbb{N}_0)^* & \xleftarrow{inord} & 1 + (\mathbb{B} \times \mathbb{N}_0) \times ((\mathbb{B} \times \mathbb{N}_0)^*)^2 \\ \uparrow hanoi & & \uparrow id + id \times hanoi^2 \\ \mathbb{B} \times \mathbb{N}_0 & \xrightarrow{strategy} & 1 + (\mathbb{B} \times \mathbb{N}_0) \times (\mathbb{B} \times \mathbb{N}_0)^2 \end{array}$$

¹⁵ Recall (3.84) concerning function $inord$.

We conclude that $F X = 1 + (\mathbb{B} \times \mathbb{N}_0) \times X^2$:

$$\begin{array}{ccc} (\mathbb{B} \times \mathbb{N}_0)^* & \xleftarrow{\text{inord}} & F (\mathbb{B} \times \mathbb{N}_0)^* \\ \text{hanoi} \uparrow & & \uparrow F \text{ hanoi} \\ \mathbb{B} \times \mathbb{N}_0 & \xrightarrow{\text{strategy}} & F (\mathbb{B} \times \mathbb{N}_0) \end{array}$$

Since $F X = B(Y, X)$ for some B , we get

$$F X = B(\mathbb{B} \times \mathbb{N}_0, X)$$

for $B(Y, X) = 1 + Y \times X^2$. Finally, from the table in (3.81) we conclude that the intermediate (virtual) structure of the *hanoi* hylomorphism is $BTree(\mathbb{B} \times \mathbb{N}_0)$:

$$\begin{array}{ccccc} & & (\mathbb{B} \times \mathbb{N}_0)^* & \xleftarrow{\text{inord}} & F (\mathbb{B} \times \mathbb{N}_0)^* \\ & \text{hanoi} \nearrow & \uparrow \langle \text{inord} \rangle & & \uparrow F \langle \text{inord} \rangle \\ & & BTree(\mathbb{B} \times \mathbb{N}_0) & \xrightleftharpoons[\text{in}]{\text{out}} & F(BTree(\mathbb{B} \times \mathbb{N}_0)) \\ & \nwarrow & \uparrow \llbracket \text{strategy} \rrbracket & & \uparrow F \llbracket \text{strategy} \rrbracket \\ & & \mathbb{B} \times \mathbb{N}_0 & \xrightarrow{\text{strategy}} & F(\mathbb{B} \times \mathbb{N}_0) \end{array}$$

F hanoi

Exercise 3.18. Show that (3.92) unfolds to (3.91) for $F X = 1 + (\mathbb{B} \times \mathbb{N}_0) \times X^2$.

□

Exercise 3.19. From the *hanoi* function (3.91) one can derive the function that gives the total number of disk movements of the puzzle:

$$\begin{aligned} nm\ 0 &= 0 \\ nm\ (n + 1) &= 2 * (nm\ n) + 1 \end{aligned}$$

That is:

$$nm\ n = \text{for odd } 0 \text{ where odd } n = 2 * n + 1 \quad (3.93)$$

Show that

$$nm\ n = 2^n - 1.$$

Hint: define $k\ n = 2^n - 1$ and solve the equation $k = \text{for odd } 0$ using the laws of catamorphisms and basic properties of arithmetics.

□

Exercise 3.20. From the pointwise version on ‘quicksort’,

$$qSort [] = []$$

$$qSort (h:t) = qSort [a \mid a \leftarrow t, a < h] ++ [h] ++ qSort [a \mid a \leftarrow t, a \geq h]$$

infer g and h in the hylo-factorization $qSort = \llbracket g \rrbracket \cdot \llbracket h \rrbracket$, knowing that the intermediate structure is a `BTree` as in the case of `hanoi`.

□

Exercise 3.21. Consider the well-known function which computes the n -th Fibonacci number:

$$\begin{aligned} fib\ 0 &= 1 \\ fib\ 1 &= 1 \\ fib\ (n+2) &= fib\ (n+1) + fib\ n \end{aligned}$$

Show that `fib` is a hylomorphism of type `LTree` (3.81),

$$fib = \llbracket count, fibd \rrbracket$$

for

$$\begin{aligned} count &= [1, add] \\ add\ (x, y) &= x + y \\ fibd\ 0 &= i_1\ () \\ fibd\ 1 &= i_1\ () \\ fibd\ (n+2) &= i_2\ (n+1, n) \end{aligned}$$

□

Exercise 3.22. Consider the following definition of the factorial function,

$$\begin{aligned} dfac\ 0 &= 1 \\ dfac\ n &= \llbracket [id, mul], dfacd \rrbracket\ (1, n) \end{aligned}$$

where

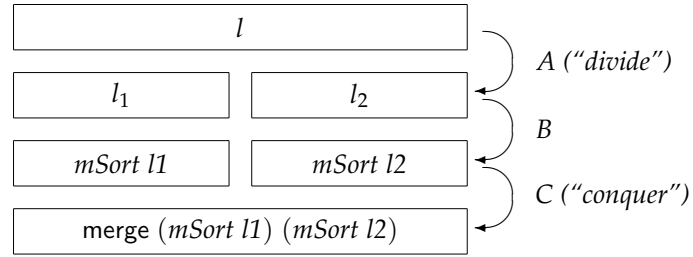
$$\begin{aligned} mul\ (x, y) &= x * y \\ dfacd\ (n, m) & \\ \quad | \ n \equiv m &= i_1\ n \\ \quad | \text{otherwise} &= i_2\ ((n, k), (k+1, m)) \text{ where } k = (n+m) \div 2 \end{aligned}$$

Derive from the above the corresponding (doubly recursive) pointwise definition of `dfac`. (This is known as the double factorial implementation of factorial.)

□

Exercise 3.23. The drawing below describes how the so-called merge sort algorithm works¹⁶:

¹⁶ Only the case of inputs with more than one element is depicted.



Define the function `merge` and then the hylomorphism

$$mSort = ([g]) \cdot [[singl, merge]]$$

(find g) knowing that its virtual data-structure is of type `LTree`. Note: the empty list should be left out of the hylomorphism and handled separately.

□

3.16 FUNCTORS AND TYPE FUNCTORS IN HASKELL

The concept of a (unary) functor is provided in HASKELL in the form of a particular class exporting the `fmap` operator:

```
class Functor f where
  fmap :: (a → b) → (f a → f b)
```

So `fmap g` encodes $F g$ once we declare `F` as an instance of `class Functor`. The most popular use of `fmap` has to do with HASKELL lists, as allowed by declaration

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

in language's *Standard Prelude*.

In order to encode the type functors we have seen so far we have to do the same concerning their declaration. For instance, should we write

```
instance Functor BTree
where fmap f = cataBTree (inBTree · (id + (f × id)))
```

concerning the binary-tree datatype of (3.81) and assuming appropriate declarations of `cataBTree` and `inBTree`, then `fmap` is overloaded and used across such binary-trees.

Bi-functors can be added easily by writing

```
class BiFunctor f where
  bmap :: (a → b) → (c → d) → (f a c → f b d)
```

Exercise 3.24. Declare all datatypes in (3.81) in HASKELL notation and turn them into HASKELL type functors, that is, define *fmap* in each case.

□

Exercise 3.25. Declare datatype (3.82) in HASKELL notation and turn it into an instance of class *BiFunctor*.

□

3.17 THE MUTUAL-RECURSION LAW

The theory developed so far for building (and reasoning about) recursive functions doesn't cope with mutual recursion. As a matter of fact, the pattern of recursion of a given cata(ana,hylo)morphism involves only the recursive function being defined, even though more than once, in general, as dictated by the relevant base functor.

It turns out that rules for handling mutual recursion are surprisingly simple to calculate. As motivation, recall section 2.10 where, by mixing products with coproducts, we obtained a result — the *exchange rule* (2.49) — which stemmed from putting together the two universal properties of product and coproduct, (2.63) and (2.65), respectively.

The question we want to address in this section is of the same brand: *what can one tell about catamorphisms which output pairs of values?* By (2.63), such catamorphisms are bound to be *splits*, as are the corresponding *genes*.¹⁷

$$\begin{array}{ccc}
 T & \xleftarrow{\text{in}} & F\,T \\
 \downarrow \langle \langle h, k \rangle \rangle & & \downarrow F\, \langle \langle h, k \rangle \rangle \\
 A \times B & \xleftarrow{\langle h, k \rangle} & F\,(A \times B)
 \end{array}$$

As we did for the exchange rule, we put (2.63) and the universal property of catamorphisms (3.67) against each other and calculate:

$$\begin{aligned}
 \langle f, g \rangle &= \langle \langle h, k \rangle \rangle \\
 &\equiv \{ \text{cata-universal (3.67)} \} \\
 \langle f, g \rangle \cdot \text{in} &= \langle h, k \rangle \cdot F\, \langle f, g \rangle \\
 &\equiv \{ \times\text{-fusion (2.26) twice} \} \\
 \langle f \cdot \text{in}, g \cdot \text{in} \rangle &= \langle h \cdot F\, \langle f, g \rangle, k \cdot F\, \langle f, g \rangle \rangle \\
 &\equiv \{ (2.64) \} \\
 &\quad \left\{ \begin{array}{l} f \cdot \text{in} = h \cdot F\, \langle f, g \rangle \\ g \cdot \text{in} = k \cdot F\, \langle f, g \rangle \end{array} \right.
 \end{aligned}$$

¹⁷ Using *T* to denote μ_F .

The rule thus obtained,

$$\begin{cases} f \cdot \text{in} = h \cdot F \langle f, g \rangle \\ g \cdot \text{in} = k \cdot F \langle f, g \rangle \end{cases} \equiv \langle f, g \rangle = (\llbracket \langle h, k \rangle \rrbracket) \quad (3.94)$$

is referred to as the *mutual recursion law* (or as “Fokkinga’s law”) and is useful in combining two mutually recursive functions f and g

$$\begin{array}{ccc} T & \xleftarrow{\text{in}} & F T \\ f \downarrow & & \downarrow F \langle f, g \rangle \\ A & \xleftarrow{h} & F(A \times B) \end{array} \quad \begin{array}{ccc} T & \xleftarrow{\text{in}} & F T \\ g \downarrow & & \downarrow F \langle f, g \rangle \\ B & \xleftarrow{k} & F(A \times B) \end{array}$$

into a single catamorphism.

When applied from left to right, law (3.94) is surprisingly useful in optimizing recursive functions in a way which saves redundant traversals of the input inductive type T . Let us take the Fibonacci function as example:

$$\begin{aligned} \text{fib } 0 &= 1 \\ \text{fib } 1 &= 1 \\ \text{fib}(n+2) &= \text{fib}(n+1) + \text{fib } n \end{aligned}$$

It can be shown — recall exercise 3.21 — that fib is a hylomorphism of type LTree (3.81). This hylo-factorization of fib tells something about its internal algorithmic structure: the *divide step* $\llbracket \text{fibd} \rrbracket$ builds a tree whose number of leaves is a Fibonacci number; the *conquer step* $\llbracket \text{count} \rrbracket$ just counts such leaves.

There is, of course, much re-calculation in this hylomorphism. Can we improve its performance? The clue is to regard the two instances of fib in the recursive branch as mutually recursive over the natural numbers. This clue is suggested not only by fib having two base cases (so, perhaps it hides two functions) but also by the lookahead $n+2$ in the recursive clause.

We start by defining a function which reduces such a lookahead by 1,

$$f \ n = \text{fib}(n+1)$$

Clearly, $f(n+1) = \text{fib}(n+2) = f \ n + \text{fib } n$ and $f \ 0 = \text{fib } 1 = 1$. Putting f and fib together,

$$\begin{aligned} f \ 0 &= 1 \\ f(n+1) &= f \ n + \text{fib } n \\ \text{fib } 0 &= 1 \\ \text{fib}(n+1) &= f \ n \end{aligned}$$

we obtain two mutually recursive functions over the natural numbers (\mathbb{N}_0) which transform into pointfree equalities

$$\begin{aligned} f \cdot [0, \text{suc}] &= [1, \text{add} \cdot \langle f, \text{fib} \rangle] \\ \text{fib} \cdot [0, \text{suc}] &= [1, f] \end{aligned}$$

over

$$\begin{array}{ccc} \mathbb{N}_0 & \xrightarrow{\cong} & 1 + \mathbb{N}_0 \\ & \text{in}=[0, \text{succ}] & \underbrace{\quad}_{F \mathbb{N}_0} \end{array} \quad (3.95)$$

Reverse $+$ -absorption (2.43) will further enable us to rewrite the above into

$$\begin{aligned} f \cdot \text{in} &= [\underline{1}, \text{add}] \cdot F \langle f, \text{fib} \rangle \\ \text{fib} \cdot \text{in} &= [\underline{1}, \pi_1] \cdot F \langle f, \text{fib} \rangle \end{aligned}$$

thus bringing functor $F f = \text{id} + f$ explicit and preparing for mutual recursion removal:

$$\begin{aligned} f \cdot \text{in} &= [\underline{1}, \text{add}] \cdot F \langle f, \text{fib} \rangle \\ \text{fib} \cdot \text{in} &= [\underline{1}, \pi_1] \cdot F \langle f, \text{fib} \rangle \\ \equiv & \{ (3.94) \} \\ \langle f, \text{fib} \rangle &= (\llbracket [\underline{1}, \text{add}], [\underline{1}, \pi_1] \rrbracket) \\ \equiv & \{ \text{exchange law (2.49)} \} \\ \langle f, \text{fib} \rangle &= (\llbracket \langle \underline{1}, \underline{1} \rangle, \langle \text{add}, \pi_1 \rangle \rrbracket) \\ \equiv & \{ \text{going pointwise and denoting } \langle f, \text{fib} \rangle \text{ by } \text{fib}' \} \\ & \left\{ \begin{array}{l} \text{fib}' 0 = (1, 1) \\ \text{fib}' (n+1) = (x+y, x) \text{ where } (x, y) = \text{fib}' n \end{array} \right. \end{aligned}$$

Since $\text{fib} = \pi_2 \cdot \text{fib}'$ we easily recover fib from fib' and obtain the intended linear version of Fibonacci, below encoded in Haskell:

```
fib n = m where
  (_, m) = fib' n
  fib' 0 = (1, 1)
  fib' (n+1) = (x+y, x) where (x, y) = fib' n
```

This version of fib is actually the semantics of the “for-loop” — recall (3.7) — one would write in an imperative language which would initialize two global variables $x, y := 1, 1$, loop over assignment $x, y := x+y, x$ and yield the result in y . In the C programming language, one would write

```
int fib(int n)
{
  int x=1; int y=1; int i;
  for (i=1; i<=n; i++) {int a=x; x=x+y; y=a;}
  return y;
};
```

where the extra variable a is required for ensuring that *simultaneous* assignment $x, y := x+y, x$ takes place in a sequential way.

Recall from section 3.1 that all \mathbb{N}_0 catamorphisms are of shape $(\llbracket k, g \rrbracket)$ and such that $(\llbracket k, g \rrbracket)n = g^n k$, where g^n is the n -th iteration of g , that is, $g^0 = id$ and $g^{n+1} = g \cdot g^n$. That is, g is the body of a “for-loop” which repeats itself n -times, starting with initial value k . Recall also that the for-loop combinator is nothing but the “fold combinator” (3.5) associated to the natural number data type.

In a sense, the mutual recursion law gives us a hint on how global variables “are born” in computer programs, out of the maths definitions themselves. Quite often more than two such variables are required in linearizing hylomorphisms by mutual recursion. Let us see an example. The question is: *how many squares can one draw on a $n \times n$ -tiled wall?* The answer is given by function

$$ns\ n \stackrel{\text{def}}{=} \sum_{i=1,n} i^2$$

that is,

$$\begin{aligned} ns\ 0 &= 0 \\ ns\ (n+1) &= (n+1)^2 + ns\ n \end{aligned}$$

in Haskell. However, this hylomorphism is inefficient because each iteration involves another hylomorphism computing square numbers.

One way of improving ns is to introduce function $bnm\ n \stackrel{\text{def}}{=} (n+1)^2$ and express this over (3.95),

$$\begin{aligned} bnm\ 0 &= 1 \\ bnm(n+1) &= 2n+3 + bnm\ n \end{aligned}$$

hoping to blend ns with bnm using the mutual recursion law. However, the same problem arises in bnm itself, which now depends on term $2n+3$. We invent $lin\ n \stackrel{\text{def}}{=} 2n+3$ and repeat the process, thus obtaining:

$$\begin{aligned} lin\ 0 &= 3 \\ lin(n+1) &= 2 + lin\ n \end{aligned}$$

By redefining

$$\begin{aligned} bnm'\ 0 &= 1 \\ bnm'(n+1) &= lin\ n + bnm'\ n \end{aligned}$$

$$\begin{aligned} ns'\ 0 &= 0 \\ ns'(n+1) &= bnm'\ n + ns'\ n \end{aligned}$$

we obtain three functions — ns' , bnm' and lin — mutually recursive over the polynomial base $F\ g = id + g$ of the natural numbers.

Exercise 3.29 below shows how to extend (3.94) to three mutually recursive functions (3.96). (From this it is easy to extend it further to

the n -ary case.) It is routine work to show that, by application of (3.96) to the above three functions, one obtains the linear version of ns which follows:

$$\begin{aligned} ns''\ n &= a \textbf{ where} \\ (a, -, -) &= aux\ n \\ aux\ 0 &= (0, 1, 3) \\ aux\ (n + 1) &= \textbf{let}\ (a, b, c) = aux\ n \textbf{ in}\ (a + b, b + c, 2 + c) \end{aligned}$$

In retrospect, note that (in general) not every system of n mutually recursive functions

$$\begin{cases} f_1 = \phi_1(f_1, \dots, f_n) \\ \vdots \\ f_n = \phi_n(f_1, \dots, f_n) \end{cases}$$

involving n functions and n functional combinators ϕ_1, \dots, ϕ_n can be handled by a suitably extended version of (3.94). This only happens if all f_i have the same “shape”, that is, if they share the same base functor F .

Exercise 3.26. Use the mutual recursion law (3.94) to show that each of the two functions

$$\begin{cases} odd\ 0 = False \\ odd(n + 1) = even\ n \end{cases} \quad \begin{cases} even\ 0 = True \\ even(n + 1) = odd\ n \end{cases}$$

checking natural number parity can be expressed as a projection of

for swap (False, True)

Encode this for-loop in C syntax.

□

Exercise 3.27. The following Haskell function computes the list of the first n natural numbers in reverse order:

$$\begin{aligned} insg\ 0 &= [] \\ insg\ (n + 1) &= (n + 1) : insg\ n \end{aligned}$$

1. Show that $insg$ can also be defined as follows:

$$\begin{aligned} insg\ 0 &= [] \\ insg\ (n + 1) &= (fsuc\ n) : insg\ n \\ fsuc\ 0 &= 1 \\ fsuc\ (n + 1) &= fsuc\ n + 1 \end{aligned}$$

2. Based on the mutual recursion law derive from such a definition the following version of $insg$ encoded as a for-loop:

$$\begin{aligned} insg &= \pi_2 \cdot insgfor \\ insgfor &= \text{for } \langle (1+) \cdot \pi_1, \text{cons} \rangle \underline{(1, [])} \end{aligned}$$

where $\text{cons } (n, m) = n : m$.

□

Exercise 3.28. The number of steps that solve the Hanoi Towers “puzzle”, for n discs, is:

$$k \ n = 2^n - 1$$

— recall exercise 3.19. Using the mutual recursion law, show that another way of computing k is

$$\begin{aligned} k &= \pi_1 \cdot g \text{ where} \\ g &= \text{for loop } (0, 1) \\ \text{loop } (k, e) &= (k + e, 2 * e) \end{aligned}$$

knowing that

$$\begin{aligned} k \ 0 &= 0 \\ k \ (n + 1) &= 2^n + k \ n \end{aligned}$$

hold (as can be easily shown) and that $2^n = \text{for } (2*) \ 1 \ n$.

□

Exercise 3.29. Show that law (3.94) generalizes to more than two mutually recursive functions, in this case three:

$$\begin{cases} f \cdot \text{in} = h \cdot F \langle f, \langle g, j \rangle \rangle \\ g \cdot \text{in} = k \cdot F \langle f, \langle g, j \rangle \rangle \\ j \cdot \text{in} = l \cdot F \langle f, \langle g, j \rangle \rangle \end{cases} \quad \equiv \quad \langle f, \langle g, j \rangle \rangle = (\langle h, \langle k, l \rangle \rangle) \quad (3.96)$$

□

Exercise 3.30. The exponential function $e^x : \mathbb{R} \rightarrow \mathbb{R}$ (where “ e ” denotes Euler’s number) can be defined in several ways, one being the calculation of Taylor series:

$$e^x = \sum_{n=0}^{\infty} \frac{x^n}{n!} \quad (3.97)$$

The following function, in Haskell,

```
exp :: Double → Integer → Double
exp x 0 = 1
exp x (n + 1) = x ↑ (n + 1) / fac (n + 1) + (exp x n)
```

computes an approximation of e^x , where the second parameter tells how many terms to compute. For instance, while $\text{exp } 1 \ 1 = 2.0$, $\text{exp } 1 \ 10$ yields 2.7182818011463845.

Function $\text{exp } x \ n$ performs badly for n larger and larger: while $\text{exp } 1 \ 100$ runs instantaneously, $\text{exp } 1 \ 1000$ takes around 9 seconds, $\text{exp } 1 \ 2000$ takes circa 33 seconds, and so on.

Decompose exp into mutually recursive functions so as to apply (3.96) and obtain the following linear version,

```
exp x n = let (e,b,c) = aux x n
          in e where
            aux x 0 = (1,2,x)
            aux x (n+1) =
              let (e,s,h) = aux x n
              in (e+h,s+1,(x/s)*h)
```

which translates directly to the encoding in C:

```
float exp(float x, int n)
{
  float h=x; float e=1; int s=2; int i;
  for (i=0; i<n+1; i++) {e=e+h; h=(x/s)*h; s++;}
  return e;
};
```

□

Exercise 3.31. Show that, for all $n \in \mathbb{N}_0$, $n = \text{suc}^n 0$. **Hint:** use cata-reflexion (3.68).

□

MUTUAL RECURSION OVER LISTS. As example of application of (3.94) for μ_F other than \mathbb{N}_0 , consider the following recursive predicate which checks whether a (non-empty) list is ordered,

$$\begin{aligned} \text{ord} : A^+ &\rightarrow 2 \\ \text{ord}[a] &= \text{TRUE} \\ \text{ord}(\text{cons}(a,l)) &= a \geq (\text{listMax } l) \wedge (\text{ord } l) \end{aligned}$$

where \geq is assumed to be a total order on datatype A and

$$\text{listMax} = ([id, \text{max}]) \quad (3.98)$$

computes the greatest element of a given list of A s:

$$\begin{array}{ccc} A^+ & \xleftarrow{[singl, \text{cons}]} & A + A \times A^+ \\ \text{listMax} \downarrow & & \downarrow id + id \times \text{listMax} \\ A & \xleftarrow{[id, \text{max}]} & A + A \times A \end{array}$$

(In the diagram, $\text{singl } a = [a]$.)

Predicate *ord* is not a catamorphism because of the presence of *listMax l* in the recursive branch. However, the following diagram depicting *ord*

$$\begin{array}{ccc}
 A^+ & \xleftarrow{[singl, cons]} & A + A \times A^+ \\
 \text{ord} \downarrow & & \downarrow id + id \times \langle listMax, ord \rangle \\
 2 & \xleftarrow{[TRUE, \alpha]} & A + A \times (A \times 2)
 \end{array}$$

(where $\alpha(a, (m, b)) \stackrel{\text{def}}{=} a \geq m \wedge b$) suggests the possibility of using the mutual recursion law. One only has to find a way of letting *listMax* depend also on *ord*, which isn't difficult: for any $A^+ \xrightarrow{g} B$, one has

$$\begin{array}{ccc}
 A^+ & \xleftarrow{[singl, cons]} & A + A \times A^+ \\
 listMax \downarrow & & \downarrow id + id \times \langle listMax, g \rangle \\
 A & \xleftarrow{[id, max \cdot (id \times \pi_1)]} & A + A \times (A \times B)
 \end{array}$$

where the extra presence of *g* is cancelled by projection π_1 .

For $B = 2$ and $g = ord$ we are in position to apply Fokkinga's law and obtain:

$$\begin{aligned}
 \langle listMax, ord \rangle &= ([id, max \cdot (id \times \pi_1)], [TRUE, \alpha]) \\
 &= \{ \text{exchange law (2.49)} \} \\
 &= ([id, TRUE], [max \cdot (id \times \pi_1), \alpha])
 \end{aligned}$$

Of course, $ord = \pi_2 \cdot \langle listMax, ord \rangle$. By denoting the above synthesized catamorphism by *aux*, we end up with the following version of *ord*:

$$ord\ l = \mathbf{let}\ (a, b) = aux\ l\ \mathbf{in}\ b$$

where

$$\begin{aligned}
 aux &: A^+ \rightarrow A \times 2 \\
 aux\ [a] &= (a, \text{True}) \\
 aux\ (\text{cons}\ (a, l)) &= (\text{max}\ (a, m), a > m \wedge b) \ \mathbf{where}\ (m, b) = aux\ l
 \end{aligned}$$

Exercise 3.32. What do the following Haskell functions do?

$$\begin{aligned}
 f_1\ [] &= [] \\
 f_1\ (h:t) &= h : (f_2\ t) \\
 f_2\ [] &= [] \\
 f_2\ (h:t) &= f_1\ t
 \end{aligned}$$

Write $f = \langle f_1, f_2 \rangle$ as a list catamorphism and encode *f* back into Haskell syntax.

□

3.18 “BANANA-SPLIT”: A COROLLARY OF THE MUTUAL-RECURSION LAW

Let $h = i \cdot F \pi_1$ and $k = j \cdot F \pi_2$ in (3.94). Then

$$\begin{aligned}
 f \cdot \text{in} &= (i \cdot F \pi_1) \cdot F \langle f, g \rangle \\
 &\equiv \{ \text{composition is associative and } F \text{ is a functor} \} \\
 f \cdot \text{in} &= i \cdot F (\pi_1 \cdot \langle f, g \rangle) \\
 &\equiv \{ \text{by } \times\text{-cancellation (2.22)} \} \\
 f \cdot \text{in} &= i \cdot F f \\
 &\equiv \{ \text{by cata-cancellation} \} \\
 f &= \langle i \rangle
 \end{aligned}$$

Similarly, from $k = j \cdot F \pi_2$ we get

$$g = \langle j \rangle$$

Then, from (3.94), we get

$$\langle \langle i \rangle, \langle j \rangle \rangle = \langle \langle i \cdot F \pi_1, j \cdot F \pi_2 \rangle \rangle$$

that is

$$\langle \langle i \rangle, \langle j \rangle \rangle = \langle (i \times j) \cdot \langle F \pi_1, F \pi_2 \rangle \rangle \quad (3.99)$$

by (reverse) \times -absorption (2.27).

This law provides us with a very useful tool for “parallel loop” inter-combination: “loops” $\langle i \rangle$ and $\langle j \rangle$ are fused together into a single “loop” $\langle (i \times j) \cdot \langle F \pi_1, F \pi_2 \rangle \rangle$. The need for this kind of calculation arises very often. Consider, for instance, the function which computes the average of a non-empty list of natural numbers,

$$\text{average} \stackrel{\text{def}}{=} (/) \cdot \langle \text{sum}, \text{length} \rangle \quad (3.100)$$

where sum and length are the expected \mathbb{N}^+ catamorphisms:

$$\begin{aligned}
 \text{sum} &= \langle [id, +] \rangle \\
 \text{length} &= \langle [1, \text{succ} \cdot \pi_2] \rangle
 \end{aligned}$$

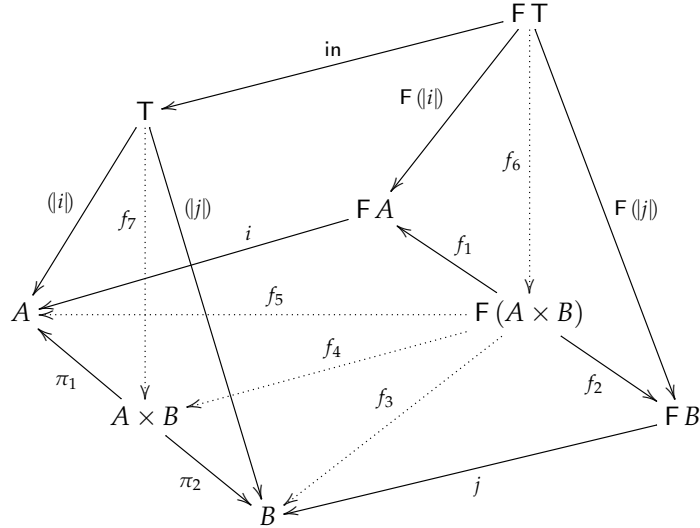
As defined by (3.100), function *average* performs two independent traversals of the argument list before division $(/)$ takes place. Banana-split will fuse such two traversals into a single one (see function *aux* below), thus leading to a function which will run “twice as fast”:

$$\begin{aligned}
 \text{average } l &= x / y \textbf{ where} \\
 (x, y) &= \text{aux } l \\
 \text{aux } [] &= (0, 0) \\
 \text{aux } (a : l) &= (a + x, y + 1) \textbf{ where } (x, y) = \text{aux } l
 \end{aligned} \quad (3.101)$$

Exercise 3.33. Calculate (3.101) from (3.100). Which of these two versions of the same function is easier to understand?

□

Exercise 3.34. The following diagram depicts “banana-split” (3.99):



Identify all functions f_1 to f_7 .

□

Exercise 3.35. Show that the standard Haskell function

$$\text{unzip } xs = (\text{map } \pi_1 \text{ } xs, \text{map } \pi_2 \text{ } xs)$$

can be defined as a catamorphism (fold) thanks to (3.99). Generalize this calculation to the generic *unzip* function over an inductive (polynomial) type T :

$$\text{unzip}_T = \langle T\pi_1, T\pi_2 \rangle$$

Suggestion: recall (3.76).

□

3.19 INDUCTIVE DATATYPE ISOMORPHISM

not yet available

3.20 BIBLIOGRAPHY NOTES

It is often the case that the expressive power of a particular programming language or paradigm is counter-productive in the sense that

too much freedom is given to programmers. Sooner or later, these will end up writing unintelligible (authorship dependent) code which will become a burden to whom has to maintain it. Such has been the case of imperative programming in the past (inc. assembly code), where the unrestricted use of `goto` instructions eventually gave place to `if-then-else`, `while` and `repeat` *structured* programming constructs.

A similar trend has been observed over the last decades at a higher programming level: arbitrary recursion and/or (side) effects have been considered harmful in functional programming. Instead, programmers have been invited to structure their code around generic program devices such as eg. *fold/unfold* combinators, which bring discipline to recursion. One witnesses progress in the sense that the loss of freedom is balanced by the increase of formal semantics and the availability of program calculi.

Such disciplined programming combinators have been extended from list-processing to other inductive structures thanks to one of the most significant advances in programming theory over the last decade: the so-called *functorial* approach to datatypes which originated mainly from [37], was popularized by [36] and reached textbook format in [10]. A comfortable basis for exploiting *polymorphism* [60], the “datatypes as functors” motto has proved beneficial at a higher level of abstraction, giving birth to *polytypism* [30].

The literature on *anas*, *catas* and *hylos* is vast (see eg. [40], [29], [20]) and it is part of a broader discipline which has become known as the *mathematics of program construction* [4]. This chapter provides an introduction to such discipline. Only the calculus of catamorphisms is presented. The corresponding theory of anamorphisms and hylo-morphisms demands further mathematical machinery (functions generalized to binary relations) and won’t be dealt with before chapter 8. The results on mutual recursion presented in this chapter, pioneered by Maarten Fokkinga [16], have been extended towards probabilistic functions [45]. They have also shown to help in program understanding and reverse engineering [59]. Recently, the whole theory has undergone significant advances through further use of category theory notions such as adjunctions¹⁸ and conjugate functors [24, 25].

¹⁸ See chapter 4.

WHY MONADS MATTER

In this chapter we present a powerful device in state-of-the-art functional programming, that of a *monad*. The monad concept is nowadays of primary importance in computing science because it makes it possible to describe computational effects as disparate as input/output, comprehension notation, state variable updating, probabilistic behaviour, context dependence, partial behaviour *etc.* in an elegant and uniform way.

Our motivation to this concept will start from a well-known problem in functional programming (and computing as a whole) — that of coping with undefined computations.

4.1 PARTIAL FUNCTIONS

Recall the function `head` that yields the first element of a finite list. Clearly, `head x` is undefined for $x = []$ because the empty list has no elements at all. As expected, the HASKELL output for `head []` is just “panic”:

```
*Main> head []
*** Exception: Prelude.head: empty list
*Main>
```

Functions such as `head` are called *partial functions* because they cannot be applied to all of their (well-typed) inputs, *i.e.*, they diverge for some of such inputs. Partial functions are very common in mathematics or programming — for other examples think of *e.g.* `tail`, and so on.

Panic is very dangerous in programming. In order to avoid this kind of behaviour one has two alternatives, either (a) ensuring that every call to `head x` is *protected* — *i.e.*, the contexts which wrap up such calls ensure *pre-condition* $x \neq []$, or (b) *raising* exceptions, *i.e.* explicit error values, as above. In the former case, mathematical proofs need to be carried out in order to ensure *safety* (that is, *pre-condition* compliance). The overall effect is that of *restricting* the domain of the partial function. In the latter case one goes the other way round, by *extending* the co-domain (vulg. range) of the function so that it accommodates exceptional outputs. In this way one might define, in HASKELL:

```
data ExtVal a = Ok a | Error
```

and then define the “extended” version of head:

```

exthead :: [a] → ExtVal a
exthead [] = Error
exthead x = Ok (head x)
    
```

Note that *ExtVal* is a *parametric* type which extends an arbitrary data type *a* with its (polymorphic) exception (or error value). It turns out that, in HASKELL, *ExtVal* is redundant because such a parametric type already exists and is called *Maybe*:

```

data Maybe a = Nothing | Just a
    
```

Clearly, the isomorphisms hold:

$$\text{ExtVal } A \cong \text{Maybe } A \cong 1 + A$$

So, we might have written the more standard code

```

exthead :: [a] → Maybe a
exthead [] = Nothing
exthead x = Just (head x)
    
```

In abstract terms, both alternatives coincide, since one may regard as *partial* every function of type

$$1 + A \xleftarrow{g} B$$

for some *A* and *B*¹.

4.2 PUTTING PARTIAL FUNCTIONS TOGETHER

Do partial functions compose? Their types won’t match in general:

$$\begin{array}{c}
 1 + B \xleftarrow{g} A \\
 \vdots \\
 1 + C \xleftarrow{f} B
 \end{array}$$

Clearly, we have to extend *f* — which is itself a partial function — to some *f'* able to accept arguments from $1 + B$:

$$\begin{array}{ccc}
 & 1 & \\
 & \downarrow i_1 & \\
 \dots & 1 + B & \xleftarrow{g} A \\
 & \uparrow i_2 & \\
 1 + C & \xleftarrow{f} B & \\
 & \nwarrow f' &
 \end{array}$$

¹ In conventional programming, every function delivering a *pointer* as result — as in e.g. the C programming language — can be regarded as one of these functions.

The most “obvious” instance of the ellipsis (...) in the diagram above is i_1 and this corresponds to what is called *strict* composition: an exception produced by the *producer* function g is propagated to the output of the *consumer* function f . We define:

$$f \bullet g \stackrel{\text{def}}{=} [i_1, f] \cdot g \quad (4.1)$$

Expressed in terms of *Maybe*, composite function $f \bullet g$ works as follows:

$$(f \bullet g)a = f'(ga)$$

where

$$\begin{aligned} f' \text{ Nothing} &= \text{Nothing} \\ f' (\text{Just } b) &= f b \end{aligned}$$

Altogether, we have the following Haskell pointwise expression for $f \bullet g$:

$$\begin{aligned} \lambda a \rightarrow f' (g a) \textbf{ where} \\ f' \text{ Nothing} &= \text{Nothing} \\ f' (\text{Just } b) &= f b \end{aligned}$$

Note that the adopted extension of f can be decomposed — by reverse $+$ -absorption (2.43) — into

$$f' = [i_1, id] \cdot (id + f)$$

as displayed in diagram

$$\begin{array}{ccc} 1 + (1 + C) & \xleftarrow{id+f} & 1 + B \xleftarrow{g} A \\ \downarrow [i_1, id] & & \vdots \\ 1 + C & \xleftarrow{f} & B \end{array}$$

All in all, we have the following version of (4.1):

$$f \bullet g \stackrel{\text{def}}{=} [i_1, id] \cdot (id + f) \cdot g$$

Does this functional composition scheme have a unit, that is, is there a function u such that

$$f \bullet u = f = u \bullet f \quad (4.2)$$

holds? Clearly, if it exists, it must bear type $1 + A \xleftarrow{u} A$. $1 + A \xleftarrow{i_2} A$ has the same type, so $u = i_2$ is a very likely solution. Let us check it:

$$\begin{aligned} f \bullet u &= f = u \bullet f \\ &\equiv \{ \text{substitutions} \} \\ &[i_1, f] \cdot i_2 = f = [i_1, i_2] \cdot f \\ &\equiv \{ \text{by } +- \text{cancellation (2.40) and } +- \text{reflection (2.41)} \} \\ &f = f = id \cdot f \\ &\equiv \{ \text{trivial} \} \\ &\text{true} \end{aligned}$$

So $f \bullet u = f = u \bullet f$ for $u = i_2$.

Exercise 4.1. Prove that property

$$f \bullet (g \bullet h) = (f \bullet g) \bullet h$$

holds, for $f \bullet g$ defined by (4.1).

□

4.3 LISTS

In contrast to partial functions, which may produce *no* output, let us now consider functions which may deliver *too many* outputs, for instance, lists of output values:

$$\begin{array}{ccc} & B^* & \xleftarrow{g} A \\ & \vdots & \\ C^* & \xleftarrow{f} B & \end{array}$$

Functions f and g do not compose but, once again, one can think of extending the consumer function (f) by mapping it along the output of the producer function (g):

$$\begin{array}{ccc} (C^*)^* & \xleftarrow{f^*} B^* & \\ \vdots & \vdots & \\ C^* & \xleftarrow{f} B & \end{array}$$

To complete the process, one has to *flatten* the nested-sequence output in $(C^*)^*$ via the obvious list-catamorphism $C^* \xleftarrow{\text{concat}} (C^*)^*$, recall $\text{concat} = (\llbracket \llbracket _, \text{conc} \rrbracket \rrbracket)$ where $\text{conc}(x, y) = x ++ y$. In summary:

$$f \bullet g \stackrel{\text{def}}{=} \text{concat} \cdot f^* \cdot g \tag{4.3}$$

as captured in the following diagram:

$$\begin{array}{ccccc} & (C^*)^* & \xleftarrow{f^*} & B^* & \xleftarrow{g} A \\ & \downarrow \text{concat} & & \vdots & \\ C^* & & \xleftarrow{f} & B & \end{array}$$

Exercise 4.2. Show that *singl* (recall exercise 3.14) is the unit u of \bullet as defined by (4.3) above.

□

Exercise 4.3. Encode in HASKELL a pointwise version of (4.3). **Hint:** start by applying (list) cata-absorption (3.77).

☐

4.4 MONADS

Both function composition schemes (4.1) and (4.3) above share the same polytypic pattern: the output of the producer function g is “*T-times*” more elaborate than the input of the consumer function f , where T is some parametric datatype: $T\ X = 1 + X$ in case of (4.1), and $T\ X = X^*$ in case of (4.3). Then a composition scheme is devised for such functions, which is displayed in

$$\begin{array}{ccccc}
T(TC) & \xleftarrow{Tf} & TB & \xleftarrow{g} & A \\
\mu \downarrow & & \vdots & & \uparrow \\
TC & \xleftarrow{f} & B & & \\
& & & \nearrow f \bullet g &
\end{array} \tag{4.4}$$

and is given by

$$f \bullet g \stackrel{\text{def}}{=} \mu \cdot \top f \cdot g \quad (4.5)$$

where $\mathsf{T} A \xleftarrow{\mu} \mathsf{T}^2 A$ is a suitable polymorphic function. (We have already seen $\mu = [i_1, id]$ in case (4.1), and $\mu = \text{concat}$ in case (4.3).)

Together with a unit function $\top A \xleftarrow{u} A$ and μ , that is

$$A \xrightarrow{u} \mathsf{T} A \xleftarrow{\mu} \mathsf{T}^2 A$$

datatype T will form a so-called *monad* type, of which $(1+)$ and $(-)^*$ are the two examples seen above. Arrow $\mu \cdot T f$ is called the *extension* of f . Functions μ and u are referred to as the monad's *multiplication* and *unit*, respectively. The monadic composition scheme (4.5) is referred to as *Kleisli composition*.

A *monadic arrow* $TB \xleftarrow{f} A$ conveys the idea of a function which produces an output of “type” B “wrapped by T ”, where datatype T describes some kind of (computational) “effect”. The monad’s unit $TB \xleftarrow{u} B$ is a primitive monadic arrow which injects (*i.e.* promotes, wraps) data *inside* such an effect.

The monad concept is nowadays of primary importance in computing science because it makes it possible to describe computational effects as disparate as input/output, state variable updating, context

dependence, partial behaviour (seen above) *etc.* in an elegant, generic and uniform way. Moreover, the monad's operators exhibit notable properties which make it possible to *reason* about such computational effects.

The remainder of this section is devoted to such properties. First of all, the properties implicit in the following diagrams will be *required* for T to be regarded as a monad:

Multiplication :

$$\begin{array}{ccc}
 T^2 A & \xleftarrow{\mu} & T^3 A \\
 \mu \downarrow & & \downarrow T\mu \\
 T A & \xleftarrow{\mu} & T^2 A
 \end{array}
 \qquad \mu \cdot \mu = \mu \cdot T\mu \qquad (4.6)$$

Unit :

$$\begin{array}{ccc}
 T^2 A & \xleftarrow{u} & T A \\
 \mu \downarrow & \swarrow id & \downarrow T u \\
 T A & \xleftarrow{\mu} & T^2 A
 \end{array}
 \qquad \mu \cdot u = \mu \cdot T u = id \qquad (4.7)$$

The simple and beautiful symmetries apparent in these diagrams will make it easy to memorize their laws and check them for particular cases. For instance, for the $(1+)$ monad, law (4.7) will read as follows:

$$[i_1, id] \cdot i_2 = [i_1, id] \cdot (id + i_2) = id$$

These equalities are easy to check.

In laws (4.6) and (4.7), the different instances of μ and u are differently typed, as these are polymorphic and exhibit natural properties:

μ -natural :

$$\begin{array}{ccc}
 A & & T A \xleftarrow{\mu} T^2 A \\
 f \downarrow & & \downarrow T f \\
 B & & T B \xleftarrow{\mu} T^2 B
 \end{array}
 \qquad T f \cdot \mu = \mu \cdot T^2 f \qquad (4.8)$$

u -natural :

$$\begin{array}{ccc}
 A & & T A \xleftarrow{u} A \\
 f \downarrow & & \downarrow f \\
 B & & T B \xleftarrow{u} B
 \end{array}
 \qquad T f \cdot u = u \cdot f \qquad (4.9)$$

The simplest of all monads is the *identity monad* $T X \stackrel{\text{def}}{=} X$, which is such that $\mu = id$, $u = id$ and $f \bullet g = f \cdot g$. So — in a sense — the *whole functional discipline* studied thus far was already *monadic*, living inside

the simplest of all monads: the identity one. Put in other words, such functional discipline can be framed into a wider discipline in which an arbitrary monad is present. Describing this is the main aim of the current chapter.

PROPERTIES INVOLVING (KLEISLI) COMPOSITION The following properties arise from the definitions and monadic properties presented above:

$$f \bullet (g \bullet h) = (f \bullet g) \bullet h \quad (4.10)$$

$$u \bullet f = f = f \bullet u \quad (4.11)$$

$$(f \bullet g) \cdot h = f \bullet (g \cdot h) \quad (4.12)$$

$$(f \cdot g) \bullet h = f \bullet (\mathbb{T} g \cdot h) \quad (4.13)$$

$$id \bullet id = \mu \quad (4.14)$$

Properties (4.10) and (4.11) are the monadic counterparts of, respectively, (2.8) and (2.10), meaning that monadic composition preserves the properties of normal functional composition. In fact, for the identity monad, these properties coincide with each other.

Above we have shown that property (4.11) holds for the list monad, recall (4.2). A general proof can be produced similarly. We select property (4.10) as an illustration of the rôle of the monadic properties:

$$\begin{aligned} & f \bullet (g \bullet h) \\ = & \{ \text{definition (4.5) twice} \} \\ & \mu \cdot \mathbb{T} f \cdot (\mu \cdot \mathbb{T} g \cdot h) \\ = & \{ \mu \text{ is natural (4.8)} \} \\ & \mu \cdot \mu \cdot \mathbb{T}^2 f \cdot \mathbb{T} g \cdot h \\ = & \{ \text{monad property (4.6)} \} \\ & \mu \cdot \mathbb{T} \mu \cdot \mathbb{T}^2 f \cdot \mathbb{T} g \cdot h \\ = & \{ \text{functor } \mathbb{T} \text{ (3.56)} \} \\ & \mu \cdot \mathbb{T} (\mu \cdot \mathbb{T} f \cdot g) \cdot h \\ = & \{ \text{definition (4.5) twice} \} \\ & (f \bullet g) \bullet h \end{aligned}$$

Clearly, this calculation generalizes that of exercise 4.1 to any monad \mathbb{T} .

Exercise 4.4. Prove the other laws above and also the following one,

$$(\mathbb{T} f) \cdot (h \bullet k) = (\mathbb{T} f \cdot h) \bullet k \quad (4.15)$$

where Kleisli composition again trades with normal composition.

□

4.5 MONADIC APPLICATION (BINDING)

We have seen above that, given a monad $A \xrightarrow{u} \mathbb{T} A \xleftarrow{\mu} \mathbb{T}^2 A$, u is the unit of Kleisli composition, $f \bullet u = f$, recall (4.11). Now, what does happen in case we Kleisli compose f with the identity id of *standard* composition? Looking at diagram (4.4) for this case,

$$\begin{array}{ccccc} \mathbb{T}(\mathbb{T} C) & \xleftarrow{\mathbb{T} f} & \mathbb{T} B & \xleftarrow{id} & \mathbb{T} B \\ \mu \downarrow & & \vdots & & \\ \mathbb{T} C & \xleftarrow{f} & B & & \end{array}$$

we realize that $f \bullet id$ accepts a value of type $\mathbb{T} B$ that is passed to $\mathbb{T} C \xleftarrow{f} B$, yielding an output of type $\mathbb{T} C$. This construction is called *binding* and denoted by $\gg=f$:

$$(\gg=f) = f \bullet id \quad (4.16)$$

Expressed pointwise, we get:²

$$x \gg=f \stackrel{\text{def}}{=} (\mu \cdot \mathbb{T} f)x \quad (4.17)$$

This operator exhibits properties that arise from its definition and the basic monadic properties, *e.g.*

$$\begin{aligned} x \gg=u & \\ & \equiv \{ \text{definition (4.17)} \} \\ & (\mu \cdot \mathbb{T} u)x \\ & \equiv \{ \text{law (4.7)} \} \\ & (id)x \\ & \equiv \{ \text{identity function} \} \\ & x \end{aligned}$$

At pointwise level, one may chain monadic compositions from left to right, *e.g.*

$$(((x \gg=f_1) \gg=f_2) \gg=\dots f_{n-1}) \gg=f_n$$

for functions $A \xrightarrow{f_1} \mathbb{T} B_1$, $B_1 \xrightarrow{f_2} \mathbb{T} B_2$, ..., $B_{n-1} \xrightarrow{f_n} \mathbb{T} B_n$.

4.6 SEQUENCING AND THE **DO**-NOTATION

Recall from above that $x \gg=f$ is the monadic *generalization* of function application $f x$, since both coincide for the identity monad. Also recall that, for $f = \underline{y}$ (the “everywhere”- y constant function) one gets $\underline{y} x = y$.

² In the case of the identity monad one has: $x \gg=f = f x$. So, $\gg=$ can be regarded as denoting *monadic function application*.

What does the corresponding monadic generalization, $x \gg y$ mean? In the standard notation, this leads to another monadic operator,

$$x \gg y \stackrel{\text{def}}{=} x \gg y \quad (4.18)$$

of type

$$(\gg) : T A \rightarrow T B \rightarrow T B$$

called “sequencing”. For instance, within the finite-list monad, one has

$$[1, 2] \gg [3, 4] = (\text{concat} \cdot [3, 4]^*)([1, 2]) = \text{concat}[[3, 4], [3, 4]] = [3, 4, 3, 4]$$

Because this operator is associative (prove this as an exercise), one may iterate it to more than two arguments and write, for instance,

$$x_1 \gg x_2 \gg \dots \gg x_n$$

This leads to the popular “**do**-notation”, which is another piece of (pointwise) notation which makes sense in a monadic context:

$$\mathbf{do} \{x_1; x_2; \dots; x_n\} \stackrel{\text{def}}{=} x_1 \gg \mathbf{do} \{x_2; \dots; x_n\}$$

for $n \geq 1$. For $n = 1$ one trivially has

$$\mathbf{do} x_1 \stackrel{\text{def}}{=} x_1$$

4.7 GENERATORS AND COMPREHENSIONS

The monadic **do**-notation paves the way to a device that is very useful in (pointwise) monadic programming. As before, we consider its (non-monadic) counterpart first. Consider for instance the expression $x + \text{sum } y$, where sum is some operator in some context, e.g. adding up all elements of a list. Nothing impedes us from “structuring” expression $x + \text{sum } y$ in the following way:

```
let  $a = \text{sum } y$ 
in  $x + a$ 
```

It turns out that the above is the same as the following monadic expression,

```
do {
   $a \leftarrow \text{sum } y;$ 
   $u (x + a)$  }
```

provided the underlying monad is the *identity* monad. Now, what does the notation $a \leftarrow \dots$ mean for an arbitrary monad $A \xrightarrow{u} T A \xleftarrow{\mu} T^2 A$?

The **do**-notation accepts a variant in which the arguments of the \gg operator are “generators” of the form

$$a \leftarrow x \quad (4.19)$$

where, for a of type A , x is an inhabitant of monadic type $\mathbb{T} A$. One may regard $a \leftarrow x$ as meaning “let a be taken from x ”. Then the **do**-notation unfolds as follows:

$$\mathbf{do} \ a \leftarrow x_1; x_2; \dots; x_n \stackrel{\text{def}}{=} x_1 \gg \lambda a. (\mathbf{do} \ x_2; \dots; x_n) \quad (4.20)$$

Of course, we should now allow for the x_i to range over terms involving variable a . For instance, by writing (again in the list-monad)

$$\mathbf{do} \ a \leftarrow [1, 2, 3]; [a^2] \quad (4.21)$$

we mean

$$\begin{aligned} & [1, 2, 3] \gg \lambda a. [a^2] \\ &= \text{concat}((\lambda a. [a^2])^* [1, 2, 3]) \\ &= \text{concat}[[1], [4], [9]] \\ &= [1, 4, 9] \end{aligned}$$

The analogy with classical set-theoretic ZF-notation, whereby one might write $\{a^2 \mid a \in \{1, 2, 3\}\}$ to describe the set of the first three perfect squares, calls for the following notation,

$$[a^2 \mid a \leftarrow [1, 2, 3]] \quad (4.22)$$

as a “shorthand” of (4.21). This is an instance of the so-called *comprehension* notation, which can be defined in general as follows:

$$[e \mid a_1 \leftarrow x_1, \dots, a_n \leftarrow x_n] = \mathbf{do} \ \{a_1 \leftarrow x_1; \dots; a_n \leftarrow x_n; u(e)\} \quad (4.23)$$

where u is the monad’s unit (4.7, 4.9).

Alternatively, comprehensions can be defined as follows, where p, q stand for arbitrary generators:

$$[t] = u \ t \quad (4.24)$$

$$[f \ x \mid x \leftarrow l] = (\mathbb{T} f) l \quad (4.25)$$

$$[t \mid p, q] = \mu [[t \mid q] \mid p] \quad (4.26)$$

Note, however, that comprehensions are not restricted to lists or sets — they can be defined for any monad \mathbb{T} thanks to the **do**-notation.

Exercise 4.5. Show that

$$(f \bullet g) a = \mathbf{do} \ \{b \leftarrow g \ a; f \ b\} \quad (4.27)$$

$$\mathbb{T} f \ x = \mathbf{do} \ \{a \leftarrow x; u \ (f \ x)\} \quad (4.28)$$

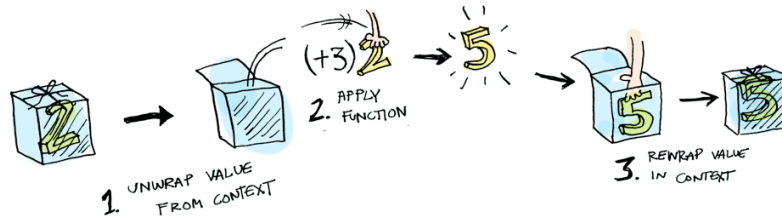
Note that the second **do** expression is equivalent to $x \gg (u \cdot f)$.

□

Exercise 4.6. Show that $x \gg= f = \mathbf{do} \{ a \leftarrow x; f a \}$ and then that $(x \gg= g) \gg= f$ is the same as $x \gg= f \bullet g$.

□

Fact (4.28) is illustrated in the cartoon³



for the computation of $\mathbf{T} (+3) x$, where $x = u\ 2$ is the \mathbf{T} -monadic object containing number 2.

4.8 MONADS IN HASKELL

In the *Standard Prelude* for HASKELL, one finds the following minimal definition of the *Monad* class,

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

where `return` refers to the unit of m , on top of which the “sequence” operator

```
(>>) :: m a -> m b -> m b
fail  :: String -> m a
```

is defined by

$$p \gg q = p \gg= \lambda _ \rightarrow q$$

as expected. This class is instantiated for finite sequences (`[]`), *Maybe* and *IO*, among others.

The μ multiplication operator is function `join` in module `Monad.hs`:

```
join :: (Monad m) => m (m a) -> m a
join x = x >>= id
```

³ Credits: see this and other helpful, artistic illustrations in http://adit.io/posts/2013-04-17-functors,_applicatives,_and_monads_in_pictures.html.

This is easily justified:

$$\begin{aligned}
 \text{join } x &= x \gg= id & (4.29) \\
 &= \{ \text{definition (4.17)} \} \\
 &\quad (\mu \cdot \top id)x \\
 &= \{ \text{functors commute with identity (3.55)} \} \\
 &\quad (\mu \cdot id)x \\
 &= \{ \text{law (2.10)} \} \\
 &\quad \mu x
 \end{aligned}$$

The following infix notation for (Kleisli) monadic composition in HASKELL uses the binding operator:

$$\begin{aligned}
 (\bullet) &:: \text{Monad } t \Rightarrow (b \rightarrow t c) \rightarrow (d \rightarrow t b) \rightarrow d \rightarrow t c \\
 (f \bullet g) a &= (g a) \gg= f
 \end{aligned}$$

Exercise 4.7. Consider the HASKELL function

$$\begin{aligned}
 \text{discollect} &:: [(a, [b])] \rightarrow [(a, b)] \\
 \text{discollect } [] &= [] \\
 \text{discollect } ((a, x) : y) &= [(a, b) \mid b \leftarrow x] ++ \text{discollect } y
 \end{aligned}$$

Knowing that finite lists form a monad where $\mu = \text{concat} = ([\text{nil}, \text{conc}])$ and $\text{conc } (x, y) = x ++ y$, derive the above pointfree code from the definition

$$\text{discollect} = \text{lstr} \bullet id \quad (4.30)$$

where $\text{lstr } (a, x) = [(a, b) \mid b \leftarrow x]$.

□

MONADIC I/O IO, a parametric datatype whose inhabitants are special values called *actions* or *commands*, is a most relevant monad. Actions perform the interconnection between HASKELL and the environment (file system, operating system). For instance,

$\text{getLine} :: \text{IO String}$

is a particular such action. Parameter *String* refers to the fact that this action “delivers” — or extracts — a string from the environment. This meaning is clearly conveyed by the type *String* assigned to symbol *l* in

$\text{do } l \leftarrow \text{getLine}; \dots l \dots$

which is consistent with typing rule for generators (4.19). Sequencing corresponds to the “;” syntax in most programming languages (e.g. C) and the **do**-notation is particularly intuitive in the IO-context.

Examples of functions delivering actions are

$$\text{FilePath} \xrightarrow{\text{readFile}} \text{IO String}$$

and

$$\text{Char} \xrightarrow{\text{putChar}} \text{IO } ()$$

— both produce I/O commands as result.

As is to be expected, the implementation of the IO monad in HASKELL

— available from the *Standard Prelude* — is not totally visible, for it is bound to deal with the intricacies of the underlying machine:

instance Monad IO where

$(\gg=) = \text{primbindIO}$

$\text{return} = \text{primretIO}$

Rather interesting is the way IO is regarded as a functor:

$$\text{fmap } f \, x = x \gg= (\text{return} \cdot f)$$

This goes the other way round, the monadic structure “helping” in defining the functor structure, everything consistent with the underlying theory:

$$\begin{aligned} x \gg= (u \cdot f) &= (\mu \cdot \text{IO}(u \cdot f))x \\ &= \{ \text{functors commute with composition} \} \\ &\quad (\mu \cdot \text{IO } u \cdot \text{IO } f)x \\ &= \{ \text{law (4.7) for } T = \text{IO} \} \\ &\quad (\text{IO } f)x \\ &= \{ \text{definition of } \text{fmap} \} \\ &\quad (\text{fmap } f) \, x \end{aligned}$$

For enjoyable reading on monadic input/output in HASKELL see [26], chapter 18.

Exercise 4.8. Extend the Maybe monad to the following “error message” exception handling datatype:

data Error a = Err String | Ok a **deriving** Show

In case of several error messages issued in a `do` sequence, how many turn up on the screen? Which ones?

□

Exercise 4.9. Recalling section 3.13, show that any inductive type with base functor

$$B(f, g) = f + F \, g$$

where F is an arbitrary functor, forms a monad for

$$\begin{aligned}\mu &= ([id, in \cdot i_2]) \\ u &= in \cdot i_1.\end{aligned}$$

Identify F for known monads such as eg. Maybe, LTree and (non-empty) lists.

□

4.9 THE STATE MONAD

The so-called *state monad* is a monad whose inhabitants are state-transitions encoding a particular brand of state-based automata known as *Mealy machines*. Given a set A (input alphabet), a set B (output alphabet) and a set of states S , a deterministic Mealy machine (DMM) is specified by a transition function of type

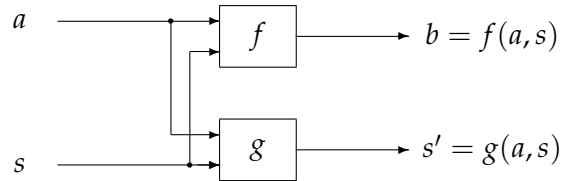
$$A \times S \xrightarrow{\delta} B \times S \quad (4.31)$$

Wherever $(b, s') = \delta(a, s)$, we say that the machine has transition

$$s \xrightarrow{a|b} s'$$

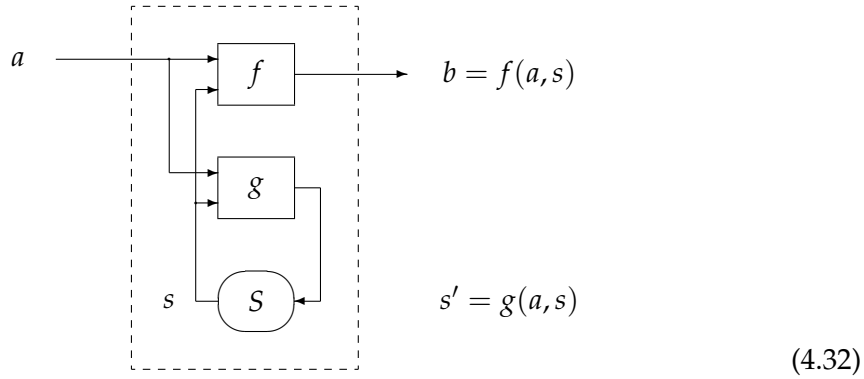
and refer to s as the *before* state, and to s' as the *after* state. Many programs that one writes in conventional programming languages such as C or Java can be regarded as DMMs.

It is clear from (4.31) that δ can be expressed as the *split* of two functions f and g — $\delta = \langle f, g \rangle$ — as depicted in the following drawing:



Note, however, that the information recorded in the state of a DMM is either meaningless to the user of the machine (as in eg. the case of states represented by numbers) or too complex to be perceived and handled explicitly (as is the case of eg. the data kept in a large database). So, it is convenient to *abstract* from it, via the “encapsulation” sug-

gested by the following, transformed, version of the previous drawing,



in which the state is no longer accessible from the outside.

Such an abstraction is nicely captured by the so-called *state monad*, in the following way: taking (4.31) and recalling (2.91), we simply transpose (ie. *curry*) δ and obtain

$$A \xrightarrow{\bar{\delta}} \underbrace{(B \times S)^S}_{(\text{St } S) B} \quad (4.33)$$

thus “shifting” the *input* state to the *output*. In this way, $\bar{\delta} a$ is a function capturing all state-transitions (and corresponding outputs) for input a . For instance, the function that *appends* a new element at the rear of a queue,

$$\text{enq}(a, s) \stackrel{\text{def}}{=} s ++ [a]$$

can be converted into a DMM by adding to it a dummy output of type 1 and then transposing:

$$\begin{aligned} \text{enqueue} &: A \rightarrow (1 \times S)^S \\ \text{enqueue } a &\stackrel{\text{def}}{=} \langle !, (++ [a]) \rangle \end{aligned} \quad (4.34)$$

Action *enqueue* performs *enq* on the state while acknowledging it by issuing an output of type 1.⁴

UNIT AND MULTIPLICATION. Let us now show that

$$(\text{St } S) A \cong (A \times S)^S \quad (4.35)$$

forms a monad. As we shall see, the fact that the *values* of this monad are functions brings the theory of exponentiation to the forefront. Thus, a review of section 2.15 is recommended at this point.

Notation \hat{f} will be used to abbreviate *uncurry* f , enabling the following variant of universal law (2.83),

$$\hat{k} = f \Leftrightarrow f = ap \cdot (k \times id) \quad (4.36)$$

⁴ A kind of “done!” message.

whose cancellation

$$\widehat{k} = ap \cdot (k \times id) \quad (4.37)$$

is written pointwise as follows:

$$\widehat{k}(c, a) = (k \ c)a \quad (4.38)$$

First of all, what is the functor behind (4.35)? Fixing the state space S , we obtain

$$\mathbb{T}X \stackrel{\text{def}}{=} (X \times S)^S \quad (4.39)$$

on objects and

$$\mathbb{T}f \stackrel{\text{def}}{=} (f \times id)^S \quad (4.40)$$

on functions, where $(_)^S$ is the exponential functor (2.87).

The unit of this monad is the transpose of the simplest of all Mealy machines — the identity:

$$\begin{aligned} u &: A \rightarrow (A \times S)^S \\ u &= \overline{id} \end{aligned} \quad (4.41)$$

Let us see what this means:

$$\begin{aligned} u &= \overline{id} \\ &\equiv \{ (2.83) \} \\ &ap \cdot (u \times id) = id \\ &\equiv \{ \text{introducing variables} \} \\ &ap(u \ a, s) = (a, s) \\ &\equiv \{ \text{definition of } ap \} \\ &(u \ a)s = (a, s) \end{aligned}$$

So, action $u \ a$ performed on state s keeps s unchanged and outputs a .

From the type of μ , for this monad,

$$((A \times S)^S \times S)^S \xrightarrow{\mu} (A \times S)^S$$

one figures out $\mu = x^S$ (recalling the exponential functor as defined by (2.87)) for some $((A \times S)^S \times S) \xrightarrow{x} (A \times S)$. This, on its turn, is easily recognized as an instance of the ap polymorphic function (2.83), which is such that $ap = \widehat{id}$, recall (2.85). Altogether, we define

$$\mu = ap^S \quad (4.42)$$

Let us inspect the behaviour of μ by checking the meaning of applying it to an action expressed as in diagram (2.91):

$$\begin{aligned}
& \mu\langle f, g \rangle = ap^S\langle f, g \rangle \\
\equiv & \quad \{ (2.87) \} \\
& \mu\langle f, g \rangle = ap \cdot \langle f, g \rangle \\
\equiv & \quad \{ \text{extensional equality (2.5)} \} \\
& \mu\langle f, g \rangle s = ap(f\ s, g\ s) \\
\equiv & \quad \{ \text{definition of } ap \} \\
& \mu\langle f, g \rangle s = (f\ s)(g\ s)
\end{aligned}$$

We find out that μ “unnests” the action inside f by applying it to the state delivered by g .

CHECKING THE MONADIC LAWS. The calculation of (4.7) is made in two parts, checking $\mu \cdot u = id$ first,

$$\begin{aligned}
& \mu \cdot u \\
= & \quad \{ \text{definitions} \} \\
& ap^S \cdot \overline{id} \\
= & \quad \{ \text{exponentials absorption (2.88)} \} \\
& \overline{ap \cdot id} \\
= & \quad \{ \text{reflection (2.85)} \} \\
& id \\
& \square
\end{aligned}$$

and then checking $\mu \cdot (\top u) = id$:

$$\begin{aligned}
& \mu \cdot (\top u) \\
= & \quad \{ (4.42, 4.40) \} \\
& ap^S \cdot (\overline{id} \times id)^S \\
= & \quad \{ \text{functor} \} \\
& (ap \cdot (\overline{id} \times id))^S \\
= & \quad \{ \text{cancellation (2.84)} \} \\
& id^S \\
= & \quad \{ \text{functor} \} \\
& id \\
& \square
\end{aligned}$$

The proof of (4.6) is also not difficult once supported by the laws of exponentials.

KLEISLI COMPOSITION. Let us calculate $f \bullet g$ for this monad:

$$\begin{aligned}
& f \bullet g \\
= & \{ (4.5) \} \\
& \mu \cdot \top f \cdot g \\
= & \{ (4.42); (4.40) \} \\
& ap^S \cdot (f \times id)^S \cdot g \\
= & \{ (-)^S \text{ is a functor} \} \\
& (ap \cdot (f \times id))^S \cdot g \\
= & \{ (4.36) \} \\
& \widehat{f}^S \cdot g \\
= & \{ \text{cancellation} \} \\
& \widehat{f}^S \cdot \widehat{g} \\
= & \{ \text{absorption (2.88)} \} \\
& \widehat{\widehat{f} \cdot \widehat{g}}
\end{aligned}$$

In summary, we have:

$$f \bullet g = \widehat{\widehat{f} \cdot \widehat{g}} \quad (4.43)$$

which can be written alternatively as

$$\widehat{f \bullet g} = \widehat{f} \cdot \widehat{g}$$

Let us use this in calculating law

$$pop \bullet push = u \quad (4.44)$$

where $push$ and pop are such that

$$\begin{aligned}
push & : A \rightarrow (1 \times S)^S \\
\widehat{push} & \stackrel{\text{def}}{=} \langle !, \widehat{(\cdot)} \rangle
\end{aligned} \quad (4.45)$$

$$\begin{aligned}
pop & : 1 \rightarrow (A \times S)^S \\
\widehat{pop} & \stackrel{\text{def}}{=} \langle head, tail \rangle \cdot \pi_2
\end{aligned} \quad (4.46)$$

for S the datatype of finite lists. We reason:

$$\begin{aligned}
& pop \bullet push \\
= & \{ (4.43) \} \\
& \widehat{\widehat{pop} \cdot \widehat{push}} \\
= & \{ (4.45, 4.46) \} \\
& \widehat{\langle head, tail \rangle \cdot \pi_2 \cdot \langle !, \widehat{(\cdot)} \rangle}
\end{aligned}$$

$$\begin{aligned}
&= \{ (2.22, 2.26) \} \\
&\quad \overline{\langle head, tail \rangle \cdot (\cdot)} \\
&= \{ out \cdot in = id \text{ (lists)} \} \\
&\quad \overline{id} \\
&= \{ (4.41) \} \\
&\quad u \\
&\square
\end{aligned}$$

BIND. The effect of binding a state transition x to a state-monadic function h is calculated in a similar way:

$$\begin{aligned}
&x \gg= h \\
&= \{ (4.17) \} \\
&\quad (\mu \cdot Th)x \\
&= \{ (4.42) \text{ and } (4.40) \} \\
&\quad (ap^S \cdot (h \times id)^S)x \\
&= \{ (-)^S \text{ is a functor} \} \\
&\quad (ap \cdot (h \times id))^S x \\
&= \{ \text{cancellation (4.37)} \} \\
&\quad \widehat{h}^S x \\
&= \{ \text{exponential functor (2.87)} \} \\
&\quad \widehat{h} \cdot x
\end{aligned}$$

Let us unfold $\widehat{h} \cdot x$ by splitting x into its components two components f and g :

$$\begin{aligned}
&\langle f, g \rangle \gg= h = \widehat{h} \cdot \langle f, g \rangle \\
&\equiv \{ \text{go pointwise} \} \\
&\quad (\langle f, g \rangle \gg= h)s = \widehat{h}(\langle f, g \rangle s) \\
&\equiv \{ (2.20) \} \\
&\quad (\langle f, g \rangle \gg= h)s = \widehat{h}(f\ s, g\ s) \\
&\equiv \{ (4.38) \} \\
&\quad (\langle f, g \rangle \gg= h)s = h(f\ s)(g\ s)
\end{aligned}$$

In summary, for a given “before state” s , $g\ s$ is the intermediate state upon which $f\ s$ runs and yields the output and (final) “after state”.

TWO PROTOTYPICAL INHABITANTS OF THE STATE MONAD: *get* AND *put*. These generic actions are defined as follows, in the PF-style:

$$get \stackrel{\text{def}}{=} \langle id, id \rangle \quad (4.47)$$

$$put \stackrel{\text{def}}{=} \overline{\langle !, \pi_1 \rangle} \quad (4.48)$$

Action *g* retrieves the data stored in the state without changing it, while *put* stores a particular value in the state. Note that *put* can also be written

$$put\ s = \langle !, \underline{s} \rangle \quad (4.49)$$

or even as

$$put\ s = update\ \underline{s} \quad (4.50)$$

where

$$update\ f \stackrel{\text{def}}{=} \langle !, f \rangle \quad (4.51)$$

updates the state via state-to-state function *f*.

The following is an example, written in Haskell, of the standard use of *get/put* in managing context data, in this case a counter. Function *decBTree* decorates each node of a *BTree* (recall this datatype from page 90) with its position in the tree:

```
decBTree Empty = return Empty
decBTree (Node (a, (t1, t2))) = do {
  n ← get;
  put (n + 1);
  x ← decBTree t1;
  y ← decBTree t2;
  return (Node ((a, n), (x, y)))
}
```

To close the chapter, we will present a strategy for deriving this kind of monadic functions.

4.10 ‘MONADIFICATION’ OF HASKELL CODE MADE EASY

There is an easy roadmap for “monadification” of Haskell code. What do we mean by *monadification*? Well, in a sense — as we shall soon see — every piece of code is monadic: we don’t notice this because the underlying monad is *invisible* (the *identity* monad). We are going to see how to make it visible taking advantage of monadic *do* notation and leaving it open for instantiation. This will bridge the gap between monads’ theory and its application to handling particular effects in concrete programming situations.

Let us take as starting point the pointwise version of *sum*, the list catamorphism that adds all numbers found in its input:

```
sum [] = 0
sum (h : t) = h + sum t
```

Notice that this code could have been written as follows

```
sum [] = id 0
sum (h : t) = let x = sum t in id (h + x)
```

using *let* notation and two instances of the identity function. Question: what is the point of such a “baroque” version of the starting, so simple piece of code? Answer:

- The *let ... in ...* notation stresses the fact that recursive call happens *earlier* than the delivery of the result.
- The *id* functions signal the exit points of the algorithm, that is, the points where it *returns* something to the caller.

Next, let us

- re-write *id* into *return*—;
- re-write *let x = ... in ...*— into *do { x <- ... ; ... }*

One will obtain the following version of *sum*:

```
msum [] = return 0
msum (h : t) = do { x <- msum t; return (h + x) }
```

Typewise, while *sum* has type $(\text{Num } a) \Rightarrow [a] \rightarrow a$, *msum* has type

$$(\text{Monad } m, \text{Num } a) \Rightarrow [a] \rightarrow m a$$

That is, *msum* is monadic — parametric on monad *m* — while *sum* is not.

There is a particular monad for which *sum* and *msum* coincide: the **identity** monad $\text{Id } X = X$. It is very easy to show that inside this monad *return* is the identity and *do* $x \leftarrow \dots$ means the same as *let* $x = \dots$, as already mentioned — enough for the pointwise versions of the two functions to be the same. Thus the “invisible” monad mentioned earlier is the identity monad.

In summary, the monadic version is *generic* in the sense that it runs on whatever monad you like, enabling you to perform *side effects* while the code runs. If you don’t need any effects then you get the “non-monadic” version as special case, as seen above. Otherwise, Haskell will automatically switch to the effects you want, depending on the monad you choose (often determined by context).

For each particular monad we may decide to add specific monadic code like *get* and *put* in the *decBTree* example, where we want to

take advantage of the state monad. As another example, check the following enrichment of *msum* with state-monadic code helping you to trace the execution of your program:

```
msum' [] = return 0
msum' (h:t) =
  do { x ← msum' t;
      print ("x= " ++ show x);
      return (h + x) }
```

Thus one obtains traces of the code in the way prescribed by the particular usage of the *print* (state monadic) function:

```
*Main> msum' [3,5,1,3,4]
"x= 0 "
"x= 4 "
"x= 7 "
"x= 8 "
"x= 13 "
*Main>
```

In the reverse direction, one may try and see what happens to monadic code upon removing all monad-specific functions and going into the identity monad once it gets monad generic. In the case of *decBTree*, for instance, we will get

```
decBTree Empty = return Empty
decBTree (Node (a, (t1, t2))) =
  do
    x ← decBTree t1;
    y ← decBTree t2;
    return (Node (a, (x, y)))
```

once *get* and *put* are removed (and therefore all instances of *n*), and then

```
decBTree Empty = Empty
decBTree (Node (a, (t1, t2))) =
  let
    x = decBTree t1
    y = decBTree t2
  in Node (a, (x, y))
```

This is the identity function on type *BTree*, recall the cata-reflection law (3.68). So, the *archetype* of (inspiration for) much monadic code is the most basic of all tree traversal functions — the identity⁵. The same could be said about imperative code of a particular class — the *recursive descent* one — much used in compiler construction, for instance.

⁵ We have seen the same kind of “inspiration” before in building type functors (3.76) which, for *f* = *id*, boil down to the identity.

Playing with effects

As it may seem from the previous examples, adding effects to produce monadic code is far from arbitrary. This can be further appreciated by defining the function that yields the smallest element of a list,

$$\begin{aligned} \text{getmin } [a] &= a \\ \text{getmin } (h:t) &= \min h (\text{getmin } t) \end{aligned}$$

which is incomplete in the sense that it does not specify the meaning of `getmin []`. As this is mathematically undefined, it should be expressed “outside the maths”, that is, as an effect. Thus, to complete the definition we first go monadic, as we did before,

$$\begin{aligned} \text{mgetmin } [a] &= \text{return } a \\ \text{mgetmin } (h:t) &= \text{do } \{x \leftarrow \text{mgetmin } t; \text{return } (\min h x)\} \end{aligned}$$

and then chose a monad in which to express the meaning of `getmin []`, for instance the `Maybe` monad

$$\begin{aligned} \text{mgetmin } [] &= \text{Nothing} \\ \text{mgetmin } [a] &= \text{return } a \\ \text{mgetmin } (h:t) &= \text{do } \{x \leftarrow \text{mgetmin } t; \text{return } (\min h x)\} \end{aligned}$$

Alternatively, we might have written

$$\text{mgetmin } [] = \text{Error "Empty input"}$$

going into the `Error` monad, or even the simpler (yet interesting) `mgetmin [] = []`, which shifts the code into the list monad, yielding singleton lists in the success case, otherwise the empty list.

Function `getmin` above is an example of a partial function, that is, a function which is undefined for some of its inputs.⁶ These functions cause much interference in functional programming, which monads help us to keep under control.

Let us see how such interference is coped with in the case of higher order functions, taking `map` as example

$$\begin{aligned} \text{map } f [] &= [] \\ \text{map } f (h:t) &= (f h) : \text{map } f t \end{aligned}$$

and supposing `f` is not a total function. How do we cope with erring evaluations of `f h`? As before, we first “letify” the code,

$$\begin{aligned} \text{map } f [] &= [] \\ \text{map } f (h:t) &= \text{let} \\ &\quad b = f h \\ &\quad x = \text{map } f t \text{ in } b : x \end{aligned}$$

⁶ Recall that function partiality was our motivation for studying monads right from the beginning of this chapter.

we go monadic in the usual way,

```
mmap f [] = return []
mmap f (h:t) = do { b ← f h; x ← mmap f t; return (b:x) }
```

and everything goes smoothly — as can be checked, the function thus built is of the expected (monadic) type:

$$\text{mmap} :: (\text{Monad } T) \Rightarrow (a \rightarrow T\ b) \rightarrow [a] \rightarrow T\ [b] \quad (4.52)$$

Run `mmap Just [1,2,3,4]`, for instance: you will obtain `Just [1,2,3,4]`. Now run `mmap print [1,2,3,4]`. You will see the items in the sequence printed sequentially.

One may wonder about the behaviour of the `mmap` for `f` the identity function: will we get an error? No, we get a well-typed function of type $[m\ a] \rightarrow m\ [a]$, for m a monad. We thus obtain the well-known monadic function sequence which evaluates each *action* in the input sequence, from left to right, collecting the results. For instance, applying this function to input sequence `[Just 1, Nothing, Just 2]` the output will be `Nothing`.

Exercise 4.10. Use the monadification technique to encode monadic function

```
filterM :: Monad m => (a -> m B) -> [a] -> m [a]
```

which generalizes the list-based `filter` function.

□

Exercise 4.11. “Reverse” the following monadic code into its non-monadic archetype:

```
f :: (Monad m) => (a -> m B) -> [a] -> m [a]
f p [] = return []
f p (h:t) = do {
  b ← p h;
  t' ← f p t;
  return (if b then h:t' else [])
}
```

Which function of the Haskell Prelude do you get by such reverse monadification?

□

4.11 MONADIC RECURSION

There is much more one could say about monadic recursive programming. In particular, one can express the code “monadification” strate-

gies of the previous section in terms of catamorphisms. As an example, recall (4.52):

$$\begin{array}{ccc}
 A & & A^* \xleftarrow{\text{in}_{A^*}} 1 + A \times A^* \\
 f \downarrow & & \text{mmap } f \downarrow \quad \quad \downarrow \text{id} + \text{id} \times \text{mmap } f \\
 \mathbb{T} B & & \mathbb{T} B^* \xleftarrow{g} 1 + A \times \mathbb{T} B^*
 \end{array}$$

How do we build g ? Clearly, the recipe given by (3.76) needs to be adapted:

$$\begin{array}{ccc}
 A & & A^* \xleftarrow{\text{in}_{A^*}} 1 + A \times A^* \\
 f \downarrow & & \text{mmap } f \downarrow \quad \quad \downarrow \text{id} + \text{id} \times \text{mmap } f \\
 \mathbb{T} B & & \mathbb{T} B^* \xleftarrow{g} 1 + A \times \mathbb{T} B^* \\
 & & \quad \quad \downarrow \text{id} + f \times \text{id} \\
 & & 1 + \mathbb{T} B \times \mathbb{T} B^*
 \end{array}$$

[return · nil, [cons]]

where

$$[f] (x, y) = \mathbf{do} \{ a \leftarrow x; b \leftarrow y; \text{return } (f(a, b)) \}$$

By defining

$$\begin{aligned}
 (|g|)^b &= (|[\text{return} \cdot f, [h]]|) \text{ where} \\
 f &= (g \cdot i_1) \\
 h &= (g \cdot i_2)
 \end{aligned}$$

we can write

$$\text{mmap } f = (|(\text{in} \cdot (\text{id} + f \times \text{id}))|)^b \quad (4.53)$$

where (recall) $\text{in} = [\text{nil}, \text{cons}]$.

Handling monadic recursion in full generality calls for technical ingredients called *commutative laws* which fall outside the current scope of this chapter.

4.12 WHERE DO MONADS COME FROM?

In the current context, a good way to find an answer this question is to recall the universal property of exponentials (2.83):

$$\begin{array}{ccc}
 B^A & & B^A \times A \xrightarrow{\text{ap}} B \\
 k = \bar{f} \uparrow & & \uparrow k \times \text{id} \\
 C & & C \times A
 \end{array}$$

f

Let us re-draw this diagram by unfolding $B^A \times A$ into the composition of two functors $G (F B)$ where $F X = X^A$ and $G X = X \times A$:

$$\begin{array}{ccc}
 F B & & G (F B) \xrightarrow{\text{ap}} B \\
 k = \bar{f} \uparrow & & \uparrow G k \\
 C & & G C
 \end{array}$$

f

(4.54)

As we already know, this establishes the (*curry/uncurry*) isomorphism

$$G C \rightarrow B \cong C \rightarrow F B \quad (4.55)$$

assuming F and G as defined above.

Note how (4.55) expresses a kind of “shunting rule” at type level: G s on the input side can be “shunted” to the output if replaced by F s. This is exactly what *curry* and *uncurry* do typewise. The corollaries of the universal property can also be expressed in terms of F and G :

- Reflection: $\overline{ap} = id$, that is, $ap = \widehat{id}$ – recall (2.85)
- Cancellation: $\widehat{id} \cdot G \bar{f} = f$ – recall (2.84)
- Fusion: $\bar{h} \cdot g = \overline{h \cdot G g}$ – recall (2.86)
- Absorption: $(F g) \cdot \bar{h} = \overline{g \cdot h}$ – recall (2.88)
- Naturality: $h \cdot \widehat{id} = \widehat{id} \cdot G (F h)$
- Functor: $F h = \overline{h \cdot ap}$
- Closed definitions: $\widehat{k} = \widehat{id} \cdot (G k)$ and $\bar{g} = (F g) \cdot \bar{id}$, the latter following from absorption.

Now observe what happens if the functor composition $G \cdot F$ is swapped: $F (G X) = (X \times A)^A$. We get the *state monad* out of this construction,

$$(G \cdot F) X = (X \times A)^A = St A X$$

— recall (4.35). Interestingly, the universal property (4.54) can be expressed also in terms of such a monad structure, as the simple calculation shows:

$$\begin{aligned} k = \bar{f} &\Leftrightarrow ap \cdot G k = f \\ &\equiv \{ \text{see above} \} \\ k &= (F f) \cdot \bar{id} \Leftrightarrow f = \widehat{k} \\ &\equiv \{ \text{swapping variables } k \text{ and } f, \text{ to match the starting diagram} \} \\ f &= (F k) \cdot \bar{id} \Leftrightarrow k = \widehat{f} \end{aligned}$$

That is,

$$k = \widehat{f} \Leftrightarrow f = \underbrace{F k \cdot \eta}_{\bar{k}} \quad \begin{array}{ccc} G B & & F (G B) \xleftarrow{\eta} B \\ k = \widehat{f} \downarrow & & \downarrow F k \quad \swarrow f \\ C & & F C \end{array}$$

for $\eta = \bar{id}$, the unit of the monad $T = F \cdot G$. To complete the definition of the T monad in this way, we recall (4.42)

$$\mu = F \widehat{id} \quad (4.56)$$

with type $(T \cdot T) X \xrightarrow{\mu} T X$, where $id : T X \rightarrow T X$.

Adjunctions

The reasoning we have made above for exponentials and the state monad generalizes for any other monad. In general, an isomorphism of shape (4.55) is called an *adjunction* of the two functors F and G , which are said to be *adjoint* to each other. One writes $G \dashv F$ and says that G is *left* adjoint and that F is *right* adjoint. Using notation $\lfloor k \rfloor$ and $\lceil k \rceil$ for the generic witnesses of the isomorphism we write

$$\begin{array}{ccc} & \xrightarrow{\lceil - \rceil} & \\ G C \rightarrow B & \cong & C \rightarrow F B \\ & \xleftarrow{\lfloor - \rfloor} & \end{array} \quad (4.57)$$

From this a monad $T = F \cdot G$ arises defined by $\eta = \lceil id \rceil$ and $\mu = F \lfloor id \rfloor$. Finally, from all this we can infer the generic version of $f \bullet g$,

$$f \bullet g = \lceil \lfloor f \rfloor \cdot \lfloor g \rfloor \rceil \quad (4.58)$$

by replaying the calculation which lead to (4.43):

$$\begin{aligned} & f \bullet g \\ = & \{ (4.5) \} \\ & \mu \cdot T f \cdot g \\ = & \{ T = F \cdot G; \mu = F \lfloor id \rfloor \} \\ & F \lfloor id \rfloor \cdot (F (G f)) \cdot g \\ = & \{ \text{functor } F \} \\ & F (\lfloor id \rfloor \cdot G f) \cdot g \\ = & \{ \text{cancellation: } \lfloor id \rfloor \cdot G f = \lfloor f \rfloor; g = \lceil \lfloor g \rfloor \rceil \} \\ & F \lfloor f \rfloor \cdot \lceil \lfloor g \rfloor \rceil \\ = & \{ \text{absorption: } (F g) \cdot \lceil h \rceil = \lceil g \cdot h \rceil \} \\ & \lceil \lfloor f \rfloor \cdot \lfloor g \rfloor \rceil \end{aligned}$$

Finally, let us see another example of a monad arising from one such adjunction (4.57). Recall exercise 2.26, on page 41, where pair / unpair witness an isomorphism similar to that of *curry/uncurry*, for pair $(f, g) = \langle f, g \rangle$ and unpair $k = (\pi_1 \cdot k, \pi_2 \cdot k)$. This can be cast into an adjunction as follows

$$\begin{aligned} k &= \text{pair } (f, g) \Leftrightarrow (\pi_1 \cdot k, \pi_2 \cdot k) = (f, g) \\ \equiv & \{ \text{see below} \} \\ k &= \text{pair } (f, g) \Leftrightarrow (\pi_1, \pi_2) \cdot (G k) = (f, g) \end{aligned}$$

where $G k = (k, k)$. Note the abuse of notation, on the righthand side, of extending function composition notation to composition of *pairs* of functions, defined in the expected way: $(f, g) \cdot (h, k) = (f \cdot h, g \cdot k)$.

Note that, for $f : A \rightarrow B$ and $g : C \rightarrow D$, the pair (f, g) has type $(A \rightarrow B) \times (C \rightarrow D)$. However, we shall abuse of notation again and declare the type $(f, g) : (A, C) \rightarrow (B, D)$.⁷ In the opposite direction, $F (f, g) = f \times g$:

$$\begin{array}{ccc}
 B \times A & (B \times A, B \times A) & \xrightarrow{(\pi_1, \pi_2)} (B, A) \\
 \uparrow k=\text{pair } (f, g) & \uparrow (k, k) & \nearrow (f, g) \\
 C & (C, C) &
 \end{array}$$

This is but another way of writing the universal property of products (2.63), since $(f, g) = (h, k) \Leftrightarrow f = h \wedge g = k$ and $\text{pair } (f, g) = \langle f, g \rangle$, recall exercise 2.26.

What is, then, the monad behind this *pairing* adjunction? It is the *pairing monad* $(F \cdot G) X = F (G X) = F (X, X) = X \times X$, where $\eta = \langle id, id \rangle$ and $\mu = \pi_1 \times \pi_2$. This monad allows us to work with pairs regarded as 2-dimensional *vectors* (y, x) . For instance, the **do**-expression

```
do { x ← (2,3); y ← (4,5); return (x + y) }
```

yields $(6, 8)$ as result in this monad — the *vectorial* sum of vectors $(2, 3)$ and $(4, 5)$. A simple encoding of this monad in Haskell is:

```
data P a = P (a, a) deriving Show
instance Functor P where
  fmap f (P (a, b)) = P (f a, f b)
instance Monad P where
  x >>= f = (μ · fmap f) x
  return a = P (a, a)
  μ :: P (P a) → P a
  μ (P (P (a, b), P (c, d))) = P (a, d)
```

Exercise 4.12. What is the vectorial operation expressed by the definition

```
op k v = do { x ← v; return (k × x) }
```

in the *pairing monad*?

□

4.13 BIBLIOGRAPHY NOTES

The use of monads in computer science started with Moggi [41], who had the idea that monads should supply the extra semantic information needed to implement the lambda-calculus theory. Haskell [31] is

⁷ Strictly speaking, we are not abusing notation but rather working on a new *category*, that is, another mathematical system where functions and objects always come in pairs. For more on categories see the standard textbook [34].

among the computer languages which make systematic use of monads for implementing effects and imperative constructs in a purely functional style.

Category theorists invented monads in the 1960's to concisely express certain aspects of universal algebra. Functional programmers invented list comprehensions in the 1970's to concisely express certain programs involving lists. Philip Wadler [60] made a great contribution to the field by showing that list comprehensions could be generalised to arbitrary monads and unify with imperative "do"-notation in case of the monad which explains imperative computations.

Monads are nowadays an essential feature of functional programming and are used in fields as diverse as language parsing [27], component-oriented programming [8], strategic programming [32], multimedia [26] and probabilistic programming [13]. Adjunctions play a major role in [24].

Part II

CALCULATING WITH RELATIONS

WHEN EVERYTHING BECOMES A RELATION

In the previous chapters, (recursive) functions were taken as a basis for expressing computations, exhibiting powerful laws for calculating programs in a functional programming style.

When writing such programs one of course follows some line of thought concerning *what* the programs should do. *What the program should do* is usually understood as the *specification* of the problem that motivates writing the program in the first place. Specifications can be quite complex in real life situations. In other situations, the complexity of the program that one writes is in strong contrast with the simplicity of the specification. Take the example of sorting, which can be specified as simply as:

Yield an ordered permutation of the input.

Where do you find, in this specification, the orientation (or inspiration) that will guide a programmer towards writing a bi-recursive program like *quicksort*?

The question is, then: are functions *enough* for one to calculate functional programs from given specifications? It is the experience in other fields of mathematics that sometimes it is easier to solve a problem of domain D if one generalizes from D to some wider domain D' . In the field of real numbers, for instance, most of trigonometric identities are easily derived (and memorized) from Euler's formula involving complex exponentials: $e^{ix} = \cos x + i(\sin x)$.

Similarly, it turns out that functional programs often become easier to calculate if one handles them in the wider mathematical domain of binary relations. At school one gets accustomed to the sentence *every function is a special case of a relation*. This chapter puts the usefulness of such a piece of common knowledge into practice.

5.1 FUNCTIONS ARE NOT ENOUGH

Consider the following fragment of a requirement posed by a (fictional) telecommunication company:

(...) *For each list of calls stored in the mobile phone (eg. numbers dialed, SMS messages, lost calls), the store operation should work in a way such that (a) the more recently a call is made the more accessible*

it is; (b) no number appears twice in a list; (c) only the last 10 entries in each list are stored.

A tentative, first implementation of the *store* operation could be

$$\begin{aligned} \text{store} &: \text{Call} \rightarrow \text{Call}^* \rightarrow \text{Call}^* \\ \text{store } c \ l &= c : l \end{aligned}$$

However, such a version of function *store* fails to preserve the *properties* required in the fragment above in case $\text{length } l = 10$, or $c \in \text{elems } l$, where elems yields the set of all elements of a finite list,

$$\text{elems} = ([\text{empty}, \text{join}]) \quad (5.1)$$

for $\text{empty } _ = \{ \}$ and $\text{join } (a, s) = \{ a \} \cup s$.

Clearly, the designer would have to *restrict* the application of *store* to input values c, l such that the given properties are preserved. This could be achieved by adding a so-called “*pre-condition*”:

$$\begin{aligned} \text{store} &: \text{Call} \rightarrow \text{Call}^* \rightarrow \text{Call}^* \\ \text{store } c \ l &= c : l \\ \text{pre length } l &< 10 \wedge \neg (c \in \text{elems } l) \end{aligned}$$

Such a pre-condition is a predicate telling the range of *acceptable* input values, to be read as a *warning* provided by the designer that the function will not meet the requirements outside such a range of input values.

Thus *store* becomes a *partial function*, that is, a function defined only for some of its inputs. Although this partiality can be regarded as a symptom that the requirements have been partly misunderstood, it turns out that *partial functions* are the rule rather than the exception in mathematics and computing. For example, in the numeric field, we know what $1/2$ means; what about $1/0$? Ruling out this case means that *division* is a partial function. In list processing, given a sequence s , what does $s !! i$ mean in case $i > \text{length } s$? — list indexing is another *partial operation* (as are *head*, *tail* and so on).

Partial functions are not new to readers of this text: in section 4.1, the *Maybe* monad was used to “totalize” partial functions. In this chapter we shall adopt another strategy to cope with partiality, and one that has extra merits: it will also cope with computational nondeterminacy and vagueness of software requirements.

It can be shown that the following evolution of *store*,

$$\text{store } c = (\text{take } 10) \cdot (c:) \cdot \text{filter } (c \neq) \quad (5.2)$$

meets all requirements above with no need for preconditions, the extra components *take 10* and *filter (c ≠)* being added to comply with requirements (c) and (b), respectively.

Implementation (5.2) alone should be regarded as example of how functional programs can be built compositionally in a requirement-driven fashion. It does not, however, give any guarantees that the

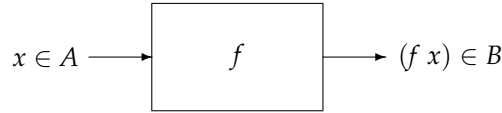
requirements are *indeed* met. How can we ensure this in the compositional way advocated in this book since its beginning? The main purpose of this chapter is to answer such a question.

5.2 FROM FUNCTIONS TO RELATIONS

The way functions are handled and expressed in standard maths books, e.g. in analysis and calculus,

$$y = f(x)$$

is indicative that, more important than the *reactive* behaviour of f ,



which was the starting point of section 2.1, mathematicians are more interested in expressing the input/output *relationship* of f , that is, the set of all pairs (y, x) such that $y = f x$. Such a set of pairs is often referred to as the “graph” of f , which can be plotted two-dimensionally in case types A and B are linearly ordered. (As is the standard case in which $A=B=\mathbb{R}$, the real numbers.)

It turns out that such a *graph* can be regarded as a special case of a *binary relation*. Take for instance the following functional declaration

$$\begin{cases} \text{succ} : \mathbb{N}_0 \rightarrow \mathbb{N}_0 \\ \text{succ } x = x + 1 \end{cases}$$

which expresses the computation rule of the *successor* function. Writing $y = \text{succ } n$ establishes the *binary relation* $y = x + 1$. This binary relation “coincides” with succ in the sense that writing

$$\begin{cases} \text{succ} : \mathbb{N}_0 \rightarrow \mathbb{N}_0 \\ y \text{ succ } x \Leftrightarrow y = x + 1 \end{cases}$$

means the same as the original definition, while making the i/o relationship explicit. Because there is *only one* y such that $y = x + 1$ we can safely drop both ys from $y \text{ succ } x \Leftrightarrow y = x + 1$, obtaining the original $\text{succ } x = x + 1$.

The new style is, however, more expressive, in the sense that it enables us to declare *genuine* binary relations, for instance

$$\begin{cases} R : \mathbb{N}_0 \rightarrow \mathbb{N}_0 \\ y R x \Leftrightarrow y \geq x + 1 \end{cases} \quad (5.3)$$

In this case, not only x and y such that $y = x + 1$ are admissible, but also $y = x + 2$, $y = x + 3$ and so on. It also enables us to express the *converse* of any function — an operation hitherto the privilege of isomorphisms only (2.16):

$$y f x \Leftrightarrow x f^\circ y \quad (5.4)$$

Converses of functions are very useful in problem solving, as we shall soon see. For instance, $\mathbb{N}_0 \xleftarrow{\text{succ}^\circ} \mathbb{N}_0$ denotes the *predecessor* relation in \mathbb{N}_0 . It is not a function because no y such that $y \text{ succ}^\circ 0$ exists — try and solve $0 = y + 1$ in \mathbb{N}_0 .

The intuitions above should suffice for us to start generalizing what we know about functions, from the preceding chapters, to binary relations. First of all, such relations are denoted by *arrows* exactly in the same way functions are. So,

we shall write $R : B \leftarrow A$, $R : A \rightarrow B$, $B \xleftarrow{R} A$ or $A \xrightarrow{R} B$
to indicate that relation R relates B -values to A -values.

That is, relations are typed in the same way as functions.

Given binary relation $R : B \leftarrow A$, writing $b R a$ (read: “ b is related to a by R ”) means the same as $a R^\circ b$, where R° is said to be the *converse* of R . In terms of grammar, R° corresponds to the *passive voice* — compare e.g.

$\underbrace{\text{John}}_b \underbrace{\text{loves}}_R \underbrace{\text{Mary}}_a$

with

$\text{Mary} \underbrace{\text{is loved by}}_{R^\circ} \text{John}$

That is, $(\text{loves})^\circ = (\text{is loved by})$. Another example:

Catherine eats the apple

— $R = (\text{eats})$, active voice — compared with

the apple is eaten by Catherine

— $R^\circ = (\text{is eaten by})$, passive voice.

Following a widespread convention, functions are denoted by lowercase characters (eg. f , g , ϕ) or identifiers starting with a lowercase characters, while uppercase letters are reserved to arbitrary relations. In the case of functions ($R := f$), $b f a$ means exactly $b = f a$. This is because functions are *univocal*, that is, no two different b and b' are such that $b f a \wedge b' f a$. In fact, the following facts hold about *any* function f :

- *Univocality* (or “left” uniqueness) —

$$b f a \wedge b' f a \Rightarrow b = b' \quad (5.5)$$

- *Leibniz principle* —

$$a = a' \Rightarrow f a = f a' \quad (5.6)$$

Clearly, not every relation obeys (5.5), for instance

$$2 < 3 \wedge 1 < 3 \not\Rightarrow 2 = 1$$

Relations obeying (5.5) will be referred to as *simple*, according to a terminology to follow shortly.

Implication (5.6) expresses the (philosophically) interesting fact that no function (observation) can be found able to distinguish between two equal objects. This is another fact true about functions which does not generalize to binary relations, as we shall see when we come back to this later.

Recapitulating: we regard function $f : A \rightarrow B$ as the binary relation which relates b to a iff $b = f a$. So,

$$b f a \text{ literally means } b = f a \quad (5.7)$$

The purpose of this chapter is to generalize from

$$\boxed{\begin{array}{c} B \xleftarrow{f} A \\ b = f a \end{array}} \quad \text{to} \quad \boxed{\begin{array}{c} B \xleftarrow{R} A \\ b R a \end{array}}$$

5.3 PRE/POST CONDITIONS

It should be noted that relations are used in virtually every body of science and it is hard to think of another way to express human knowledge in philosophy, epistemology and common life, as suggestively illustrated in figure 5.2. This figure is also illustrative of another popular ingredient when using relations — the *arrows* drawn to denote relationships.¹

In real life, “everything appears to be a relation”. This has lead software theorists to invent linguistic layouts for relational specification, leading to so-called *specification languages*. One such language, today historically relevant, is the language of the Vienna Development Method (VDM). In this notation, the relation described in (5.3) will be written:

$$\begin{array}{l} R (x : \mathbb{N}_0) y : \mathbb{N}_0 \\ \text{post } y \geq x + 1 \end{array}$$

where the clause prefixed by *post* is said to be a post-condition. The format also includes pre-conditions, if necessary. Such is the case of the following pre / post -styled specification of the operation that extracts an arbitrary element from a set:

$$\begin{array}{l} \text{Pick } (x : PA) (r : A, y : PA) \\ \text{pre } x \neq \{ \} \\ \text{post } r \in x \wedge y = x - \{r\} \end{array} \quad (5.8)$$

¹ Our extensive use of arrows to denote relations in the sequel is therefore rooted on common, informal practice. Unfortunately, mathematicians do not follow such practice and insist on regarding relations just as sets of pairs.

Explicit definition of *max* function:

$$\begin{aligned} \max &: \mathbb{Z} \times \mathbb{Z} \rightarrow \mathbb{Z} \\ \max(i, j) &= \text{if } i \leq j \text{ then } j \text{ else } i \end{aligned}$$

Its *implicit specification*:

$$\begin{aligned} \max(i: \mathbb{Z}, j: \mathbb{Z}) \ r: \mathbb{Z} \\ \text{post } r \in \{i, j\} \wedge i \leq r \wedge j \leq r \end{aligned}$$

Of a different nature is the following *pre/post*-pair:

$$\begin{aligned} \text{Sqrt} &: (i: \mathbb{R}) \ r: \mathbb{R} \\ \text{post } r^2 &= i \end{aligned}$$

Here the *specifier* is telling the *implementer* that either solution $r = +\sqrt{i}$ or $r = -\sqrt{i}$ will do.³ Indeed, *square root* is not a function, it is the binary relation:

$$r \text{ Sqrt } i \Leftrightarrow r^2 = i \quad (5.10)$$

We proceed with a thorough study of the concept of a binary relation, by analogy with a similar study carried out about functions in chapter 2.

5.4 RELATIONAL COMPOSITION AND CONVERSE

Such as functions, relations can be combined via composition ($R \cdot S$), defined as follows:

$$\begin{array}{c} B \xleftarrow{R} A \xleftarrow{S} C \\ \quad \quad \quad \curvearrowright \\ \quad \quad \quad R \cdot S \end{array} \quad b(R \cdot S)c \equiv \langle \exists a : b R a : a S c \rangle \quad (5.11)$$

Example: $\text{Uncle} = \text{Brother} \cdot \text{Parent}$, expanding to

$$u \text{ Uncle } c \equiv \langle \exists p :: u \text{ Brother } p \wedge p \text{ Parent } c \rangle$$

An explanation on the \exists -notation is on demand: \exists is an instance of a so-called *quantifier*, a main ingredient of formal logic. In this book we follow the so-called *Eindhoven quantifier* notation, whereby expressions of the form

$$\langle \forall x : P : Q \rangle$$

mean

“for **all** x in the range P , Q holds”

where P and Q are logical expressions involving x ; and expressions of the form

$$\langle \exists x : P : Q \rangle$$

³ This aspect of formal specification is called *vagueness*.

mean

“for **some** x in the range P , Q holds”.

Note how the symbols \exists and \forall “twist” letters E (exists) and A (all), respectively. P is known as the *range* of the quantification and Q as the quantified *term*.⁴ This logical notation enjoys a well-known set of properties, some of which are given in appendix A.2. As an example, by application of the \exists -trading rule (A.2), predicate $\langle \exists a :: b R a \wedge a S c \rangle$ in (5.11) can be written $\langle \exists a : b R a : a S c \rangle$.

Note how (5.11) *removes* \exists and bound variable a when applied from right to left. This is an example of conversion from pointwise to point-free notation, since “point” a also disappears. Indeed, we shall try and avoid lengthy, complex \forall, \exists -formulae by converting them to *pointfree* notation, as is the case in (5.11) once relational composition is used.

A simple calculation shows (5.11) to instantiate to (2.6) for the special case where R and S are functions, $R, S := f, g$:

$$\begin{aligned}
 b(f \cdot g)c &\equiv \langle \exists a :: b f a \wedge a g c \rangle \\
 &\equiv \{ \text{functions are univocal (simple) relations} \} \\
 &\quad \langle \exists a :: b = f a \wedge a = g c \rangle \\
 &\equiv \{ \exists\text{-trading rule (A.2)} \} \\
 &\quad \langle \exists a : a = g c : b = f a \rangle \\
 &\equiv \{ \exists\text{-“one-point” rule (A.6)} \} \\
 &\quad b = f (g c) \\
 &\quad \square
 \end{aligned}$$

Like its functional version (2.8), relation composition is associative:

$$R \cdot (S \cdot P) = (R \cdot S) \cdot P \quad (5.12)$$

Everywhere $T = R \cdot S$ holds, the replacement of T by $R \cdot S$ will be referred to as a “factorization” and that of $R \cdot S$ by T as “fusion”. Every relation $B \xleftarrow{R} A$ admits two trivial factorizations,

$$\begin{cases} R = R \cdot id_A \\ R = id_B \cdot R \end{cases} \quad (5.13)$$

where, for every X , id_X is the identity relation relating every element of X with itself (2.9). In other words: the identity (equality) *relation* coincides with the identity *function*.

In section 2.7 we introduced a very special case of function f — isomorphism — which has a converse f° such that (2.16) holds. A major advantage of generalizing functions to relations is that *every* relation $A \xrightarrow{R} B$ has a converse $A \xleftarrow{R^\circ} B$ defined by

$$b R a \Leftrightarrow a R^\circ b \quad (5.14)$$

⁴ In particular, Q or P can be universally False or True. Assertions of the form $\langle \forall x : \text{True} : Q \rangle$ or $\langle \exists x : \text{True} : Q \rangle$ are abbreviated to $\langle \forall x :: Q \rangle$ or $\langle \exists x :: Q \rangle$, respectively.

— the *passive voice* written relationally, as already mentioned. Two important properties of converse follow: it is an involution

$$(R^\circ)^\circ = R \quad (5.15)$$

and it commutes with composition in a contravariant way:

$$(R \cdot S)^\circ = S^\circ \cdot R^\circ \quad (5.16)$$

Converses of functions enjoy a number of properties from which the following is singled out as a way to introduce / remove them from logical expressions:

$$b(f^\circ \cdot R \cdot g)a \equiv (f b)R(g a) \quad (5.17)$$

For instance, the consequent of implication (5.6) could have been written $a(f^\circ \cdot id \cdot f)a'$, or even simpler as $a(f^\circ \cdot f)a'$, as it takes very little effort to show:

$$\begin{aligned} & a(f^\circ \cdot id \cdot f)a' \\ \equiv & \quad \{ (5.17) \} \\ & (f a)id(f a') \\ \equiv & \quad \{ b f a \equiv b = f a \} \\ & (f a) = id(f a') \\ \equiv & \quad \{ (2.9) \} \\ & f a = f a' \\ & \square \end{aligned}$$

Exercise 5.1. Let $sq \ x = x^2$ be the function that computes the square of a real number. Use (5.17) to show that (5.10) reduces to

$$Sqrt = sq^\circ$$

in relational pointfree notation.

□

Exercise 5.2. Give an implicit definition of function $f \ x = x^2 - 1$ in the form of a post-condition not involving subtraction. Then re-write it without variables using (5.17).

□

5.5 RELATIONAL EQUALITY

Recall that function equality (2.5) is established by extensionality:

$$f = g \text{ iff } \langle \forall a : a \in A : f a = g a \rangle$$

Also recall that $f = g$ only makes sense iff both functions have the same type, say $A \rightarrow B$. Can we do the same for relations? The relational generalization of (2.5) will be

$$R = S \text{ iff } \langle \forall a, b : a \in A \wedge b \in B : b R a \Leftrightarrow b S a \rangle \quad (5.18)$$

Since \Leftrightarrow is bi-implication, we can replace the term of the quantification by

$$(b R a \Rightarrow b S a) \wedge (b S a \Rightarrow b R a)$$

Now, what does $b R a \Rightarrow b S a$ mean? It simply captures relational *inclusion*,

$$R \subseteq S \text{ iff } \langle \forall a, b :: b R a \Rightarrow b S a \rangle \quad (5.19)$$

whose righthand side can also be written

$$\langle \forall a, b : b R a : b S a \rangle$$

by \forall -trading (A.1). Note the same pointwise-pointfree move when one reads (5.19) from right to left: \forall, a and b disappear.

Altogether, (5.18) can be written in less symbols as follows:

$$R = S \quad \equiv \quad R \subseteq S \wedge S \subseteq R \quad (5.20)$$

This way of establishing relational equality is usually referred to as *circular inclusion*. Note that relational inclusion (5.19) is a partial order: it is *reflexive*, since

$$R \subseteq R \quad (5.21)$$

holds for every R ; it is *transitive*, since for all R, S, T

$$R \subseteq S \wedge S \subseteq T \Rightarrow R \subseteq T \quad (5.22)$$

holds; and it is *antisymmetric*, as established by circular-inclusion (5.20) itself. Circular-inclusion is also jocosely known as the “ping-pong” method for establishing $R = S$: first calculate $R \subseteq S$ (“ping”) and then $S \subseteq R$ (“pong”). This can be performed in one go by adopting the following calculation layout:

$$\begin{array}{l} R \subseteq \dots \\ \subseteq S \\ \subseteq \dots \\ \subseteq R \\ \square \end{array}$$

This has the advantage of making apparent that not only R and S are the same, but also that every two steps in the circular reasoning are so (just choose a different start and stop point in the “circle”).

Circular inclusion (5.20) is not the only way to establish relational equality. A less obvious, but very useful way of calculating the equality of two relations is the method of *indirect equality*:

$$R = S \equiv \langle \forall X :: (X \subseteq R \Leftrightarrow X \subseteq S) \rangle \quad (5.23)$$

$$\equiv \langle \forall X :: (R \subseteq X \Leftrightarrow S \subseteq X) \rangle \quad (5.24)$$

The reader unaware of this way of indirectly setting algebraic equalities will recognize that the same pattern of indirection is used when establishing set equality via the membership relation, cf.

$$A = B \equiv \langle \forall x :: x \in A \Leftrightarrow x \in B \rangle$$

The typical layout of using any of these rules is the following:

$$\left\{ \begin{array}{l} X \subseteq R \\ \equiv \{ \dots \} \\ X \subseteq \dots \\ \equiv \{ \dots \} \\ X \subseteq S \\ :: \{ \text{indirect equality (5.23)} \} \\ R = S \\ \square \end{array} \right.$$

This proof method is very powerful and we shall make extensive use of it in the sequel. (The curious reader can have a quick look at section 5.9 for a simple illustration.)

RELATIONAL TYPES. From this point onwards we shall regard the type $B \leftarrow A$ as including not only all functions $f : A \rightarrow B$ but also all relations of the same type, $R : A \rightarrow B$. This is far more than we had before! In particular, type $A \rightarrow B$ includes:

- the *bottom* relation $B \xleftarrow{\perp} A$, which is such that, for all b, a ,

$$b \perp a \equiv \text{FALSE}$$

- the *topmost* relation $B \xleftarrow{\top} A$, which is such that, for all b, a ,

$$b \top a \equiv \text{TRUE}$$

The former is referred to as the void, or *empty* relation. The latter is known as the universal, or *coexistence* relation. Clearly, for every R ,

$$\perp \subseteq R \subseteq \top \quad (5.25)$$

and

$$R \cdot \perp = \perp \cdot R = \perp \quad (5.26)$$

hold. By (5.25) and (5.20), writing $R = \perp$ (respectively, $R = \top$) is the same as writing $R \subseteq \perp$ (respectively, $\top \subseteq R$).

A relation $B \xleftarrow{V} A$ is said to be a *vector* if either A or B are the singleton type 1. Relation $1 \xleftarrow{X} A$ is said to be a *row-vector*; clearly, $X \subseteq !$. Relation $B \xleftarrow{Z} 1$ is said to be a *column-vector*; clearly, $Z \subseteq !^\circ$.⁵ A relation of type $1 \xleftarrow{S} 1$ is called a *scalar*.

Last but not least, note that in a relational setting types $B \leftarrow A$ and B^A do not coincide — B^A is the type of all *functions* from A to B , while $B \leftarrow A$ is the type of all *relations* from A to B . Clearly, $B^A \subseteq B \leftarrow A$.

5.6 DIAGRAMS

As happens with functions, the arrow notation adopted for functions makes it possible to express relational formulæ using diagrams. This is a major ingredient of the relational method because it provides a graphical way of picturing relation types and relational constraints.

Paths in diagrams are built by arrow chaining, which corresponds to relational composition $R \cdot S$ (5.11), meaning “... is R of some S of ...” in natural language.

Assertions of the form $X \subseteq Y$ where X and Y are relation compositions can be represented graphically by rectangle-shaped diagrams, as is the case in

$$\begin{array}{ccc} \text{Descriptor} & \xleftarrow{FT} & \text{Handle} \\ \text{path} \downarrow & \subseteq & \downarrow \top \\ \text{Path} & \xleftarrow{FS^\circ} & \text{File} \end{array} \quad (5.27)$$

in the context of modelling a file-system. Relation FS models a *file store* (a table mapping file system paths to the respective files), FT is the *open-file descriptor* table (holding the information about the files that are currently open⁶), function $path$ yields the path of a file descriptor and \top is the largest possible relation between file-handles and files, as seen above. The diagram depicts the constraint:

$$path \cdot FT \subseteq FS^\circ \cdot \top \quad (5.28)$$

What does (5.28) mean, then, in predicate logic?

-
- ⁵ The column and row qualifiers have to do with an analogy with vectors in linear algebra.
- ⁶ Open files are manipulated by the file system via open file descriptor data structures, which hold various relevant metadata (e.g. current position within the file). Such descriptors are identified by file handles which the file system provides to applications that manipulate files. This indirection layer avoids unnecessary coupling between applications and the details of the file system implementation.

FROM DIAGRAMS TO LOGIC. We reason, using definitions (5.19,5.11) and the laws of the predicate calculus given in appendix A.2:

$$\begin{aligned}
& path \cdot FT \subseteq FS^\circ \cdot \top \\
\equiv & \quad \{ \text{'at most' ordering (5.19)} \} \\
& \langle \forall p, h : p(path \cdot FT)h : p(FS^\circ \cdot \top)h \rangle \\
\equiv & \quad \{ \text{composition (5.11) ; } path \text{ is a function} \} \\
& \langle \forall p, h : \langle \exists d : p = path\ d : d\ FT\ h \rangle : p(FS^\circ \cdot \top)h \rangle \\
\equiv & \quad \{ \text{quantifier calculus — } \exists\text{-trading (A.2)} \} \\
& \langle \forall p, h : \langle \exists d : d\ FT\ h : p = path\ d \rangle : p(FS^\circ \cdot \top)h \rangle \\
\equiv & \quad \{ \text{quantifier calculus — } \forall\text{-nesting (A.7)} \} \\
& \langle \forall h :: \langle \forall p : \langle \exists d : d\ FT\ h : p = path\ d \rangle : p(FS^\circ \cdot \top)h \rangle \rangle \\
\equiv & \quad \{ \text{quantifier calculus — splitting rule (A.13)} \} \\
& \langle \forall h :: \langle \forall d : d\ FT\ h : \langle \forall p : p = path\ d : p(FS^\circ \cdot \top)h \rangle \rangle \rangle \\
\equiv & \quad \{ \text{quantifier calculus — } \forall\text{-nesting (A.7)} \} \\
& \langle \forall d, h : d\ FT\ h : \langle \forall p : p = path\ d : p(FS^\circ \cdot \top)h \rangle \rangle \\
\equiv & \quad \{ \text{quantifier calculus — } \forall\text{-one-point rule (A.5)} \} \\
& \langle \forall d, h : d\ FT\ h : (path\ d)(FS^\circ \cdot \top)h \rangle
\end{aligned}$$

We still have to unfold term $(path\ d)(FS^\circ \cdot \top)h$:

$$\begin{aligned}
& (path\ d)(FS^\circ \cdot \top)h \\
\equiv & \quad \{ \text{composition (5.11)} \} \\
& \langle \exists x :: (path\ d)FS^\circ x \wedge x\top h \rangle \\
\equiv & \quad \{ \text{converse ; } x\top h \text{ always holds} \} \\
& \langle \exists x :: x\ FS\ (path\ d) \wedge \text{TRUE} \rangle \\
\equiv & \quad \{ \text{trivia} \} \\
& \langle \exists x :: x\ FS\ (path\ d) \rangle
\end{aligned}$$

In summary, $path \cdot FT \subseteq FS^\circ \cdot \top$ unfolds into

$$\langle \forall d, h : d\ FT\ h : \langle \exists x :: x\ FS\ (path\ d) \rangle \rangle \quad (5.29)$$

Literally:

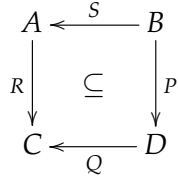
If h is the handle of some open-file descriptor d , then this holds the path of some existing file x .

In fewer words:

Non-existing files cannot be opened (referential integrity).

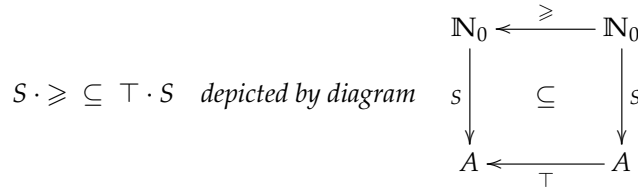
Thus we see how relation diagrams “hide” logically quantified formulae capturing properties of the systems one wishes to describe.

Compared with the commutative diagrams of previous chapters, a diagram



is said to be *semi-commutative* because $Q \cdot P \subseteq R \cdot S$ is not forced to hold, only $R \cdot S \subseteq Q \cdot P$ is. In case both hold, the \subseteq symbol is dropped, cf. (5.20).

Exercise 5.3. Let $a S n$ mean: “student a is assigned number n ”. Using (5.11) and (5.19), check that assertion



means that numbers are assigned to students in increasing order.

□

5.7 TAXONOMY OF BINARY RELATIONS

The Leibniz principle about functions (5.6) can now be simplified thanks to equivalence (5.19), as shown next:

$$\begin{aligned}
 & \langle \forall a, a' :: a = a' \Rightarrow f a = f a' \rangle \\
 \equiv & \quad \{ \text{introduction of } id; \text{ consequent as calculated already} \} \\
 & \langle \forall a, a' :: a = id a' \Rightarrow a(f^\circ \cdot f)a' \rangle \\
 \equiv & \quad \{ b f a \text{ means the same as } b = f a \} \\
 & \langle \forall a, a' :: a id a' \Rightarrow a(f^\circ \cdot f)a' \rangle \\
 \equiv & \quad \{ (5.19) \} \\
 & id \subseteq f^\circ \cdot f
 \end{aligned} \tag{5.30}$$

A similar calculation will reduce univocality (5.5) to

$$f \cdot f^\circ \subseteq id \tag{5.31}$$

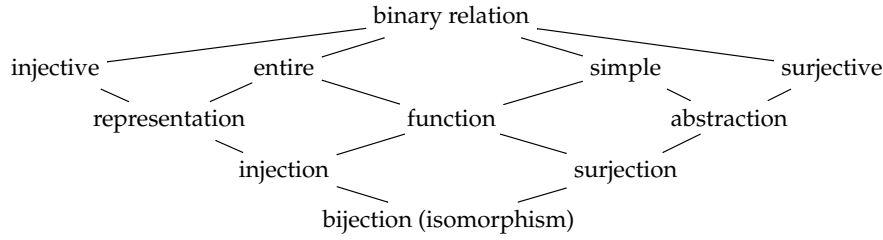


Figure 5.3.: Binary relation taxonomy

Thus a function f is characterized by comparing $f^\circ \cdot f$ and $f \cdot f^\circ$ with the identity.⁷

The exact characterization of functions as special cases of relations is achieved in terms of converse, which is in fact of paramount importance in establishing the whole taxonomy of binary relations depicted in figure 5.3. First, we need to define two important notions: given a relation $B \xleftarrow{R} A$, the *kernel* of R is the relation $A \xleftarrow{\ker R} A$ defined by:

$$\ker R = R^\circ \cdot R \quad (5.32)$$

Clearly, $a' \ker R a$ holds between any two sources a and a' which have (at least) a common target c such that $c R a'$ and $c R a$. We can also define its dual, $B \xleftarrow{\text{img } R} B$, called the *image* of R , defined by:⁸

$$\text{img } R = R \cdot R^\circ \quad (5.33)$$

From (5.15, 5.16) one immediately draws:

$$\ker (R^\circ) = \text{img } R \quad (5.34)$$

$$\text{img } (R^\circ) = \ker R \quad (5.35)$$

Kernel and image lead to the four top criteria of the taxonomy of figure 5.3:

	<i>Reflexive</i>	<i>Coreflexive</i>
$\ker R$	entire R	injective R
$\text{img } R$	surjective R	simple R

(5.36)

In words: a relation R is said to be *entire* (or total) iff its kernel is reflexive and to be *simple* (or functional) iff its image is coreflexive. Dually, R is *surjective* iff R° is entire, and R is *injective* iff R° is simple.

⁷ As we shall see in section 5.13, relations larger than the identity ($id \subseteq R$) are said to be *reflexive* and relations at most the identity ($R \subseteq id$) are said to be *coreflexive* or *partial identities*.

⁸ These operators are relational extensions of two concepts familiar from set theory: the image of a function f , which corresponds to the set of all y such that $\langle \exists x :: y = f x \rangle$, and the kernel of f , which is the equivalence relation $b \ker f a \Leftrightarrow (f b) = (f a)$. (See exercise 5.8 later on.)

Representing binary relations by Boolean matrices gives us a simple, graphical way of detecting properties such as simplicity, surjectiveness, and so on. Let the enumerated types $A = \{a_1, a_2, a_3, a_4, a_5\}$ and $B = \{b_1, b_2, b_3, b_4, b_5\}$ be given. Two examples of relations of type $A \rightarrow B$ are given in figure 5.4 — the leftmost and the rightmost, which we shall refer to as R and S , respectively.⁹ The matrix representing R is:

	a_1	a_2	a_3	a_4	a_5
b_1	0	1	0	0	0
b_2	1	0	0	0	0
b_3	0	0	1	1	0
b_4	0	0	0	0	1
b_5	0	0	0	0	0

(5.37)

The 1 addressed by b_2 and a_1 means that $b_2 R a_1$ holds, that between b_1 and a_2 means $b_1 R a_2$, and so on and so forth. Then, R is:

- *simple* because there is *at most* one 1 in every column
- *entire* because there is *at least* one 1 in every column
- not *injective* because there is *more than* one 1 in some row
- not *surjective* because some row (the last) has no 1s.

So this relation is a *function* that is neither an injection nor a surjection.

Let us now have a look at the matrix that represents $S : A \rightarrow B$:

	a_1	a_2	a_3	a_4	a_5
b_1	0	1	0	0	0
b_2	1	0	0	0	0
b_3	0	0	0	1	0
b_4	0	0	0	0	1
b_5	0	0	1	0	0

Now every row and every column has *exactly* one 1 — this tells us that S is not only a function but in fact a bijection. Looking at the matrix that represents $S^\circ : A \leftarrow B$,

	b_1	b_2	b_3	b_4	b_5
a_1	0	1	0	0	0
a_2	1	0	0	0	0
a_3	0	0	0	0	1
a_4	0	0	1	0	0
a_5	0	0	0	1	0

we realize that it also is a function, in fact another bijection. This gives us a rule of thumb for (constructively) checking for bijections (isomorphisms):

$$\text{A relation } f \text{ is a bijection iff its converse } f^\circ \text{ is a function } g \quad (5.38)$$

⁹ Credits: <http://www.matematikaria.com/unit/injective-surjective-bijective.html>. Note that we enumerate a_1, a_2, \dots from the top to the bottom.

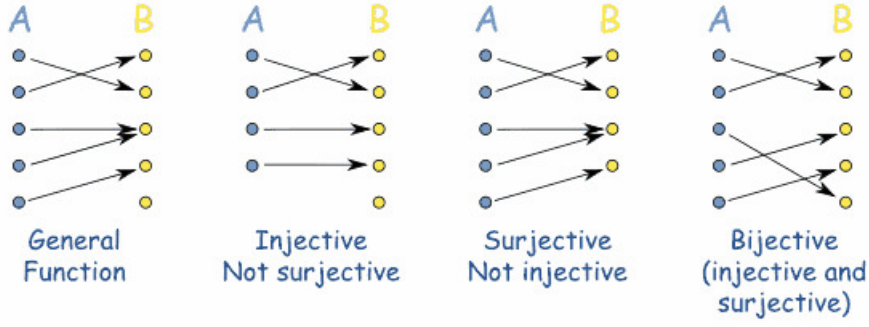


Figure 5.4.: Four binary relations.

Then g is also a bijection since $f^\circ = g \Leftrightarrow f = g^\circ$.¹⁰ Recall how some definitions of isomorphisms given before, e.g. (2.92), are nothing but applications of this rule $f^\circ = g$, once written pointwise with the help of (5.17):

$$f b = a \Leftrightarrow b = g a$$

Bijections (isomorphisms) are reversible functions — they don't lose any information. By contrast, $! : A \rightarrow 1$ (2.58) and indeed all constant functions $\underline{c} : A \rightarrow C$ (2.12) lose all the information contained in their inputs, recall (2.14). This property is actually more general,

$$\underline{c} \cdot R \subseteq \underline{c} \quad (5.39)$$

for all suitably typed R .

In the same way $! : A \rightarrow 1$ is always a constant function — in fact the *unique* possible function of its type, $f : 1 \rightarrow A$ is bound to be a constant function too, for any choice of a target value in non-empty A . Because there are as many such functions as elements if A , functions $\underline{a} : 1 \rightarrow A$ are referred to as *points*. These two situations correspond to isomorphisms $1^A \cong 1$ (2.97) and $A^1 \cong A$ (2.98), respectively. Two short-hands are introduced for the constant functions

$$\text{true} = \underline{\text{True}} \quad (5.40)$$

$$\text{false} = \underline{\text{False}} \quad (5.41)$$

Exercise 5.4. Prove (5.38) by completing:

$$\begin{aligned} & f \text{ and } f^\circ \text{ are functions} \\ \equiv & \{ \dots \} \\ & (id \subseteq \ker f \wedge \text{img } f \subseteq id) \wedge (id \subseteq \ker (f^\circ) \wedge \text{img } (f^\circ) \subseteq id) \\ \equiv & \{ \dots \} \\ & \vdots \\ \equiv & \{ \dots \} \\ & f \text{ is a bijection} \end{aligned}$$

¹⁰ The interested reader may go back to (2.18,2.19) at this point and check these rules in the light of (5.38).

□

Exercise 5.5. Compute, for the relations in figure 5.4, the kernel and the image of each relation. Why are all these relations functions? (NB: note that the types are not all the same.)

□

Exercise 5.6. Recall the definition of a constant function (2.12),

$$\begin{aligned} \underline{k} &: A \rightarrow K \\ \underline{k}a &= k \end{aligned}$$

where K is assumed to be non-empty. Show that $\ker \underline{k} = \top$ and compute which relations are defined by the expressions

$$\underline{b} \cdot \underline{c}^\circ, \quad \text{img } \underline{k} \quad (5.42)$$

Finally, show that (5.39) holds.

□

Exercise 5.7. Resort to (5.34,5.35) and (5.36) to prove the following rules of thumb:

- converse of injective is simple (and vice-versa) (5.43)

- converse of entire is surjective (and vice-versa) (5.44)

□

Exercise 5.8. Given a function $B \xleftarrow{f} A$, calculate the pointwise version

$$b(\ker f)a \equiv f b = f a \quad (5.45)$$

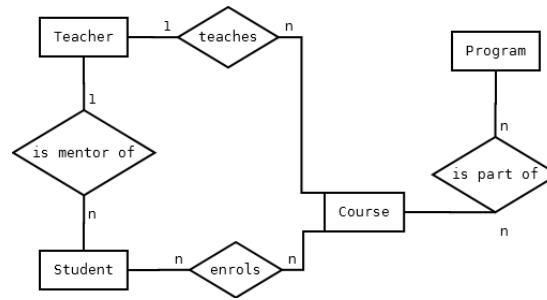
of $\ker f$. What is the outcome of the same exercise for $\text{img } f$?

□

ENTITY-RELATIONSHIP DIAGRAMS In the tradition of relational databases, so-called *entity-relationship* (ER) diagrams have become popular as an informal means for capturing the properties of the relationships involved in a particular database design.

Consider the following example of one such diagram:¹¹

11 Credits: <https://dba.stackexchange.com/questions>.



In the case of relation

$$Teacher \xleftarrow{\text{is mentor of}} Student$$

the drawing tells not only that some teacher may mentor more than one student, but also that a given student has exactly one mentor. So *is mentor of* is a *simple* relation (figure 5.3).

The possibility $n = 0$ allows for students with no mentor. Should this possibility be ruled out ($n > 1$), the relation would become also *entire*, i.e. a function. Then

$$t \text{ is mentor of } s$$

could be written

$$t = \text{is mentor of } s$$

— recall (5.7) — meaning:

$$t \text{ is the mentor of student } s.$$

That is, *is mentor of* would become an *attribute* of *Student*. Note how definite article “the” captures the presence of functions in normal speech. “The” means not only determinism (one and only one output) but also definedness (there is always one such output). In the case of *is mentor of* being simple but not entire, we have to say:

$$t \text{ is the mentor of student } s, \text{ if any.}$$

Exercise 5.9. Complete the exercise of declaring in $A \xrightarrow{R} B$ notation the other relations of the ER-diagram above and telling which properties in Figure 5.3 are required for such relations.

□

5.8 FUNCTIONS, RELATIONALLY

Among all binary relations, functions play a central role in relation algebra — as can be seen in figure 5.3. Recapitulating, a *function* f is a binary relation such that

Pointwise	Pointfree	
"Left" Uniqueness		
$b f a \wedge b' f a \Rightarrow b = b'$	$\text{img } f \subseteq \text{id}$	(f is simple)
Leibniz principle		
$a = a' \Rightarrow f a = f a'$	$\text{id} \subseteq \ker f$	(f is entire)

It turns out that *any* function f enjoys the following properties, known as *shunting rules*:

$$f \cdot R \subseteq S \equiv R \subseteq f^\circ \cdot S \quad (5.46)$$

$$R \cdot f^\circ \subseteq S \equiv R \subseteq S \cdot f \quad (5.47)$$

These will prove extremely useful in the sequel. Another very useful fact is the function *equality rule*:

$$f \subseteq g \equiv f = g \equiv f \supseteq g \quad (5.48)$$

Rule (5.48) follows immediately from (5.46,5.47) by "cyclic inclusion" (5.20):

$$\begin{aligned}
& f \subseteq g \\
& \equiv \{ \text{natural-id (2.10)} \} \\
& f \cdot \text{id} \subseteq g \\
& \equiv \{ \text{shunting on } f \text{ (5.46)} \} \\
& \text{id} \subseteq f^\circ \cdot g \\
& \equiv \{ \text{shunting on } g \text{ (5.47)} \} \\
& \text{id} \cdot g^\circ \subseteq f^\circ \\
& \equiv \{ \text{converses; identity} \} \\
& g \subseteq f
\end{aligned}$$

Then:

$$\begin{aligned}
& f = g \\
& \equiv \{ \text{cyclic inclusion (5.20)} \} \\
& f \subseteq g \wedge g \subseteq f \\
& \equiv \{ \text{above} \} \\
& f \subseteq g \\
& \equiv \{ \text{above} \} \\
& g \subseteq f \\
& \square
\end{aligned}$$

Exercise 5.10. Infer $\text{id} \subseteq \ker f$ (f is entire) and $\text{img } f \subseteq \text{id}$ (f is simple) from shunting rules (5.46) and (5.47).

□

Exercise 5.11. For $R := f$, the property (5.39) “immediately” coincides with (2.14). Why?

□

FUNCTION DIVISION. Given two functions $B \xrightarrow{g} C \xleftarrow{f} A$, we can compose f with the converse of g . This turns out to be a very frequent pattern in relation algebra, known as the *division* of f by g :

$$\frac{f}{g} = g^\circ \cdot f \quad \text{cf.} \quad \begin{array}{ccc} & \xleftarrow{\frac{f}{g}} & \\ B & & A \\ & \searrow g \quad \swarrow f & \\ & C & \end{array} \quad (5.49)$$

That is,

$$b \frac{f}{g} a \Leftrightarrow g b = f a$$

Think of the sentence:

Mary lives where John was born.

This can be expressed by a division:

$$\text{Mary} \frac{\text{birthplace}}{\text{residence}} \text{John} \Leftrightarrow \text{residence Mary} = \text{birthplace John}$$

Thus $R = \frac{\text{birthplace}}{\text{residence}}$ is the relation “... resides in the birthplace of ...”. In general,

$b \frac{f}{g} a$ means “the g of b is the f of a ”.

This combinator enjoys a number of interesting properties, for instance:

$$\frac{f}{id} = f \quad (5.50)$$

$$\left(\frac{f}{g} \right)^\circ = \frac{g}{f} \quad (5.51)$$

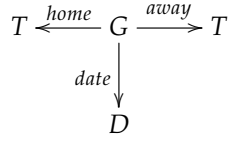
$$\frac{f \cdot h}{g \cdot k} = k^\circ \cdot \frac{f}{g} \cdot h \quad (5.52)$$

$$\frac{f}{f} = \ker f \quad (5.53)$$

$$a \neq b \Leftrightarrow \frac{a}{b} = \perp \quad (5.54)$$

Function division is a special case of the more general, and important, concept of relational division, a topic that shall be addressed in section 5.19.

Exercise 5.12. The teams (T) of a football league play games (G) at home or away, and every game takes place in some date:



Moreover, (a) No team can play two games on the same date; (b) All teams play against each other but not against themselves; (c) For each home game there is another game away involving the same two teams. Show that

$$\text{id} \subseteq \frac{\text{away}}{\text{home}} \cdot \frac{\text{away}}{\text{home}} \quad (5.55)$$

captures one of the requirements above — which?

□

Exercise 5.13. Check the properties of function division given above.

□

5.9 MEET AND JOIN

Like sets, two relations of the same type, say $B \xleftarrow{R,S} A$, can be intersected or joined in the obvious way:

$$b (R \cap S) a \equiv b R a \wedge b S a \quad (5.56)$$

$$b (R \cup S) a \equiv b R a \vee b S a \quad (5.57)$$

$R \cap S$ is usually called *meet* (intersection) and $R \cup S$ is called *join* (union). They lift pointwise conjunction and disjunction, respectively, to the pointfree level. Their meaning is nicely captured by the following *universal* properties:¹²

$$X \subseteq R \cap S \equiv X \subseteq R \wedge X \subseteq S \quad (5.58)$$

$$R \cup S \subseteq X \equiv R \subseteq X \wedge S \subseteq X \quad (5.59)$$

Meet and join have the expected properties, e.g. *associativity*

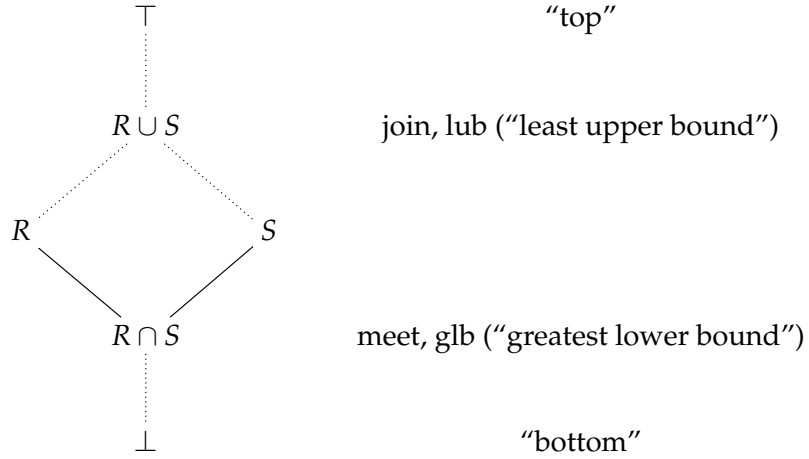
$$(R \cap S) \cap T = R \cap (S \cap T)$$

¹² Recall the generic notions of *greatest lower bound* and *least upper bound*, respectively.

proved next by indirect equality (5.23):

$$\begin{aligned}
 & X \subseteq (R \cap S) \cap T \\
 \equiv & \quad \{ \cap\text{-universal (5.58) twice} \} \\
 & (X \subseteq R \wedge X \subseteq S) \wedge X \subseteq T \\
 \equiv & \quad \{ \wedge \text{ is associative} \} \\
 & X \subseteq R \wedge (X \subseteq S \wedge X \subseteq T) \\
 \equiv & \quad \{ \cap\text{-universal (5.58) twice} \} \\
 & X \subseteq R \cap (S \cap T) \\
 \therefore & \quad \{ \text{indirection (5.23)} \} \\
 & (R \cap S) \cap T = R \cap (S \cap T) \\
 & \square
 \end{aligned}$$

In summary, type $B \leftarrow A$ forms a lattice:



DISTRIBUTIVE PROPERTIES. As it will be proved later, *composition* distributes over *union*

$$R \cdot (S \cup T) = (R \cdot S) \cup (R \cdot T) \quad (5.60)$$

$$(S \cup T) \cdot R = (S \cdot R) \cup (T \cdot R) \quad (5.61)$$

while distributivity over *intersection* is side-conditioned:

$$(S \cap Q) \cdot R = (S \cdot R) \cap (Q \cdot R) \Leftrightarrow \begin{cases} Q \cdot \text{img } R \subseteq Q \\ \vee \\ S \cdot \text{img } R \subseteq S \end{cases} \quad (5.62)$$

$$R \cdot (Q \cap S) = (R \cdot Q) \cap (R \cdot S) \Leftrightarrow \begin{cases} (\ker R) \cdot Q \subseteq Q \\ \vee \\ (\ker R) \cdot S \subseteq S \end{cases} \quad (5.63)$$

Properties (5.60,5.61) express the *bilinearity* of relation composition with respect to relational join. These, and properties such as e.g.

$$(R \cap S)^\circ = R^\circ \cap S^\circ \quad (5.64)$$

$$(R \cup S)^\circ = R^\circ \cup S^\circ \quad (5.65)$$

will be shown to derive from a general construction that will be explained in section 5.18.

Exercise 5.14. Show that

$$R \cap \perp = \perp \quad (5.66)$$

$$R \cap \top = R \quad (5.67)$$

$$R \cup \top = \top \quad (5.68)$$

$$R \cup \perp = R \quad (5.69)$$

using neither (5.56) nor (5.57).

□

Exercise 5.15. Prove the union simplicity rule:

$$M \cup N \text{ is simple} \quad \equiv \quad M, N \text{ are simple and } M \cdot N^\circ \subseteq id \quad (5.70)$$

Using converses, derive from (5.70) the corresponding rule for injective relations.

□

Exercise 5.16. Prove the distributive property:

$$g^\circ \cdot (R \cap S) \cdot f = g^\circ \cdot R \cdot f \cap g^\circ \cdot S \cdot f \quad (5.71)$$

□

Exercise 5.17. Let $\text{bag} : A^* \rightarrow \mathbb{N}_0^A$ be the function that, given a finite sequence (list), indicates the number of occurrences of its elements, for instance,

$$\text{bag} [a, b, a, c] \ a = 2$$

$$\text{bag} [a, b, a, c] \ b = 1$$

$$\text{bag} [a, b, a, c] \ c = 1$$

Let $\text{ord} : A^* \rightarrow \mathbb{B}$ be the obvious predicate assuming a total order predefined in A . Finally, let $\text{true} = \underline{\text{True}}$ (5.40). Having defined

$$S = \frac{\text{bag}}{\text{bag}} \cap \frac{\text{true}}{\text{ord}} \quad (5.72)$$

identify the type of S and, going pointwise and simplifying, tell which operation is specified by S .

□

Exercise 5.18. Derive the distributive properties:

$$k^\circ \cdot (f \cup g) = \frac{f}{k} \cup \frac{g}{k} \quad , \quad k^\circ \cdot (f \cap g) = \frac{f}{k} \cap \frac{g}{k} \quad (5.73)$$

□

5.10 RELATIONAL THINKING

Binary relations provide a natural way of describing real life situations. Relation algebra can be used to reason about such formal descriptions. This can be achieved using suitable relational combinators (and their laws), in the *pointfree* style.

Let us see a simple example of such a *relational thinking* taking one of the PROPOSITIONES AD ACUENDOS IUUENES (“Problems to Sharpen the Young”) proposed by abbot Alcuin of York († 804) as case study. Alcuin states his puzzle in the following way, in Latin:

XVIII. PROPOSITIO DE HOMINE ET CAPRA ET LVPO. *Homo quidam debebat ultra fluuium transferre lupum, capram, et fasciculum cauli. Et non potuit aliam nauem inuenire, nisi quae duos tantum ex ipsis ferre ualebat. Praeceptum itaque ei fuerat, ut omnia haec ultra illaesa omnino transferret. Dicat, qui potest, quomodo eis illaesis transire potuit?*

Our starting point will be the following (rather free) translation of the above to English:

XVIII. FOX, GOOSE AND BAG OF BEANS PUZZLE. *A farmer goes to market and purchases a fox, a goose, and a bag of beans. On his way home, the farmer comes to a river bank and hires a boat. But in crossing the river by boat, the farmer could carry only himself and a single one of his purchases - the fox, the goose or the bag of beans. (If left alone, the fox would eat the goose, and the goose would eat the beans.) Can the farmer carry himself and his purchases to the far bank of the river, leaving each purchase intact?*

We wish to describe the essence of this famous puzzle, which is the *guarantee* that

under no circumstances does the fox eat the goose or the goose eat the beans.

Clearly, we need two data types:

Being = { *Farmer, Fox, Goose, Beans* }
Bank = { *Left, Right* }

Then we identify a number of relations involving such data:

$$\begin{array}{ccc} \textit{Being} & \xrightarrow{\textit{Eats}} & \textit{Being} \\ & \text{where} \downarrow & \\ & \textit{Bank} & \xrightarrow{\textit{CROSS}} \textit{Bank} \end{array} \quad (5.74)$$

Clearly, $\text{cross Left} = \text{Right}$ and $\text{cross Right} = \text{Left}$. So cross is its own inverse and therefore a bijection (5.38). Relation Eats can be described by the Boolean matrix:

$$\text{Eats} = \begin{array}{c|cccc} & \text{Fox} & \text{Goose} & \text{Beans} & \text{Farmer} \\ \hline \text{Fox} & 0 & 1 & 0 & 0 \\ \text{Goose} & 0 & 0 & 1 & 0 \\ \text{Beans} & 0 & 0 & 0 & 0 \\ \text{Farmer} & 0 & 0 & 0 & 0 \end{array} \quad (5.75)$$

Relation $\text{where} : \text{Being} \rightarrow \text{Bank}$ is necessarily a function because:

- everyone is somewhere in a bank (where is entire)
- no one can be in both banks at the same time (where is simple)

Note that there are only two constant functions of type $\text{Being} \rightarrow \text{Bank}$, Right and Left . The puzzle consists in changing from the state $\text{where} = \text{Right}$ to the state $\text{where} = \text{Left}$, for instance, without violating the property that *nobody eats anybody*. How does one record such a property? We need two auxiliary relations capturing, respectively:

- Being at the same bank:

$$\text{SameBank} = \ker \text{where}$$

- Risk of somebody eating somebody else:

$$\text{CanEat} = \text{SameBank} \cap \text{Eats}$$

Then “starvation” is ensured by the *Farmer’s* presence at the same bank:

$$\text{CanEat} \subseteq \text{SameBank} \cdot \underline{\text{Farmer}} \quad (5.76)$$

By (5.46), this “starvation” property (5.76) converts to:

$$\text{where} \cdot \text{CanEat} \subseteq \text{where} \cdot \underline{\text{Farmer}}$$

In this version, (5.76) can be depicted as a diagram

$$\begin{array}{ccc} \text{Being} & \xleftarrow{\text{CanEat}} & \text{Being} \\ \text{where} \downarrow & \subseteq & \downarrow \underline{\text{Farmer}} \\ \text{Bank} & \xleftarrow{\text{where}} & \text{Being} \end{array} \quad (5.77)$$

which “reads” in a nice way:

where (somebody) CanEat (somebody else) (that’s) where (the) Farmer (is).

Diagram (5.27) given earlier can now be identified as another example of assertion expressed relationally. Diagrams of this kind capture properties of data models that one wishes to hold at any time during the lifetime of the system being described. Such properties are commonly referred to as *invariants* and their preservation by calculation will be the main aim of chapter 7.

Exercise 5.19. Calculate the following pointwise version of the “starvation” property (5.77) by introducing quantifiers and simplifying:

$$\langle \forall b', b : b' \text{ Eat } b : \text{ where } b' = \text{ where } b \Rightarrow \text{ where } b' = \text{ where Farmer} \rangle$$

□

Exercise 5.20. Recalling property (5.39), show that the “starvation” property (5.77) is satisfied by any of the two constant functions that model the start or end states of the Alcuin puzzle.

□

5.11 MONOTONICITY

As expected, relational composition is monotonic:

$$\frac{\begin{array}{c} R \subseteq S \\ T \subseteq U \end{array}}{(R \cdot T) \subseteq (S \cdot U)} \quad (5.78)$$

Indeed, all relational combinators studied so far are also monotonic, namely

$$R \subseteq S \Rightarrow R^\circ \subseteq S^\circ \quad (5.79)$$

$$R \subseteq S \wedge U \subseteq V \Rightarrow R \cap U \subseteq S \cap V \quad (5.80)$$

$$R \subseteq S \wedge U \subseteq V \Rightarrow R \cup U \subseteq S \cup V \quad (5.81)$$

hold.

Monotonicity and transitivity (5.22) are important properties for reasoning about a given relational inclusion $R \subseteq S$. In particular, the following rules are of help by relying on a “mid-point” relation M , $R \subseteq M \subseteq S$ (analogy with interval arithmetics).

- Rule A — *lowering the upper side*:

$$\begin{array}{c} R \subseteq S \\ \Leftarrow \{ M \subseteq S \text{ is known ; transitivity of } \subseteq \text{ (5.22) } \} \\ R \subseteq M \end{array}$$

Then proceed with $R \subseteq M$.

- Rule B — *raising the lower side*:

$$\begin{array}{c}
 R \subseteq S \\
 \Leftarrow \quad \{ R \subseteq M \text{ is known; transitivity of } \subseteq \} \\
 M \subseteq S
 \end{array}$$

Then proceed with $M \subseteq S$.

The following proof of shunting property (5.46) combines these rules with monotonicity and circular implication:

$$\begin{array}{c}
 R \subseteq f^\circ \cdot S \\
 \Leftarrow \quad \{ id \subseteq f^\circ \cdot f ; \text{raising the lower-side} \} \\
 f^\circ \cdot f \cdot R \subseteq f^\circ \cdot S \\
 \Leftarrow \quad \{ \text{monotonicity of } (f^\circ \cdot) \} \\
 f \cdot R \subseteq S \\
 \Leftarrow \quad \{ f \cdot f^\circ \subseteq id ; \text{lowering the upper-side} \} \\
 f \cdot R \subseteq f \cdot f^\circ \cdot S \\
 \Leftarrow \quad \{ \text{monotonicity of } (f \cdot) \} \\
 R \subseteq f^\circ \cdot S
 \end{array}$$

Thus the equivalence in (5.46) is established by circular implication.

Rules A and B should be used only where other proof techniques (notably indirect equality) fail. They assume judicious choice of the mid-point relation M , at each step. The choice of an useless M can drive the proof nowhere.

Exercise 5.21. *Unconditional distribution laws*

$$\begin{array}{lcl}
 (P \cap Q) \cdot S & = & (P \cdot S) \cap (Q \cdot S) \\
 R \cdot (P \cap Q) & = & (R \cdot P) \cap (R \cdot Q)
 \end{array}$$

will hold provide one of R or S is simple and the other injective. Tell which, justifying.

□

Exercise 5.22. *Prove that relational composition preserves all relational classes in the taxonomy of figure 5.3.*

□

5.12 RULES OF THUMB

Quite often, involved reasoning in logic arguments can be replaced by simple and elegant calculations in relation algebra that arise thanks to smart “rules of thumb”. We have already seen two such rules, (5.43) and (5.44). Two others are:

$$\text{- smaller than injective (simple) is injective (simple)} \quad (5.82)$$

$$\text{- larger than entire (surjective) is entire (surjective)} \quad (5.83)$$

Let us see these rules in action in trying to infer what can be said of two functions f and r such that

$$f \cdot r = id$$

holds. On the one hand,

$$\begin{aligned} & f \cdot r = id \\ \equiv & \quad \{ \text{equality of functions} \} \\ & f \cdot r \subseteq id \\ \equiv & \quad \{ \text{shunting} \} \\ & r \subseteq f^\circ \end{aligned}$$

Since f is simple, f° is injective and so is r because “smaller than injective is injective”. On the other hand,

$$\begin{aligned} & f \cdot r = id \\ \equiv & \quad \{ \text{equality of functions} \} \\ & id \subseteq f \cdot r \\ \equiv & \quad \{ \text{shunting} \} \\ & r^\circ \subseteq f \end{aligned}$$

Since r is entire, r° is surjective and so is f because “larger than surjective is surjective”. We conclude that f is surjective and r is injective wherever $f \cdot r = id$ holds. Since both are functions, we furthermore conclude that

f is an *abstraction* and r is a *representation*

— cf. Figure 5.3.

The reason for this terminology can now be explained. Given $f : A \leftarrow C$ and $r : C \leftarrow A$ such that $f \cdot r = id$, that is, for all $a \in A$, $f(r a) = a$, think of C as a domain of *concrete* objects and of A as a domain of *abstract* data. For instance, let $A = \mathbb{B}$ and $C = \mathbb{N}_0$. Then define

$$\begin{cases} r : \mathbb{B} \rightarrow \mathbb{N}_0 \\ r b = \text{if } b \text{ then } k \text{ else } 0 \end{cases}$$

(where k is any natural number different from 0) and

$$\begin{cases} f : \mathbb{B} \leftarrow \mathbb{N}_0 \\ f\ n = \text{if } n = 0 \text{ then False else True} \end{cases}$$

Clearly, by the definitions of f and r :

$$\begin{aligned} f(r\ b) &= \text{if } (\text{if } b \text{ then } k \text{ else } 0) = 0 \text{ then False else True} \\ &\equiv \{ \text{conditional-fusion rule (2.71)} \} \\ f(r\ b) &= \text{if } (\text{if } b \text{ then } k = 0 \text{ else True}) \text{ then False else True} \\ &\equiv \{ k = 0 \text{ is always false} \} \\ f(r\ b) &= \text{if } (\text{if } b \text{ then False else True}) \text{ then False else True} \\ &\equiv \{ \text{pointwise definition of } \neg b \} \\ f(r\ b) &= \text{if } \neg b \text{ then False else True} \\ &\equiv \{ \text{trivial} \} \\ &\quad b \end{aligned}$$

That is, r represents the Booleans True and False by natural numbers while f abstracts from such real numbers back to Booleans. r being injective means $r\ \text{False} \neq r\ \text{True}$, that is, the Boolean information is not lost in the representation.¹³ f being surjective means that any Boolean is representable. Note that $r \cdot f = \text{id}$ does not hold: $r(f\ 1) = r\ \text{True} = k$ and $k \neq 1$ in general.

The abstraction/representation pair (f, r) just above underlies the way Booleans are handled in programming languages such as C, for instance. Experienced programmers will surely agree that often what is going on in the code they write are processes of representing information using primitive data structures available from the adopted programming language. For instance, representing finite sets by finite lists corresponds to the *abstraction* given by *elems* (5.1).

Exercise 5.23. Recalling exercise 5.17, complete the definition of

$$\begin{aligned} \text{bag } []\ a &= 0 \\ \text{bag } (h : t)\ a &= \text{let } b = \text{bag } t \text{ in if } \dots \end{aligned}$$

Is this function an abstraction or a representation? Justify your answer informally.

□

Exercise 5.24. Show that:

- $R \cap S$ is injective (simple) provided one of R or S is so
- $R \cup S$ is entire (surjective) provided one of R or S is so.

□

¹³ That is, r causes *no confusion* in the representation process.

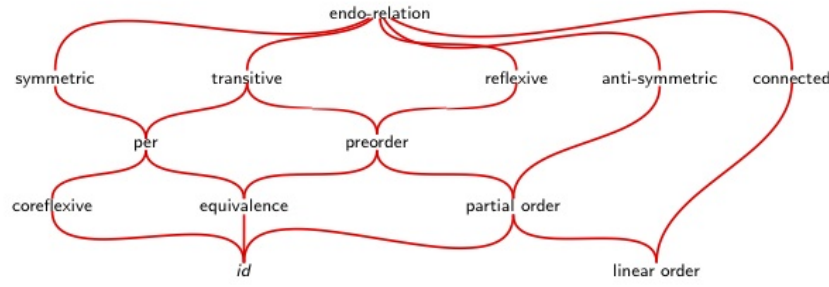


Figure 5.5.: Taxonomy of endorelations.

5.13 ENDO-RELATIONS

Relations in general are of type $A \rightarrow B$, for some A and B . In the special case that $A = B$ holds, a relation $R : A \rightarrow A$ is said to be an *endo-relation*, or a *graph*. The $A = B$ coincidence gives room for some extra terminology, extending some already given. Besides an endo-relation $A \xleftarrow{R} A$ being

$$\text{reflexive:} \quad \text{iff } id \subseteq R \quad (5.84)$$

$$\text{coreflexive:} \quad \text{iff } R \subseteq id \quad (5.85)$$

it can also be:

$$\text{transitive:} \quad \text{iff } R \cdot R \subseteq R \quad (5.86)$$

$$\text{symmetric:} \quad \text{iff } R \subseteq R^\circ (\equiv R = R^\circ) \quad (5.87)$$

$$\text{anti-symmetric:} \quad \text{iff } R \cap R^\circ \subseteq id \quad (5.88)$$

$$\text{irreflexive:} \quad \text{iff } R \cap id = \perp \quad (5.89)$$

$$\text{connected:} \quad \text{iff } R \cup R^\circ = \top \quad (5.90)$$

By combining these criteria, endo-relations $A \xleftarrow{R} A$ can further be classified as in figure 5.5. In summary:

- *Preorders* are reflexive and transitive orders.
Example: $\text{age } y \leq \text{age } x$.
- *Partial orders* are anti-symmetric preorders
Example: $y \subseteq x$ where x and y are sets.
- *Linear orders* are connected partial orders
Example: $y \leq x$ in \mathbb{N}_0
- *Equivalences* are symmetric preorders
Example: $\text{age } y = \text{age } x$.¹⁴

¹⁴ Kernels of functions are always equivalence relations, see exercise 5.25.

- *Pers* are partial equivalences

Example: $y \text{ IsBrotherOf } x$.

Preorders are normally denoted by asymmetric symbols such as e.g. $y \sqsubseteq x, y \leq x$. In case of a function f such that

$$f \cdot (\sqsubseteq) \subseteq (\leq) \cdot f \quad (5.91)$$

we say that f is monotonic. Indeed, this is equivalent to

$$a \sqsubseteq b \Rightarrow (f a) \leq (f b)$$

once shunting (5.46) takes place, and variables are added and handled via (5.17). Another frequent situation is that of two functions f and g such that

$$f \subseteq (\leq) \cdot g \quad (5.92)$$

This converts to the pointwise

$$\langle \forall a :: f a \leq g a \rangle$$

that is, f is *always at most* g for all possible inputs. The following abbreviation is often used to capture this ordering on functions induced by a pre-order (\leq) on their outputs:

$$f \dot{\leq} g \text{ iff } f \subseteq (\leq) \cdot g \quad (5.93)$$

For instance, $f \dot{\leq} id$ means $f a \leq a$ for all inputs a .

CLOSURE OPERATORS Given a partial order (\leq), a function f is said to be a *closure operator* iff

$$(\leq) \cdot f = f^\circ \cdot (\leq) \cdot f \quad (5.94)$$

holds. Let us write the same with points — via (5.17) —, for all x, y :

$$y \leq f x \Leftrightarrow f x \leq f y \quad (5.95)$$

Clearly, for $(\geq) = (\leq)^\circ$, (5.94) can also be written

$$f^\circ \cdot (\geq) = f^\circ \cdot (\geq) \cdot f$$

Any of these alternatives is an elegant way of defining a closure operator f , in so far it can be shown to be equivalent to the conjunction of three facts about f : (a) f is monotonic; (b) $id \dot{\leq} f$ and (c) $f = f \cdot f$.

As an example, consider the function that *closes* a finite set of natural numbers by filling in the intermediate numbers, e.g. $f \{4, 2, 6\} = \{2, 3, 4, 5, 6\}$. Clearly, $x \subseteq f x$. If you apply f again, you get

$$f \{2, 3, 4, 5, 6\} = \{2, 3, 4, 5, 6\}$$

This happens because f is a closure operator.

Exercise 5.25. Knowing that property

$$f \cdot f^\circ \cdot f = f \quad (5.96)$$

holds for every function f , prove that $\ker f = \frac{f}{f}$ (5.53) is an equivalence relation.
 \square

Exercise 5.26. From $\ker ! = \top$ and (5.96) infer

$$\top \cdot R \subseteq \top \cdot S \quad \Leftrightarrow \quad R \subseteq \top \cdot S \quad (5.97)$$

Conclude that $(\top \cdot)$ is a closure operator.

\square

Exercise 5.27. Generalizing the previous exercise, show that pre/post-composition with functional kernels are closure operations:

$$S \cdot \ker f \subseteq R \cdot \ker f \quad \equiv \quad S \subseteq R \cdot \ker f \quad (5.98)$$

$$\ker f \cdot S \subseteq \ker f \cdot R \quad \equiv \quad S \subseteq \ker f \cdot R \quad (5.99)$$

\square

Exercise 5.28. Consider the relation

$$b R a \Leftrightarrow \text{team } b \text{ is playing against team } a$$

Is this relation: reflexive? irreflexive? transitive? anti-symmetric? symmetric? connected?

\square

Exercise 5.29. Expand criteria (5.86) to (5.90) to pointwise notation.

\square

Exercise 5.30. A relation R is said to be co-transitive or dense iff the following holds:

$$\langle \forall b, a : b R a : \langle \exists c : b R c : c R a \rangle \rangle \quad (5.100)$$

Write the formula above in PF notation. Find a relation (eg. over numbers) which is co-transitive and another which is not.

\square

Exercise 5.31. Check which of the following properties,

transitive, symmetric, anti-symmetric, connected

hold for the relation *Eats* (5.75) of the Alcuin puzzle.

□

Exercise 5.32. Show that (5.55) of exercise 5.12 amounts to forcing relation $\text{home} \cdot \text{away}^\circ$ to be symmetric.

□

5.14 RELATIONAL PAIRING

Recall from sections 2.8 and 2.9 that functions can be composed in parallel and in alternation, giving rise to so-called *products* and *coproducts*. Does a product diagram like (2.23),

$$\begin{array}{ccccc} A & \xleftarrow{\pi_1} & A \times B & \xrightarrow{\pi_2} & B \\ & \searrow f & \uparrow \langle f, g \rangle & \nearrow g & \\ & & C & & \end{array}$$

make sense when f e g are generalized to relations R and S ? We start from definition (2.20),

$$\langle f, g \rangle c \stackrel{\text{def}}{=} (f c, g c)$$

to try and see what such a generalization could mean. The relational, pointwise expression of function $\langle f, g \rangle$ is

$$y = \langle f, g \rangle c$$

which can be rephrased to $(a, b) = \langle f, g \rangle c$, knowing that $\langle f, g \rangle$ is of type $C \rightarrow A \times B$ in (2.23). We reason:

$$\begin{aligned} (a, b) &= \langle f, g \rangle c \\ \equiv & \{ \langle f, g \rangle c = (f c, g c); \text{equality of pairs} \} \\ & \left\{ \begin{array}{l} a = f c \\ b = g c \end{array} \right. \\ \equiv & \{ y = f x \Leftrightarrow y f x \} \\ & \left\{ \begin{array}{l} a f c \\ b g c \end{array} \right. \\ & \square \end{aligned}$$

By in-lining the conjunction expressed by the braces just above, one gets

$$(a, b) \langle f, g \rangle c \Leftrightarrow a f c \wedge b g c$$

which proposes the generalization:

$$(a, b) \langle R, S \rangle c \Leftrightarrow a R c \wedge b S c \quad (5.101)$$

Recalling the projections $\pi_1(a, b) = a$ and $\pi_2(a, b) = b$, let us try and remove variables a, b and c from the above, towards a closed definition of $\langle R, S \rangle$:

$$\begin{aligned} & (a, b) \langle R, S \rangle c \Leftrightarrow a R c \wedge b S c \\ \equiv & \quad \{ \pi_1(a, b) = a \text{ and } \pi_2(a, b) = b \} \\ & (a, b) \langle R, S \rangle c \Leftrightarrow \pi_1(a, b) R c \wedge \pi_2(a, b) S c \\ \equiv & \quad \{ (5.17) \text{ twice} \} \\ & (a, b) \langle R, S \rangle c \Leftrightarrow (a, b) (\pi_1^\circ \cdot R) c \wedge (a, b) (\pi_2^\circ \cdot S) c \\ \equiv & \quad \{ (5.56) \} \\ & (a, b) \langle R, S \rangle c \Leftrightarrow (a, b) (\pi_1^\circ \cdot R \cap \pi_2^\circ \cdot S) c \\ \equiv & \quad \{ (5.19) \} \\ & \langle R, S \rangle = \pi_1^\circ \cdot R \cap \pi_2^\circ \cdot S \end{aligned} \quad (5.102)$$

We proceed to investigating what kind of universal property $\langle R, S \rangle$, defined by $\pi_1^\circ \cdot R \cap \pi_2^\circ \cdot S$, satisfies. The strategy is to use indirect equality:

$$\begin{aligned} & X \subseteq \langle R, S \rangle \\ \equiv & \quad \{ (5.102) \} \\ & X \subseteq \pi_1^\circ \cdot R \cap \pi_2^\circ \cdot S \\ \equiv & \quad \{ (5.58) \} \\ & \left\{ \begin{array}{l} X \subseteq \pi_1^\circ \cdot R \\ X \subseteq \pi_2^\circ \cdot S \end{array} \right. \\ \equiv & \quad \{ \text{shunting} \} \\ & \left\{ \begin{array}{l} \pi_1 \cdot X \subseteq R \\ \pi_2 \cdot X \subseteq S \end{array} \right. \end{aligned}$$

In summary, the universal property of $\langle R, S \rangle$ is:

$$X \subseteq \langle R, S \rangle \Leftrightarrow \left\{ \begin{array}{l} \pi_1 \cdot X \subseteq R \\ \pi_2 \cdot X \subseteq S \end{array} \right. \quad (5.103)$$

For functions, $X, R, S := k, f, g$ it can be observed that (5.103) coincides with (2.63). But otherwise, the corollaries derived from (5.103) are different from those that emerge from (2.63). For instance, cancellation becomes:

$$\left\{ \begin{array}{l} \pi_1 \cdot \langle R, S \rangle \subseteq R \\ \pi_2 \cdot \langle R, S \rangle \subseteq S \end{array} \right.$$

This tells us that pairing R with S has the (side) effect of deleting from R all those inputs for which S is undefined (and vice-versa), since output pairs require that *both* relations respond to the input. Thus, for relations, laws such as the \times -fusion rule (2.26) call for a side-condition:

$$\begin{aligned} \langle R, S \rangle \cdot T &= \langle R \cdot T, S \cdot T \rangle \\ &\Leftarrow R \cdot (\text{img } T) \subseteq R \vee S \cdot (\text{img } T) \subseteq S \end{aligned} \quad (5.104)$$

Clearly,

$$\langle R, S \rangle \cdot f = \langle R \cdot f, S \cdot f \rangle \quad (5.105)$$

holds, since $\text{img } f \subseteq \text{id}$. Moreover, the *absorption* law (2.27) remains unchanged,

$$(R \times S) \cdot \langle P, Q \rangle = \langle R \cdot P, S \cdot Q \rangle \quad (5.106)$$

where $R \times S$ is defined in the same way as for functions:

$$R \times S = \langle R \cdot \pi_1, S \cdot \pi_2 \rangle \quad (5.107)$$

As generalization of (5.105) and also immediate by monotonicity,

$$\langle R, S \rangle \cdot T = \langle R \cdot T, S \cdot T \rangle$$

holds for T simple.

Because (5.103) is not the universal property of a product, we tend to avoid talking about relational *products* and rather talk about relational *pairing* instead.¹⁵ In spite of the weaker properties, relational pairing has interesting laws, namely

$$\langle R, S \rangle^\circ \cdot \langle X, Y \rangle = (R^\circ \cdot X) \cap (S^\circ \cdot Y) \quad (5.108)$$

that will prove quite useful later on.

Exercise 5.33. Derive from (5.108) the following properties:

$$\frac{f}{g} \cap \frac{h}{k} = \frac{\langle f, h \rangle}{\langle g, k \rangle} \quad (5.109)$$

$$(5.110)$$

$$\ker \langle R, S \rangle = \ker R \cap \ker S \quad (5.111)$$

$\langle R, \text{id} \rangle$ is always injective, for whatever R

□

Exercise 5.34. Recalling (5.38), prove that $\text{swap} = \langle \pi_2, \pi_1 \rangle$ (2.32) is its own converse and therefore a bijection.

□

¹⁵ Relational products do exist but are not obtained by $\langle R, S \rangle$. For more about this see section 5.23 later on.

Exercise 5.35. *Derive from the laws studied thus far the following facts about relational pairing:*

$$id \times id = id \quad (5.112)$$

$$(R \times S) \cdot (P \times Q) = (R \cdot P) \times (S \cdot Q) \quad (5.113)$$

□

5.15 RELATIONAL COPRODUCTS

Let us now show that, in contrast with products, coproducts extend perfectly from functions to relations, that is, universal property (2.65) extends to

$$X = [R, S] \Leftrightarrow \begin{cases} X \cdot i_1 = R \\ X \cdot i_2 = S \end{cases} \quad (5.114)$$

where $X : A + B \rightarrow C$, $R : A \rightarrow C$ and $S : B \rightarrow C$ are binary relations. First of all, we need to understand what $[R, S]$ means. Our starting point is $+$ -cancellation, recall (2.40):

$$\begin{aligned} & \begin{cases} [g, h] \cdot i_1 = g \\ [g, h] \cdot i_2 = h \end{cases} \\ \equiv & \quad \{ \text{equality of functions} \} \\ & \begin{cases} g \subseteq [g, h] \cdot i_1 \\ h \subseteq [g, h] \cdot i_2 \end{cases} \\ \equiv & \quad \{ \text{shunting followed by (5.57)} \} \\ & g \cdot i_1^\circ \cup h \cdot i_2^\circ \subseteq [g, h] \end{aligned}$$

On the other hand:

$$\begin{aligned} & \begin{cases} [g, h] \cdot i_1 = g \\ [g, h] \cdot i_2 = h \end{cases} \\ \equiv & \quad \{ \text{equality of functions} \} \\ & \begin{cases} [g, h] \cdot i_1 \subseteq g \\ [g, h] \cdot i_2 \subseteq h \end{cases} \\ \Rightarrow & \quad \{ \text{monotonicity} \} \\ & \begin{cases} [g, h] \cdot i_1 \cdot i_1^\circ \subseteq g \cdot i_1^\circ \\ [g, h] \cdot i_2 \cdot i_2^\circ \subseteq h \cdot i_2^\circ \end{cases} \\ \Rightarrow & \quad \{ \text{monotonicity (5.81) and distribution (5.60)} \} \\ & [g, h] \cdot (i_1 \cdot i_1^\circ \cup i_2 \cdot i_2^\circ) \subseteq g \cdot i_1^\circ \cup h \cdot i_2^\circ \end{aligned}$$

$$\begin{aligned} &\equiv \{ \text{img } i_1 \cup \text{img } i_2 = id, \text{ more about this below} \} \\ [g, h] &\subseteq g \cdot i_1^\circ \cup h \cdot i_2^\circ \end{aligned}$$

Altogether, we obtain:

$$[g, h] = g \cdot i_1^\circ \cup h \cdot i_2^\circ$$

Note how this matches with (2.37), once variables are introduced:

$$c [g, h] x \Leftrightarrow \langle \exists a : x = i_1 a : c = g a \rangle \vee \langle \exists b : x = i_2 b : c = h b \rangle$$

Fact

$$\text{img } i_1 \cup \text{img } i_2 = id \quad (5.115)$$

assumed above is a property stemming from the construction of coproducts,

$$A + B \stackrel{\text{def}}{=} \{ i_1 a \mid a \in A \} \cup \{ i_2 b \mid b \in B \}$$

since i_1 and i_2 are the *only* constructors of data of type $A + B$. Another property implicit in this construction is:

$$i_1^\circ \cdot i_2 = \perp \quad (5.116)$$

equivalent to its converse $i_2^\circ \cdot i_1 = \perp$. It spells out that, for any $a \in A$ and $b \in B$, $i_1 a = i_2 b$ is impossible.¹⁶ In other words, the union is a *disjoint* one.

Let us now generalize the above to relations instead of functions,

$$[R, S] = R \cdot i_1^\circ \cup S \cdot i_2^\circ \quad (5.117)$$

and show that (5.114) holds. First of all,

$$\begin{aligned} X &= R \cdot i_1^\circ \cup S \cdot i_2^\circ \\ \Rightarrow &\{ \text{compose both sides with } i_1 \text{ and simplify; similarly for } i_2 \} \\ X \cdot i_1 &= R \wedge X \cdot i_2 = S \end{aligned}$$

The simplifications arise from i_1 and i_2 being injections, so their kernels are identities. On the other hand, $i_1^\circ \cdot i_2 = \perp$ and $i_2^\circ \cdot i_1 = \perp$, as seen above. The converse implication (\Leftarrow) holds:

$$\begin{aligned} X &= R \cdot i_1^\circ \cup S \cdot i_2^\circ \\ \equiv &\{ (5.115) \} \\ X \cdot (\text{img } i_1 \cup \text{img } i_2) &= R \cdot i_1^\circ \cup S \cdot i_2^\circ \\ \equiv &\{ \text{distribution} \} \\ X \cdot \text{img } i_1 \cup X \cdot \text{img } i_2 &= R \cdot i_1^\circ \cup S \cdot i_2^\circ \\ \Leftarrow &\{ \text{Leibniz} \} \\ X \cdot i_1 \cdot i_1^\circ &= R \cdot i_1^\circ \wedge X \cdot i_2 \cdot i_2^\circ = S \cdot i_2^\circ \\ \Leftarrow &\{ \text{monotonicity} \} \\ X \cdot i_1 &= R \wedge X \cdot i_2 = S \end{aligned}$$

□

¹⁶ Note that in (2.36) this is ensured by always choosing two different tags $t_1 \neq t_2$.

Thus (5.114) holds in general, for relations:

$$(B + C) \rightarrow A \begin{array}{c} \xrightarrow{[-,-]^\circ} \\ \cong \\ \xleftarrow{[-,-]} \end{array} (B \rightarrow A) \times (C \rightarrow A) \quad (5.118)$$

A most useful consequence of this is that all results known for coproducts of functions are valid for relational coproducts. In particular, relational direct sum

$$R + S = [i_1 \cdot R, i_2 \cdot S] \quad (5.119)$$

can be defined satisfying (2.43), (2.44) etc with relations replacing functions. Moreover, the McCarthy conditional (2.70) can be extended to relations in the expected way:

$$p \rightarrow R, S \stackrel{\text{def}}{=} [R, S] \cdot p? \quad (5.120)$$

The property for sums (coproducts) corresponding to (5.108) for products is:

$$[R, S] \cdot [T, U]^\circ = (R \cdot T^\circ) \cup (S \cdot U^\circ) \quad (5.121)$$

This *divide-and-conquer* rule is essential to *parallelizing* relation composition by so-called *block decomposition*.

Finally, the *exchange law* (2.49) extends to relations,

$$[\langle R, S \rangle, \langle T, V \rangle] = \langle [R, T], [S, V] \rangle \quad (5.122)$$

cf.

$$\begin{array}{ccccc} A & \xrightarrow{i_1} & A + B & \xleftarrow{i_2} & B \\ & \searrow & \downarrow T & \swarrow & \\ R \downarrow & & C \times D & & \downarrow V \\ C & \xleftarrow{\pi_1} & C \times D & \xrightarrow{\pi_2} & D \end{array}$$

For the proof see the following exercise.

Exercise 5.36. Relying on both (5.114) and (5.105) prove (5.122). Moreover, prove

$$(R + S)^\circ = R^\circ + S^\circ \quad (5.123)$$

□

Exercise 5.37. From (5.117) prove (5.121). Then show that

$$\text{img } [R, S] = \text{img } R \cup \text{img } S \quad (5.124)$$

follows immediately from (5.121).

□

Exercise 5.38. Prove that the coproduct $[R, S]$ is injective iff both R, S are injective and $R^\circ \cdot S = \perp$.

□

Exercise 5.39. Prove:

$$\frac{f}{g} \times \frac{h}{k} = \frac{f \times h}{g \times k} \quad (5.125)$$

$$\frac{f}{g} + \frac{h}{k} = \frac{f + h}{g + k} \quad (5.126)$$

□

5.16 ON KEY-VALUE DATA MODELS

Simple relations abstract what is currently known as the *key-value-pair* data model in modern databases.¹⁷ In this setting, given a *simple* relation $K \xrightarrow{S} V$, K is regarded as a type of data *keys* and V as a type of data *values*.

By pairing (5.102) such key-value-pairs one obtains more elaborate stores. Conversely, one may use projections to select particular key-attribute relationships from key-value stores. Note that keys and values can be *anything* (that is, of any type) and, in particular, they can be compound, for instance

$$\underbrace{\text{PartitionKey} \times \text{SortKey}}_K \rightarrow \underbrace{\text{Type} \times \dots}_V$$

in the example of figure 5.6.¹⁸

The example furthermore shows how keys and values can structure themselves even further. In particular, “*schema is defined per item*” indicates that the values may be of coproduct types, something like $\text{Title} \times (1 + \text{Author} \times (1 + \text{Date} \times \dots))$, for instance. Although the simplicity of the columnar model suggested by the key-value principle is somewhat sacrificed in the example, this shows how expressive *simple* relations involving *product* and *coproduct* types are.

One of the standard variations of the key-value model is to equip keys with time-stamps indicating *when* the pair was inserted or modified in the store, for instance

$$\text{Student} \times \text{Course} \times \text{Time} \rightarrow \text{Result} \quad (5.127)$$

¹⁷ For example, Hbase, Amazon DynamoDB, and so on, are examples of database systems that use the key-value pair data model.

¹⁸ Credits: <https://aws.amazon.com/nosql/key-value/>.

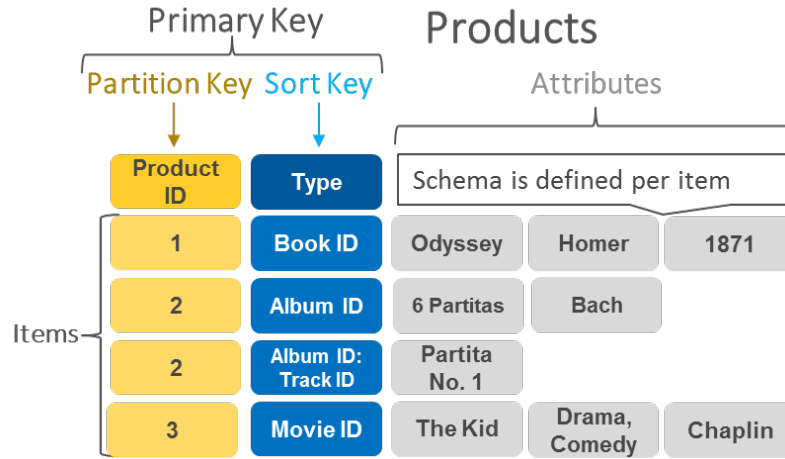


Figure 5.6.: Key-value data model instance.

telling the possibly different results of students in exams of a particular course. This combination of the key-value model with that of *temporal* (also called *historical*) databases is very powerful.

The relational combinators studied in this book apply naturally to key-value-pair storage processing and offer themselves as a powerful, pointfree high-level language for handling such data in a “noSQL” style.

5.17 WHAT ABOUT RELATIONAL “CURRYING”?

Recall isomorphism (2.93),

$$(C^B)^A \begin{array}{c} \xrightarrow{\text{uncurry}} \\ \cong \\ \xleftarrow{\text{curry}} \end{array} C^{A \times B}$$

that is at the core of the way functions are handled in functional programming. Does this isomorphism hold when functions are generalized to relations, something like...

$$A \times B \rightarrow C \cong A \rightarrow \dots?$$

Knowing that the type $A \times B \rightarrow C$ of relations is far larger than $C^{A \times B}$, it can be anticipated that the isomorphism will not extend to relations in the same way. In fact, a rather simpler one happens instead, among relations:

$$A \times B \rightarrow C \begin{array}{c} \xrightarrow{\text{trans}} \\ \cong \\ \xleftarrow{\text{untrans}} \end{array} A \rightarrow C \times B \quad (5.128)$$

This tells us the (obvious, but very useful) fact that relations involving product types can be reshaped in any way we like, leftwards or rightwards.

It is quite convenient to overload the notation used for functions and write \bar{R} to denote *trans* R and \hat{R} to denote *untrans* R . Then the isomorphism above is captured by universal property,¹⁹

$$\begin{array}{ccc} C \times B & & (C \times B) \times B \xrightarrow{\epsilon} C \\ \bar{R} \uparrow & & \uparrow \bar{R} \times id \\ A & & A \times B \end{array} \quad \begin{array}{c} \nearrow R \end{array}$$

where

$$\bar{R} = \langle R, \pi_2 \rangle \cdot \pi_1^\circ \quad \begin{array}{ccc} C \times B & & \\ \bar{R} \uparrow & \nwarrow \langle R, \pi_2 \rangle & \\ A & \xrightarrow{\pi_1^\circ} & A \times B \end{array} \quad (5.129)$$

that is

$$(c, b) \bar{R} a \equiv c R (a, b)$$

Moral: every n -ary relation can be expressed as a binary relation; moreover, where each particular attribute is placed (input/output) is irrelevant.

By *converse duality*, $(\hat{S})^\circ = \overline{(S^\circ)}$, we obtain the definition of relational “uncurrying”:

$$\hat{S} = \pi_1 \cdot \langle S^\circ, \pi_2 \rangle^\circ$$

Then

$$\epsilon = \hat{id} = \pi_1 \cdot \langle id, \pi_2 \rangle^\circ.$$

With points:

$$c_2 \in ((c_1, b_1), b_2) \equiv c_2 = c_1 \wedge b_1 = b_2$$

THE “PAIRING WHEEL” RULE The flexibility offered by (5.128) means that, in relation algebra, the information altogether captured by the three relations M , P and Q in

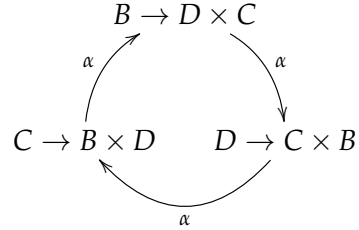
$$\begin{array}{ccc} & B & \\ & \uparrow M & \\ & A & \\ Q \swarrow & & \searrow P \\ C & & D \end{array} \quad (5.130)$$

can be aggregated in several ways, namely

¹⁹ Compare with (2.84).

$$\begin{aligned}
B &\xrightarrow{\langle P, Q \rangle \cdot M^\circ} D \times C \\
D &\xrightarrow{\langle Q, M \rangle \cdot P^\circ} C \times B \\
C &\xrightarrow{\langle M, P \rangle \cdot Q^\circ} B \times C
\end{aligned}$$

all isomorphic to each other:



The rotation among relations and types justifies the name “pairing wheel” given to (5.130). Isomorphism α holds in the sense that every entry of one of the aggregates is uniquely represented by another entry in any other aggregate, for instance:

$$\begin{aligned}
&(d, c) (\langle P, Q \rangle \cdot M^\circ) b \\
&= \{ \text{composition ; pairing} \} \\
&\langle \exists a : d P a \wedge c Q a : a M^\circ b \rangle \\
&= \{ \text{converse; } \wedge \text{ is associative and commutative} \} \\
&\langle \exists a :: (c Q a \wedge b M a) \wedge a P^\circ d \rangle \\
&= \{ \text{composition ; pairing} \} \\
&(c, b) (\langle Q, M \rangle \cdot P^\circ) d
\end{aligned}$$

Thus: $\alpha (\langle P, Q \rangle \cdot M^\circ) = (\langle Q, M \rangle \cdot P^\circ)$.

Exercise 5.40. Express α in terms of trans (5.128) and its converse (5.129).

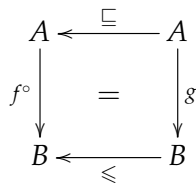
□

5.18 GALOIS CONNECTIONS

Recall from section 5.13 that a preorder is a reflexive and transitive relation. Given two preorders \leq and \sqsubseteq , one may relate arguments and results of pairs of suitably typed functions f and g in a particular way,

$$f^\circ \cdot \sqsubseteq = \leq \cdot g \quad (5.131)$$

as in the diagram:



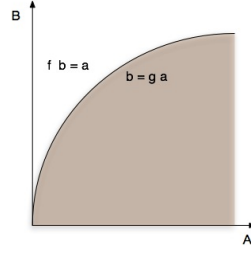


Figure 5.7.: Graphical interpretation of equation (5.131): (a) relation

$B \xleftarrow{(\leq) \cdot g} A$ is the “area” below function g wrt. \leq ; (b) relation $B \xrightarrow{f \cdot (\sqsubseteq)} A$ is the “area” above function f wrt. \sqsubseteq , to the right (oriented 90°); (c) f and g are such that these areas are the same.

In this very special situation, f, g are said to be *Galois connected*. We write

$$f \vdash g \quad (5.132)$$

as abbreviation of (5.131) when the two preorders \sqsubseteq, \leq are implicit from the context. Another way to represent this is:

$$(A, \sqsubseteq) \begin{array}{c} \xrightarrow{g} \\ \xleftarrow{f} \end{array} (B, \leq)$$

Function f (resp. g) is referred to as the *left* (resp. *right*) adjoint of the connection. By introducing variables in both sides of (5.131) via (5.17), we obtain, for all x and y

$$(f x) \sqsubseteq y \quad \equiv \quad x \leq (g y) \quad (5.133)$$

In particular, the two preorders in (5.131) can be the identity id , in which case (5.131) reduces to $f^\circ = g$, that is, f and g are each-other inverses — i.e., isomorphisms. Therefore, the Galois connection concept is a generalization of the concept of isomorphism.²⁰

Quite often, the two adjoints are *sections* of binary operators. Recall that, given a binary operator $a \theta b$, its two sections $(a\theta)$ and (θb) are unary functions f and g such that, respectively:

$$f = (a\theta) \quad \equiv \quad f b = a \theta b \quad (5.134)$$

$$g = (\theta b) \quad \equiv \quad g a = a \theta b \quad (5.135)$$

Galois connections in which the two preorders are relation inclusion ($\leq, \sqsubseteq := \subseteq, \subseteq$) and whose adjoints are sections of relational combinators are particularly interesting because they express universal properties about such combinators. Table 3 lists some connections that are relevant for this book.

²⁰ Interestingly, every Galois connection is on its turn a special case of an *adjunction*, recall (4.57). Just promote the adjoints f and g in (5.133) to functors, and replace the preorder symbols by arrows. This “syntactic trick” can be taken as a rough sketch of a formal, categorical argument that we shall skip for the time being.

$(f X) \subseteq Y \equiv X \subseteq (g Y)$			
Description	f	g	Obs.
converse	$(-)^{\circ}$	$(-)^{\circ}$	
<i>shunting rule</i>	$(h \cdot)$	$(h^{\circ} \cdot)$	h is a function
<i>“converse” shunting rule</i>	$(\cdot h^{\circ})$	$(\cdot h)$	h is a function
difference	$(- - R)$	$(R \cup)$	
implication	$(R \cap -)$	$(R \Rightarrow -)$	

Table 3.: Sample of Galois connections in the relational calculus. The general formula given on top is a logical equivalence universally quantified on S and R . It has a left part involving *left adjoint* f and a right part involving *right adjoint* g .

It is remarkably easy to recover known properties of the relation calculus from table 3. For instance, the first row yields

$$X^{\circ} \subseteq Y \equiv X \subseteq Y^{\circ} \quad (5.136)$$

since $f = g = (-)^{\circ}$ in this case. Thus converse is its own self adjoint. From this we derive

$$R \subseteq S \equiv R^{\circ} \subseteq S^{\circ} \quad (5.137)$$

by making $X, Y := R, S^{\circ}$ and simplifying by involution (5.15). Moreover, the entry marked “*shunting rule*” in the table leads to

$$h \cdot X \subseteq Y \equiv X \subseteq h^{\circ} \cdot Y$$

for all h, X and Y . By taking converses, one gets another entry in table 3, namely

$$X \cdot h^{\circ} \subseteq Y \equiv X \subseteq Y \cdot h$$

These are the equivalences (5.46) and (5.47) that we have already met, popularly known as “shunting rules”.

The fourth and fifth rows in the table are Galois connections that respectively introduce two new relational operators — relational *difference* $S - R$ and relational *implication* $R \Rightarrow S$ — as a *left adjoint* and an *right adjoint*, respectively:

$$X - R \subseteq Y \equiv X \subseteq Y \cup R \quad (5.138)$$

$$R \cap X \subseteq Y \equiv X \subseteq R \Rightarrow Y \quad (5.139)$$

There are *many* advantages in describing the meaning of relational operators by Galois connections. Further to the systematic tabulation of operators (of which table 3 is just a sample), the concept of a Galois connection is a *generic* one, which offers a rich algebra of *generic* properties, namely:

- both adjoints f and g in a Galois connection are monotonic;

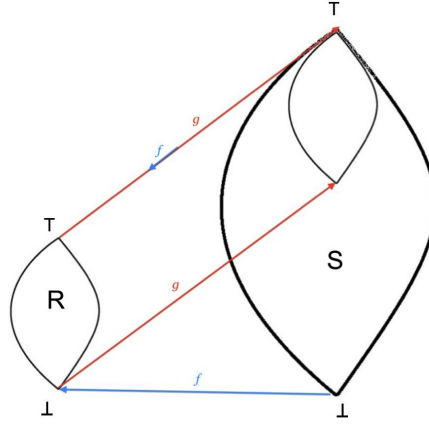


Figure 5.8.: Left-perfect Galois connection $f \vdash g$ involving two lattices S and R .

- left adjoint f distributes with join and right-adjoint g distributes with meet, wherever these exist:

$$f(b \sqcup b') = (f b) \vee (f b') \quad (5.140)$$

$$g(a \wedge a') = (g a) \sqcap (g a') \quad (5.141)$$

- left adjoint f preserves infima and right-adjoint g preserves suprema, wherever these exist:²¹

$$f \perp = \perp \quad (5.142)$$

$$g \top = \top \quad (5.143)$$

- two cancellation laws hold,

$$(f \cdot g)a \leq a \quad \text{and} \quad b \sqsubseteq (g \cdot f)b \quad (5.144)$$

respectively known as *left-cancellation* and *right-cancellation*.

- Semi-inverse properties:

$$f = f \cdot g \cdot f \quad (5.145)$$

$$g = g \cdot f \cdot g \quad (5.146)$$

It may happen that a cancellation law holds up to equality, for instance $f(g a) = a$, in which case the connection is said to be *perfect* on the particular side. The picture of a left-perfect Galois connection $f \vdash g$ is given in figure 5.8.²²

²¹ In these case both orders will form a so-called *lattice* structure.

²² Adapted from [4].

Let us take for instance Galois connection (5.138) as example. Following the general rules above, we get *for free*: the monotonicity of $(- - R)$,

$$X \subseteq Z \Rightarrow X - R \subseteq Z - R$$

the monotonicity of $(- \cup R)$,

$$X \subseteq Z \Rightarrow X \cup R \subseteq Z \cup R$$

the distribution of $(- - R)$ over *join*,

$$(X \cup Y) - R = (X - R) \cup (Y - R)$$

the distribution of $(- \cup R)$ over *meet*,

$$(X \cap Y) \cup R = (X \cup R) \cap (Y \cup R)$$

the preservation of infima by $(- - R)$,

$$\perp - R = \perp$$

the preservation of suprema by $(- \cup R)$,

$$\top \cup R = \top$$

left-cancellation ($Y := X - R$),

$$X \subseteq (X - R) \cup R$$

right-cancellation ($X := Y \cup R$),

$$(Y \cup R - R) \subseteq Y$$

and finally the semi-inverse properties:

$$X - ((X - R) \cup R) = X - R$$

$$((X \cup R) - R) \cup R = X \cup R$$

The reader is invited to extract similar properties from the other connections listed in table 3. Altogether, we get 50 properties out of this table! Such is the power of *generic* concepts in mathematics.

Two such connections were deliberately left out from table 3, which play a central role in relation algebra and will deserve a section of their own — section 5.19.

Exercise 5.41. Show that $R - S \subseteq R$, $R - \perp = R$ and $R - R = \perp$ hold.

□

Exercise 5.42. Infer

$$b(R \Rightarrow S)a \equiv (b R a) \Rightarrow (b S a) \quad (5.147)$$

from the Galois connection

$$R \cap X \subseteq Y \quad \equiv \quad X \subseteq (R \Rightarrow Y) \quad (5.148)$$

Suggestion: note that $b(R \Rightarrow S)a$ can be written $\text{id} \subseteq \underline{b}^\circ \cdot (R \Rightarrow S) \cdot \underline{a}$ (check this!). Then proceed with (5.148) and simplify.

□

Exercise 5.43. (Lexicographic orders) The lexicographic chaining of two relations R and S is defined by:

$$R ; S = R \cap (R^\circ \Rightarrow S) \quad (5.149)$$

Show that (5.149) is the same as stating the universal property:

$$X \subseteq (R ; S) \equiv X \subseteq R \wedge X \cap R^\circ \subseteq S$$

□

Exercise 5.44. Let students in a course have two numeric marks,

$$\mathbb{N}_0 \xleftarrow{\text{mark1}} \text{Student} \xrightarrow{\text{mark2}} \mathbb{N}_0$$

and define the preorders:

$$\leq_{\text{mark1}} = \text{mark1}^\circ \cdot \leq \cdot \text{mark1}$$

$$\leq_{\text{mark2}} = \text{mark2}^\circ \cdot \leq \cdot \text{mark2}$$

Spell out in pointwise notation the meaning of lexicographic ordering

$$\leq_{\text{mark1}} ; \leq_{\text{mark2}}$$

□

NEGATION We define $\neg R = R \Rightarrow \perp$ since $b(\neg R)a \Leftrightarrow \neg(b R a)$. Clearly, $\neg \top = \perp$. It can also be shown that

$$R \cup \neg R = \top \quad (5.150)$$

holds and therefore:

$$\top - R \subseteq R \Rightarrow \perp \quad (5.151)$$

From the Galois connection of $R \Rightarrow S$ and through the usual rule of indirect equality, one immediately infers the so-called *de Morgan law*,

$$\neg(R \cup S) = (\neg R) \cap (\neg S) \quad (5.152)$$

and other expected properties analogous to logic negation. One of the most famous rules for handling negated relations is the so-called *Schröder's rule*:

$$\neg Q \cdot S^\circ \subseteq \neg R \Leftrightarrow R^\circ \cdot \neg Q \subseteq \neg S \quad (5.153)$$

Exercise 5.45. Assuming

$$f^\circ \cdot (R \Rightarrow S) \cdot g = (f^\circ \cdot R \cdot g) \Rightarrow (f^\circ \cdot S \cdot g) \quad (5.154)$$

and (5.151), prove:

$$\underline{c}^\circ \cdot (\top - \underline{c}) = \perp \quad (5.155)$$

□

5.19 RELATION DIVISION

However intimidating it may sound, structuring a calculus in terms of Galois connections turns out to be a great simplification, leading to *rules* that make the reasoning closer to school algebra. Think for instance the rule used at school to reason about whole division of two natural numbers x and y ,

$$z \times y \leq x \equiv z \leq x \div y \quad (y > 0) \quad (5.156)$$

assumed universally quantified in all its variables. Pragmatically, it expresses a “shunting” rule which enables one to trade between a whole division in the upper side of an inequality and a multiplication in the lower side. This rule is easily identified as the Galois connection

$$\underbrace{z(\times y)}_f \leq x \Leftrightarrow z \leq \underbrace{x(\div y)}_g.$$

where multiplication is the left adjoint and division is the right adjoint: $(\times y) \vdash (\div y)$, for $y \neq 0$.²³

As seen in the previous section, many properties of (\times) and (\div) can be inferred from (5.156), for instance the cancellation $(x \div y) \times y \leq x$ — just replace z by $x \div y$ and simplify, and so on.

A parallel with relation algebra could be made by trying a rule similar to (5.156),

$$Z \cdot Y \subseteq X \equiv Z \subseteq X/Y \quad (5.157)$$

which suggests that, like integer multiplication, relational composition has an right adjoint, denoted X / Y . The question is: does such a *relation division* operator actually exist? Proceeding with the parallel, note that, in the same way

$$z \times y \leq x \equiv z \leq x \div y$$

means that $x \div y$ is the largest *number* which multiplied by y approximates x , (5.157) means that X/Y is the largest *relation* Z which, precomposed with Y , approximates X .

²³ This connection is perfect on the lower side since $(z \times y) \div y = z$.

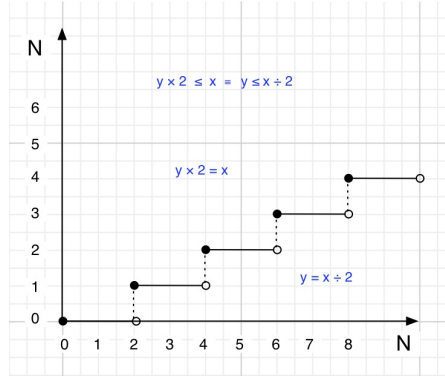
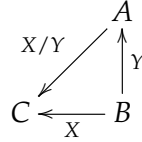


Figure 5.9.: Picturing Galois connection $(\times 2) \dashv (\div 2)$ as in figure 5.7. $f = (\times 2)$ is the left adjoint of $g = (\div 2)$. The area below $g = (\div 2)$ is the same as the area above $f = (\times 2)$. $f = (\times 2)$ is not surjective. $g = (\div 2)$ is not injective.

What is the pointwise meaning of X/Y ? Let us first of all equip (5.157) with a type diagram:

$$Z \cdot Y \subseteq X \equiv Z \subseteq X/Y$$

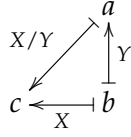


Then we calculate:²⁴

$$\begin{aligned}
 & c \ (X/Y) \ a \\
 \equiv & \quad \{ \text{introduce points } C \xleftarrow{c} 1 \text{ and } A \xleftarrow{a} 1 ; (5.17) \} \\
 & x(\underline{c}^\circ \cdot (X/Y) \cdot \underline{a})x \\
 \equiv & \quad \{ \forall\text{-one-point (A.5)} \} \\
 & x' = x \Rightarrow x'(\underline{c}^\circ \cdot (X/Y) \cdot \underline{a})x \\
 \equiv & \quad \{ \text{go pointfree (5.19)} \} \\
 & id \subseteq \underline{c}^\circ \cdot (X/Y) \cdot \underline{a} \\
 \equiv & \quad \{ \text{shunting rules} \} \\
 & \underline{c} \cdot \underline{a}^\circ \subseteq X/Y \\
 \equiv & \quad \{ \text{universal property (5.157)} \} \\
 & \underline{c} \cdot \underline{a}^\circ \cdot Y \subseteq X \\
 \equiv & \quad \{ \text{now shunt } \underline{c} \text{ back to the right} \} \\
 & \underline{a}^\circ \cdot Y \subseteq \underline{c}^\circ \cdot X \\
 \equiv & \quad \{ \text{back to points via (5.17)} \} \\
 & \langle \forall b : a \ Y \ b : c \ X \ b \rangle
 \end{aligned}$$

²⁴ Following the strategy suggested in exercise 5.42.

In summary:

$$c (X/Y) a \equiv \langle \forall b : a Y b : c X b \rangle \quad (5.158)$$


In words: in the same way relation *composition* hides an *existential* quantifier (5.11), *relation division* (5.158) hides a *universal* one. Let us feel what (5.158) means through an example: let

$$\begin{aligned} a Y b &= \text{passenger } a \text{ choses flight } b \\ c X b &= \text{company } c \text{ operates flight } b \end{aligned}$$

Then (5.158) yields : whenever a choses a flight b it turns out that b is operated by company c . So:

$$c (X/Y) a = \text{company } c \text{ is the only one trusted by passenger } a, \text{ that is, } a \text{ only flies } c.$$

Therefore, (5.157) captures, in a rather eloquent way, the duality between universal and existential quantification. It is no wonder, then, that the relational equivalent to $(x \div y) \times y \leq x$ above is

$$(X/S) \cdot S \subseteq X$$

This *cancellation* rule, very often used in practice, unfolds to

$$\langle \forall b : a S b : c X b \rangle \wedge a S b' \Rightarrow c X b'$$

i.e. to the well-known device in logic known as *modus ponens*: $((S \rightarrow X) \wedge S) \rightarrow X$.

There is one important difference between (5.156) and (5.157): while multiplication in (5.156) is commutative, and thus writing $z \times y$ or $y \times z$ is the same, writing $Z \cdot Y$ or $Y \cdot Z$ makes a lot of difference because composition is not commutative in general. The dual division operator is obtained by taking converses over (5.157):

$$\begin{aligned} Y \cdot Z &\subseteq X \\ &\equiv \{ \text{converses} \} \\ Z^\circ \cdot Y^\circ &\subseteq X^\circ \\ &\equiv \{ \text{division (5.157)} \} \\ Z^\circ &\subseteq X^\circ / Y^\circ \\ &\equiv \{ \text{converses} \} \\ Z &\subseteq \underbrace{(X^\circ / Y^\circ)^\circ}_{Y \setminus X} \end{aligned}$$

In summary:

$$X \cdot Z \subseteq Y \Leftrightarrow Z \subseteq X \setminus Y \quad (5.159)$$

Once variables are added to $Y \setminus X$ we get:

$$a(X \setminus Y)c \equiv \langle \forall b : b X a : b Y c \rangle \quad (5.160)$$

Thus we are ready to add two more rows to table 3:

$(f X) \subseteq Y \equiv X \subseteq (g Y)$			
Description	f	g	Obs.
Left-division	$(R \cdot)$	$(R \setminus)$	read “ R under ...”
Right-division	$(\cdot R)$	$(/ R)$	read “...over R ”

As example of left division consider the relation $a \in x$ between a set x and each of its elements a :

$$A \xleftarrow{\in} PA \quad (5.161)$$

Then inspect the meaning of relation $PA \xleftarrow{\in \setminus \in} PA$ using (5.160):

$$x_1 (\in \setminus \in) x_2 \Leftrightarrow \langle \forall a : a \in x_1 : a \in x_2 \rangle$$

We conclude that quotient $PA \xleftarrow{\in \setminus \in} PA$ expresses the inclusion relation among sets.

Relation division gives rise to a number of combinators in relation algebra that are very useful in problem specification. We review some of these below.

Exercise 5.46. Prove the equalities

$$R \cdot f = R / f^\circ \quad (5.162)$$

$$f \setminus R = f^\circ \cdot R \quad (5.163)$$

$$R / \perp = \top \quad (5.164)$$

$$R / id = R \quad (5.165)$$

$$R \setminus (f^\circ \cdot S) = f \cdot R \setminus S \quad (5.166)$$

$$R \setminus \top \cdot S = ! \cdot R \setminus ! \cdot S \quad (5.167)$$

$$R / (S \cup P) = R / S \cap R / P \quad (5.168)$$

□

SYMMETRIC DIVISION Given two arbitrary relations R and S typed as in the diagram below, define the *symmetric division* $\frac{S}{R}$ of S by R by:

$$b \frac{S}{R} c \equiv \langle \forall a :: a R b \Leftrightarrow a S c \rangle \quad \begin{array}{ccc} & \xleftarrow{\frac{S}{R}} & \\ B & & C \\ & \searrow R \quad \swarrow S & \\ & A & \end{array} \quad (5.169)$$

That is, $b \frac{S}{R} c$ means that b and c are related to exactly the same outputs (in A) by R and by S . Another way of writing (5.169) is $b \frac{S}{R} c \equiv \{a \mid a R b\} = \{a \mid a S c\}$ which is the same as

$$b \frac{S}{R} c \equiv \Lambda R b = \Lambda S c \quad (5.170)$$

where Λ is the *power transpose* operator²⁵ which maps a relation $Q : Y \leftarrow X$ to the set valued function $\Lambda Q : X \rightarrow \mathcal{P} Y$ such that $\Lambda Q x = \{y \mid y Q x\}$. Another way to define $\frac{S}{R}$ is

$$\frac{S}{R} = R \setminus S \cap R^\circ / S^\circ \quad (5.171)$$

which factors symmetric division into the two asymmetric divisions $R \setminus S$ (5.159) and R / S (5.157) already studied above. Moreover, for $R, S := f, g$, definition (5.171) instantiates to $\frac{f}{g}$ as defined by (5.49). By (5.159, 5.157), (5.171) is equivalent to the universal property:

$$X \subseteq \frac{S}{R} \equiv R \cdot X \subseteq S \wedge S \cdot X^\circ \subseteq R \quad (5.172)$$

From the definitions above a number of standard properties arise:

$$\left(\frac{S}{R}\right)^\circ = \frac{R}{S} \quad (5.173)$$

$$\frac{S}{R} \cdot \frac{Q}{S} \subseteq \frac{Q}{R} \quad (5.174)$$

$$f^\circ \cdot \frac{S}{R} \cdot g = \frac{S \cdot g}{R \cdot f} \quad (5.175)$$

$$id \subseteq \frac{R}{R} \quad (5.176)$$

Thus $\frac{R}{R}$ is always an *equivalence relation*, for any given R . Furthermore,

$$R = \frac{R}{R} \equiv R \text{ is an equivalence relation} \quad (5.177)$$

holds. Also note that, even in the case of functions, (5.174) remains an inclusion,

$$\frac{f}{g} \cdot \frac{h}{f} \subseteq \frac{h}{g} \quad (5.178)$$

since:

$$\begin{aligned} & \frac{f}{g} \cdot \frac{h}{f} \subseteq \frac{h}{g} \\ \Leftrightarrow & \{ \text{factor } \frac{id}{g} \text{ out} \} \\ & f \cdot \frac{h}{f} \subseteq h \\ \Leftrightarrow & \{ \text{factor } h \text{ out} \} \end{aligned}$$

²⁵ See section 5.24 for more details about this operator.

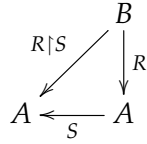
$$\begin{aligned}
& f \cdot \frac{id}{f} \subseteq id \\
& \equiv \{ \text{shunting rule (5.47)} \} \\
& f \subseteq f \\
& \equiv \{ \text{trivial} \} \\
& true \\
& \square
\end{aligned}$$

From (5.178) it follows that $\frac{f}{f}$ is always transitive. By (5.173) it is symmetric and by (5.30) it is reflexive. Thus $\frac{f}{f}$ is an *equivalence relation*.

RELATION SHRINKING Given relations $R : A \leftarrow B$ and $S : A \leftarrow A$, define $R \upharpoonright S : A \leftarrow B$, pronounced “ R shrunk by S ”, by

$$X \subseteq R \upharpoonright S \equiv X \subseteq R \wedge X \cdot R^\circ \subseteq S \quad (5.179)$$

cf. diagram:



This states that $R \upharpoonright S$ is the largest part of R such that, if it yields an output for an input x , it must be a maximum, with respect to S , among all possible outputs of x by R . By indirect equality, (5.179) is equivalent to the closed definition:

$$R \upharpoonright S = R \cap S / R^\circ \quad (5.180)$$

(5.179) can be regarded as a Galois connection between the set of all *subrelations* of R and the set of *optimization criteria* (S) on its outputs.

Combinator $R \upharpoonright S$ also makes sense when R and S are finite, relational data structures (eg. tables in a database). Consider, for instance, the following example of $R \upharpoonright S$ in a *data-processing* context: given

<i>Examiner</i>	<i>Mark</i>	<i>Student</i>
<i>Smith</i>	10	<i>John</i>
<i>Smith</i>	11	<i>Mary</i>
<i>Smith</i>	15	<i>Arthur</i>
<i>Wood</i>	12	<i>John</i>
<i>Wood</i>	11	<i>Mary</i>
<i>Wood</i>	15	<i>Arthur</i>

and wishing to “choose the best mark” for each student, project over *Mark, Student* and optimize over the \geq ordering on *Mark*:

$$\left(\begin{array}{c|c} \textit{Mark} & \textit{Student} \\ \hline 10 & \textit{John} \\ 11 & \textit{Mary} \\ 12 & \textit{John} \\ 15 & \textit{Arthur} \end{array} \right) \upharpoonright_{\geq} = \left(\begin{array}{c|c} \textit{Mark} & \textit{Student} \\ \hline 11 & \textit{Mary} \\ 12 & \textit{John} \\ 15 & \textit{Arthur} \end{array} \right)$$

Relational shrinking can be used in many other contexts. Consider, for instance, a sensor recording temperatures (T), $T \xleftarrow{S} \mathbb{N}_0$, where data in \mathbb{N}_0 are “time stamps”. Suppose one wishes to filter out repeated temperatures, keeping the first occurrences only. This can be specified by:

$$T \xleftarrow{nub S} \mathbb{N}_0 = (S^\circ \upharpoonright \leq)^\circ$$

That is, *nub* is the function that removes all duplicates while keeping the first instances.

Among the properties of shrinking [43] we single out the two *fusion* rules:

$$(S \cdot f) \upharpoonright R = (S \upharpoonright R) \cdot f \quad (5.181)$$

$$(f \cdot S) \upharpoonright R = f \cdot (S \upharpoonright (f^\circ \cdot R \cdot f)) \quad (5.182)$$

Some more basic properties are: “chaotic optimization”,

$$R \upharpoonright \top = R \quad (5.183)$$

“impossible optimization”

$$R \upharpoonright \perp = \perp \quad (5.184)$$

and “brute force” determinization:

$$R \upharpoonright id = \text{largest deterministic fragment of } R \quad (5.185)$$

$R \upharpoonright id$ is the extreme case of the fact which follows:

$$R \upharpoonright S \text{ is simple} \Leftarrow S \text{ is anti-symmetric} \quad (5.186)$$

Thus anti-symmetric criteria always lead to determinism, possibly at the sacrifice of totality. Also, for R simple:

$$R \upharpoonright S = R \quad \equiv \quad \text{img } R \subseteq S \quad (5.187)$$

Thus (for functions):

$$f \upharpoonright S = f \quad \Leftarrow \quad S \text{ is reflexive} \quad (5.188)$$

The distribution of shrinking by join,

$$(R \cup S) \upharpoonright Q = (R \upharpoonright Q) \cap Q/S^\circ \cup (S \upharpoonright Q) \cap Q/R^\circ \quad (5.189)$$

has a number of corollaries, namely a *conditional rule*,

$$(p \rightarrow R, T) \upharpoonright S = p \rightarrow (R \upharpoonright S), (p \upharpoonright S) \quad (5.190)$$

the *distribution* over alternatives (5.114),

$$[R, S] \upharpoonright U = [R \upharpoonright U, S \upharpoonright U] \quad (5.191)$$

and the “*function competition*” rule:

$$(f \cup g) \upharpoonright S = (f \cap S \cdot g) \cup (g \cap S \cdot f) \quad (5.192)$$

(Recall that $S/g^\circ = S \cdot g$.)

Putting universal properties (5.172,5.179) together we get, by indirect equality,

$$\frac{R}{g} = g^\circ \cdot (R \upharpoonright id) \quad (5.193)$$

$$\frac{f}{R} = (R \upharpoonright id)^\circ \cdot f \quad (5.194)$$

capturing a relationship between shrinking and symmetric division: knowing that $R \upharpoonright id$ is the deterministic fragment of R , we see how the *vagueness* of arbitrary R replacing either f or g in $\frac{f}{g}$ is forced to shrink.

Exercise 5.47. Use shrinking and other relational combinators to calculate, from a relation of type (5.127), the relation of type $Student \times Course \rightarrow Result$ that tells the final results of all exams. (**NB:** assume $Time = \mathbb{N}_0$ ordered by (\leq) .)

□

RELATION OVERRIDING Another operator enabled by relation division is the relational *overriding* combinator,

$$R \dagger S = S \cup R \cap \perp / S^\circ \quad (5.195)$$

which yields the relation which contains the whole of S and that part of R where S is undefined — read $R \dagger S$ as “ R overridden by S ”.

It is easy to show that $\perp \dagger S = S$, $R \dagger \perp = R$ and $R \dagger R = R$ hold. From (5.195) we derive, by indirect equality, the universal property:

$$X \subseteq R \dagger S \equiv X \subseteq R \cup S \wedge (X - S) \cdot S^\circ = \perp \quad (5.196)$$

The following property establishes a relationship between overriding and the McCarthy conditional:

$$p \rightarrow g, f = f \dagger (g \cdot \Phi_p) \quad (5.197)$$

Notation Φ_p is explained in the next section.

Exercise 5.48. Show that

$$R \dagger f = f$$

holds, arising from (5.196,5.138) — where f is a function, of course.

□

Exercise 5.49. On June 23rd, 1991, E.W. Dijkstra wrote one of his famous notes — EWD1102-5 — entitled: “Why preorders are beautiful”. The main result of his six page long manuscript is:

A binary relation is a pre-order iff $R = R / R$ holds.

The proof of this result becomes even shorter (and perhaps even more beautiful) once expressed in relation algebra. Fill in the ellipses in the following calculation of such a result:

$$\begin{aligned}
& R = R / R \\
\equiv & \quad \{ \dots\dots\dots \} \\
& \left\{ \begin{array}{l} X \subseteq R \Leftrightarrow X \cdot R \subseteq R \\ X \subseteq R \Leftrightarrow X \cdot R \subseteq R \end{array} \right. \\
\Rightarrow & \quad \{ \dots\dots\dots \} \\
& \left\{ \begin{array}{l} id \subseteq R \Leftrightarrow R \subseteq R \\ R \subseteq R \Leftrightarrow R \cdot R \subseteq R \end{array} \right. \\
\equiv & \quad \{ \dots\dots\dots \} \\
& \left\{ \begin{array}{l} id \subseteq R \\ R \cdot R \subseteq R \end{array} \right. \\
\equiv & \quad \{ \dots\dots\dots \} \\
& \left\{ \begin{array}{l} id \subseteq R \wedge (R / R) \cdot R \subseteq R \\ R \subseteq R / R \end{array} \right. \\
\Rightarrow & \quad \{ \dots\dots\dots \} \\
& \left\{ \begin{array}{l} R / R \subseteq R \\ R \subseteq R / R \end{array} \right. \\
\equiv & \quad \{ \dots\dots\dots \} \\
& R = R / R \\
& \square
\end{aligned}$$

That is,

$$R = R / R \equiv \left\{ \begin{array}{l} id \subseteq R \\ R \cdot R \subseteq R \end{array} \right.$$

□

5.20 PREDICATES ALSO BECOME RELATIONS

Recall from (5.49) the notation $\frac{f}{g} = g^\circ \cdot f$ and define, given a predicate $p : A \rightarrow \mathbb{B}$, the relation $\Phi_p : A \rightarrow A$ as follows:²⁶

$$\Phi_p = id \cap \frac{true}{p} \quad (5.198)$$

By (5.49), Φ_p is the *coreflexive* relation which represents predicate p as a binary relation,

$$y \Phi_p x \Leftrightarrow y = x \wedge p x \quad (5.199)$$

²⁶ Recall that *true* is the *constant* function yielding True for every argument (5.40).

as can be easily checked. From (5.198) one gets the limit situations:²⁷

$$\Phi_{true} = id \quad (5.200)$$

$$\Phi_{false} = \perp \quad (5.201)$$

Moreover,

$$\Phi_{p \wedge q} = \Phi_p \cap \Phi_q \quad (5.202)$$

$$\Phi_{p \vee q} = \Phi_p \cup \Phi_q \quad (5.203)$$

$$\Phi_{\neg p} = id - \Phi_p \quad (5.204)$$

follow immediately from (5.199) and from (5.39) one infers $\frac{true}{p} \cdot R \subseteq \frac{true}{p}$ for any R . In particular, $\frac{true}{p} \cdot \top = \frac{true}{p}$ since $\frac{true}{p} \subseteq \frac{true}{p} \cdot \top$ always holds. Then, by distributive property (5.62):

$$\Phi_p \cdot \top = \frac{true}{p} \quad (5.205)$$

Moreover, the following two properties hold:

$$R \cdot \Phi_p = R \cap \top \cdot \Phi_p \quad (5.206)$$

$$\Phi_q \cdot R = R \cap \Phi_q \cdot \top \quad (5.207)$$

We check (5.207):²⁸

$$\begin{aligned} & \Phi_q \cdot R \\ = & \{ (5.109) ; (5.198) \} \\ & \frac{id \vee true}{id \vee p} \cdot R \\ = & \{ (5.104) \} \\ & (id \vee p)^\circ \cdot (R \vee true) \\ = & \{ (5.108) \} \\ & R \cap \frac{true}{p} \\ = & \{ (5.205) \} \\ & R \cap \Phi_p \cdot \top \\ & \square \end{aligned}$$

Note the meaning of (5.206) and (5.207):

$$b (R \cdot \Phi_p) a \Leftrightarrow b R a \wedge (p a)$$

$$b (\Phi_q \cdot R) a \Leftrightarrow b R a \wedge (q b)$$

So (5.206) — resp. (5.207) — restricts R to inputs satisfying p — resp. outputs satisfying q .

²⁷ $\Phi_{false} = \perp$ arises from (5.54) since $\text{True} \neq \text{False}$.

²⁸ The other is obtained from (5.207) by taking converses.

Below we show how to use relation restriction and overriding in specifying the operation that, in the Alcuin puzzle — recall (5.74)

$$\begin{array}{ccc} \text{Being} & \xrightarrow{\text{Eats}} & \text{Being} \\ & \downarrow \text{where} & \\ & \text{Bank} & \xrightarrow{\text{cross}} \text{Bank} \end{array}$$

— specifies the move of *Beings* to the other bank:

$$\text{carry who where} = \text{where} \dagger (\text{cross} \cdot \text{where} \cdot \Phi_{\in \text{who}})$$

By (5.197) this simplifies to a McCarthy conditional:

$$\text{carry who where} = (\in \text{who}) \rightarrow \text{cross} \cdot \text{where} , \text{ where} \quad (5.208)$$

In pointwise notation, *carry* is the function:

$$\begin{aligned} \text{carry who where } b = \\ \text{if } b \in \text{who} \text{ then } \text{cross } m \text{ else } m \\ \text{where } m = \text{where } b \end{aligned}$$

Note the type $\text{carry} : \text{PBeing} \rightarrow \text{Bank}^{\text{Being}} \rightarrow \text{Bank}^{\text{Being}}$.

A notable property of coreflexive relations is that their composition coincides with their meet:

$$\Phi_q \cdot \Phi_p = \Phi_q \cap \Phi_p \quad (5.209)$$

In consequence, composing a coreflexive with itself yields that very same coreflexive: $\Phi_p \cdot \Phi_p = \Phi_p$. (5.209) follows from (5.206,5.207):

$$\begin{aligned} & \Phi_q \cdot \Phi_p \\ = & \{ R = R \cap R \} \\ & \Phi_q \cdot \Phi_p \cap \Phi_q \cdot \Phi_p \\ = & \{ (5.206, 5.207) \} \\ & \Phi_q \cap \top \cdot \Phi_p \cap \Phi_p \cap \Phi_p \cdot \top \\ = & \{ \text{since } \Phi_p \subseteq \top \cdot \Phi_p \text{ and } \Phi_q \subseteq \Phi_q \cdot \top \} \\ & \Phi_q \cap \Phi_p \\ & \square \end{aligned}$$

EQUALIZERS The definition of Φ_p (5.187) can be regarded as a particular case of an *equalizer*: given two functions $B \xleftarrow{f,g} A$, the equalizer of f and g is the relation $eq(f, g) = id \cap \frac{f}{g}$. By indirect equality,

$$X \subseteq eq(f, g) \Leftrightarrow X \subseteq id \wedge g \cdot X \subseteq f$$

That is, $eq(f, g)$ is the largest coreflexive X that restricts g so that f and g yield the same outputs.

Clearly, $eq(f, f) = id$. Note that an equalizer can be empty, cf. e.g. $eq(true, false) = \perp$.

Exercise 5.50. Based on (5.71) show that

$$g^\circ \cdot \Phi_p \cdot f = \frac{f}{g} \cap \frac{true}{p \cdot g} \quad (5.210)$$

holds.²⁹

□

5.21 GUARDS, COREFLEXIVES AND THE MCCARTHY CONDITIONAL

From the definition of a McCarthy conditional (2.70) we obtain $p? = p \rightarrow i_1, i_2$ and then $p? = i_2 \dagger i_1 \cdot \Phi_p$ by (5.197). A third way to express the guard $p?$ is

$$p? = i_1 \cdot \Phi_p \cup i_2 \cap (\perp / (i_1 \cdot \Phi_p)^\circ) \quad (5.211)$$

by (5.195), which simplifies to:

$$p? = [\Phi_p, \Phi_{\neg p}]^\circ \quad (5.212)$$

To prove (5.212) note that $\perp / (i_1 \cdot \Phi_p)^\circ = \perp / \Phi_p$, immediate by the laws of S / R and shunting. Then, $\perp / \Phi_p = \top \cdot \Phi_{\neg p}$. Here one only needs to check:

$$\begin{aligned} & \perp / \Phi_p \subseteq \top \cdot \Phi_{\neg p} \\ \equiv & \quad \left\{ \frac{\neg p}{true} = \frac{p}{false} \right\} \\ & \perp / \Phi_p \subseteq \frac{p}{false} \\ \equiv & \quad \{ \text{going pointwise} \} \\ & \langle \forall y, x : y (\perp / \Phi_p) x : p x = \text{False} \rangle \\ \equiv & \quad \{ (5.159); (5.199) \} \\ & \langle \forall y, x : p x \Rightarrow \text{False} : p x = \text{False} \rangle \\ \equiv & \quad \{ \text{trivial} \} \\ & true \\ & \square \end{aligned}$$

Finally, back to (5.211):

$$\begin{aligned} p? &= i_1 \cdot \Phi_p \cup i_2 \cap \top \cdot \Phi_{\neg p} \\ \equiv & \quad \{ (5.206); \text{converses} \} \end{aligned}$$

²⁹ Both sides of the equality mean $g b = f a \wedge p (g b)$.

$$\begin{aligned}
(p?)^\circ &= \Phi_p \cdot i_1^\circ \cup \Phi_{\neg p} \cdot i_2^\circ \\
&\equiv \{ (5.121) \} \\
p? &= [\Phi_p, \Phi_{\neg p}]^\circ \\
&\square
\end{aligned}$$

Exercise 5.51. From (5.211) infer

$$p \rightarrow R, S = R \cap \frac{p}{\text{true}} \cup S \cap \frac{p}{\text{false}} \quad (5.213)$$

and therefore $p \rightarrow R, S \subseteq R \cup S$. Furthermore, derive (2.78) from (5.213) knowing that $\text{true} \cup \text{false} = \top$.

□

DOMAIN AND RANGE Suppose one computes $\ker \langle R, id \rangle$ instead of $\ker R$. Since $\ker \langle R, id \rangle = \ker R \cap id$ (5.111), coreflexive relation is obtained. This is called the *domain* of R , written:

$$\delta R = \ker \langle R, id \rangle \quad (5.214)$$

Since³⁰

$$\top \cdot R \cap id = R^\circ \cdot R \cap id \quad (5.215)$$

domain can be also defined by

$$\delta R = \top \cdot R \cap id \quad (5.216)$$

Dually, one can define the *range* of R as the domain of its converse:

$$\rho R = \delta R^\circ = \text{img } R \cap id \quad (5.217)$$

For functions, range and image coincide, since $\text{img } f \subseteq id$ for any f . For injective relations, domain and kernel coincide, since $\ker R \subseteq id$ in such situations. These two operators can be shown to be characterized by two Galois connections, as follows:

$(f X) \subseteq Y \equiv X \subseteq (g Y)$			
Description	f	g	Obs.
domain	δ	$(\top \cdot)$	left \subseteq restricted to coreflexives
range	ρ	$(\cdot \top)$	left \subseteq restricted to coreflexives

Let us show that indeed

$$\delta X \subseteq Y \equiv X \subseteq \top \cdot Y \quad (5.218)$$

$$\rho R \subseteq Y \equiv R \subseteq Y \cdot \top \quad (5.219)$$

³⁰ (5.215) follows from $id \cap \top \cdot R \subseteq R^\circ \cdot R$ which can be easily checked pointwise.

hold, where variable Y ranges over coreflexive relations only. We only derive (5.218), from which (5.219) is obtained taking converses. We rely on the definition just given and on previously defined connections:

$$\begin{aligned}
 & \delta X \subseteq Y \\
 \equiv & \quad \{ (5.216) \} \\
 & \top \cdot X \cap id \subseteq Y \\
 \equiv & \quad \{ \text{two Galois connections} \} \\
 & X \subseteq \top \setminus (id \Rightarrow Y) \\
 \equiv & \quad \{ \top \setminus (id \Rightarrow Y) = \top \cdot Y, \text{ see below} \} \\
 & X \subseteq \top \cdot Y \\
 & \square
 \end{aligned}$$

To justify the hint above, first note that $\top \cdot Y \subseteq id \Rightarrow Y$, for Y coreflexive — recall (5.198) and (5.205). Then:

$$\begin{aligned}
 & \top \setminus (id \Rightarrow Y) \subseteq \top \cdot Y \\
 \Leftarrow & \quad \{ \text{monotonicity ; rule “raise-the-lower-side”} \} \\
 & \top \setminus (\top \cdot Y) \subseteq \top \cdot Y \\
 \equiv & \quad \{ (5.167) ; f \cdot f^\circ \cdot f = f \text{ for } f := ! \text{ (twice)} \} \\
 & ! \setminus ! \cdot Y \subseteq \top \cdot Y \\
 \equiv & \quad \{ f \setminus R = f^\circ \cdot R ; \top = \ker ! \} \\
 & \top \cdot Y \subseteq \top \cdot Y \\
 & \square
 \end{aligned}$$

Note the left-cancellation rule of the δ connection:

$$R \subseteq \top \cdot \delta R \quad (5.220)$$

From this the following domain/range elimination rules follow:

$$\top \cdot \delta R = \top \cdot R \quad (5.221)$$

$$\rho R \cdot \top = R \cdot \top \quad (5.222)$$

$$\delta R \subseteq \delta S \equiv R \subseteq \top \cdot S \quad (5.223)$$

Proof of (5.221):

$$\begin{aligned}
 & \top \cdot \delta R = \top \cdot R \\
 \equiv & \quad \{ \text{circular inclusion} \} \\
 & \top \cdot \delta R \subseteq \top \cdot R \wedge \top \cdot R \subseteq \top \cdot \delta R \\
 \equiv & \quad \{ (5.97) \text{ twice} \} \\
 & \delta R \subseteq \top \cdot R \wedge R \subseteq \top \cdot \delta R \\
 \equiv & \quad \{ \text{cancelation (5.220)} \}
 \end{aligned}$$

$$\begin{aligned}
& \delta R \subseteq \top \cdot R \\
& \equiv \{ \delta R = \top \cdot R \cap id \text{ (5.216)} \} \\
& \text{true} \\
& \square
\end{aligned}$$

Rule (5.222) follows by dualization (converses) and (5.223) follows from (5.218) and (5.221). More facts about domain and range:

$$\delta (R \cdot S) = \delta (\delta R \cdot S) \quad (5.224)$$

$$\rho (R \cdot S) = \rho (R \cdot \rho S) \quad (5.225)$$

$$R = R \cdot \delta R \quad (5.226)$$

$$R = \rho R \cdot R \quad (5.227)$$

Last but not least: given predicate q and function f ,

$$\Phi_{(q \cdot f)} = \delta (\Phi_q \cdot f) \quad (5.228)$$

holds. Proof:

$$\begin{aligned}
& \Phi_{(q \cdot f)} \\
& = \{ (5.198) \} \\
& \quad id \cap \frac{true}{q \cdot f} \\
& = \{ \text{since } \frac{f}{f} \text{ is reflexive (5.30)} \} \\
& \quad id \cap \frac{f}{f} \cap \frac{true \cdot f}{q \cdot f} \\
& = \{ (5.109) ; \text{products} \} \\
& \quad id \cap \frac{(id \vee true) \cdot f}{(id \vee q) \cdot f} \\
& = \{ (5.49) ; (5.109) \} \\
& \quad id \cap f^\circ \cdot (id \cap \frac{true}{q}) \cdot f \\
& = \{ (5.198) \} \\
& \quad id \cap f^\circ \cdot \Phi_q \cdot f \\
& = \{ \delta R = id \cap R^\circ \cdot R \} \\
& \quad \delta (\Phi_q \cdot f) \\
& \square
\end{aligned}$$

Exercise 5.52. Recalling (5.206), (5.207) and other properties of relation algebra, show that: (a) (5.218) and (5.219) can be re-written with R replacing \top ; (b) $\Phi \subseteq \Psi \equiv ! \cdot \Phi \subseteq ! \cdot \Psi$ holds.³¹

³¹ Thus coreflexives can be represented by *vectors* and vice-versa.

□

5.22 DIFUNCTIONALITY

A relation R is said to be *difunctional* or *regular* wherever $R \cdot R^\circ \cdot R = R$ holds, which amounts to $R \cdot R^\circ \cdot R \subseteq R$ since the converse inclusion always holds.

The class of difunctional relations is vast. \top and \perp are difunctional, and so are all coreflexive relations, as is easy to check. It also includes all simple relations, since $R \cdot R^\circ = \text{img } R \subseteq \text{id}$ wherever R is simple. Moreover, divisions of functions are difunctional because every symmetric division is so, as is easy to check by application of laws (5.174) and (5.173):

$$\begin{aligned}
 & \frac{f}{g} \cdot \left(\frac{f}{g} \right)^\circ \cdot \frac{f}{g} \subseteq \frac{f}{g} \\
 \Leftarrow & \quad \{ (5.51); (5.178) \} \\
 & \frac{f}{g} \cdot \frac{f}{f} \subseteq \frac{f}{g} \\
 \Leftarrow & \quad \{ (5.178) \} \\
 & \frac{f}{g} \subseteq \frac{f}{g} \\
 & \square
 \end{aligned}$$

For $g = \text{id}$ above we get that any function f being difunctional can be expressed by $f \cdot \frac{f}{f} = f$.

Recall that an equivalence relation can always be represented by the kernel of some function, typically by $R = \frac{\Delta R}{\Delta R}$. So equivalence relations are difunctional. The following rule is of practical relevance:

A difunctional relation that is reflexive and symmetric necessarily is transitive, and therefore an equivalence relation.

Proof (of transitivity):

$$\begin{aligned}
 & R \cdot R \subseteq R \\
 \equiv & \quad \{ R \text{ is difunctional} \} \\
 & R \cdot R \subseteq R \cdot R^\circ \cdot R \\
 \equiv & \quad \{ R \text{ is symmetric} \} \\
 & R \cdot R \subseteq R \cdot R \cdot R \\
 \Leftarrow & \quad \{ \text{monotonicity} \} \\
 & \text{id} \subseteq R \\
 & \square
 \end{aligned}$$

Difunctional relations are also called *regular*, *rational* or *uniform*. First, some intuition about what “regularity” means: a regular (difunctional) relation is such that, wherever two inputs have a common image, then they have *exactly the same* set of images. In other words, the image sets of two different inputs are either disjoint or the same. As a counterexample, take the following relation, represented as a matrix with inputs taken from set $\{a_1, \dots, a_5\}$ and outputs delivered into set $\{b_1, \dots, b_5\}$:

R	a_1	a_2	a_3	a_4	a_5	
b_1	0	0	1	0	1	
b_2	0	0	0	0	0	
b_3	0	1	0	0	0	
b_4	0	1	0	1	0	
b_5	0	0	0	1	0	

(5.229)

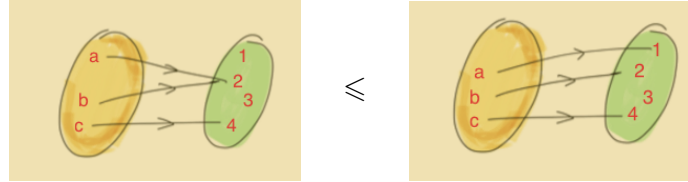
Concerning inputs a_3 and a_5 , regularity holds; but sets $\{b_3, b_4\}$ and $\{b_4, b_5\}$ — the images of a_2 and a_4 , respectively — are neither disjoint nor the same: so R isn’t regular. It would become so if e.g. b_4 were dropped from both image sets or one of b_3 or b_5 were replaced for the other in the corresponding image set.

5.23 OTHER ORDERINGS ON RELATIONS

THE INJECTIVITY PREORDER The kernel relation $\ker R = R^\circ \cdot R$ measures the level of *injectivity* of R according to the preorder

$$R \leq S \equiv \ker S \subseteq \ker R \quad (5.230)$$

telling that R is *less injective* or *more defined* (entire) than S . For instance:



This ordering is surprisingly useful in formal specification because of its properties. For instance, it is upper-bounded by relation *pairing*, recall (5.103):

$$\langle R, S \rangle \leq X \equiv R \leq X \wedge S \leq X \quad (5.231)$$

Cancellation of (5.231) means that *pairing* always *increases injectivity*:

$$R \leq \langle R, S \rangle \text{ and } S \leq \langle R, S \rangle. \quad (5.232)$$

(5.232) unfolds to $\ker \langle R, S \rangle \subseteq (\ker R) \cap (\ker S)$, confirming (5.111). The following injectivity *shunting law* arises as a Galois connection:

$$R \cdot g \leq S \equiv R \leq S \cdot g^\circ \quad (5.233)$$

Restricted to *functions*, (\leq) is *universally* bounded by

$$! \leq f \leq id$$

where (recall) $1 \xleftarrow{!} A$ is the unique function of its type, where 1 is the singleton type. Moreover,

- A function is *injective* iff $id \leq f$. Thus $\langle f, id \rangle$ is always *injective* (5.232).
- Two functions $f \in g$ are said to be *complementary* wherever $id \leq \langle f, g \rangle$.
- Any relation R can be factored into the composition $f \cdot g^\circ$ of two complementary functions f and g .³²

For instance, the *projections* $\pi_1 (a, b = a)$, $\pi_2 (a, b = b)$ are complementary since $\langle \pi_1, \pi_2 \rangle = id$ (2.32).

As illustration of the use of this ordering in formal specification, suppose one writes

$$room \leq \langle lect, slot \rangle$$

in the context of the data model

$$\begin{array}{ccccc} Teacher & \xleftarrow{lect} & Class & \xrightarrow{room} & Room \\ & & \downarrow slot & & \\ & & TD & & \end{array}$$

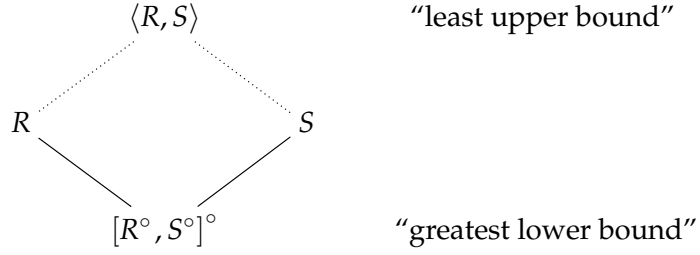
where TD abbreviates time and date. What are we telling about this model by writing $room \leq \langle lect, slot \rangle$? We unfold this constraint in the expected way:

$$\begin{aligned} & room \leq \langle lect, slot \rangle \\ \equiv & \{ (5.230) \} \\ & \ker \langle lect, slot \rangle \subseteq \ker room \\ \equiv & \{ (5.111); (5.53) \} \\ & \frac{lect}{lect} \cap \frac{slot}{slot} \subseteq \frac{room}{room} \\ \equiv & \{ \text{going pointwise, for all } c_1, c_2 \in Class \} \\ & (lect\ c_1 = lect\ c_2 \wedge slot\ c_1 = slot\ c_2) \Rightarrow (room\ c_1 = room\ c_2) \end{aligned}$$

This $room \leq \langle lect, slot \rangle$ constrains the model in the sense of imposing that a given lecturer cannot be in two different rooms at the same time. c_1 and c_2 are classes shared by different courses, possibly of different degrees. In the standard terminology of database theory this is called a *functional dependency*, see exercises 5.55 and 5.56 in the sequel.

³² This remarkable factorization is known as a *tabulation* of R [10].

Interestingly, the injectivity preorder not only has least upper bounds but also greatest lower bounds,



that is,

$$X \leq [R^\circ, S^\circ]^\circ \Leftrightarrow X \leq R \wedge X \leq S \quad (5.234)$$

as the calculation shows:

$$\begin{aligned} & X \leq [R^\circ, S^\circ]^\circ \\ \equiv & \quad \{ \text{injectivity preorder ; } \ker R^\circ = \text{img } R \} \\ & \text{img } [R^\circ, S^\circ] \subseteq \ker X \\ \equiv & \quad \{ (5.124) \} \\ & R^\circ \cdot R \cup S^\circ \cdot S \subseteq \ker X \\ \equiv & \quad \{ \text{kernel; } \cdot \cup \text{--universal} \} \\ & \ker R \subseteq \ker X \wedge \ker S \subseteq \ker X \\ \equiv & \quad \{ \text{injectivity preorder (twice)} \} \\ & X \leq R \wedge X \leq S \\ & \square \end{aligned}$$

Note the meaning of the glb of R and S ,

$$x [R^\circ, S^\circ]^\circ a \Leftrightarrow \langle \exists b : x = i_1 b : b R a \rangle \vee \langle \exists c : x = i_2 c : c R a \rangle$$

since $[R^\circ, S^\circ]^\circ = i_1 \cdot R \cup i_2 \cdot S$. This is the most injective relation that is less injective than R and S because it just “collates” the outputs of both relations without confusing them.³³

Exercise 5.53. Show that $R^\circ \cdot S = \perp$ is necessary for the coproduct of two injective relations R and S to be injective:

$$id \leq [R, S] \Leftrightarrow id \leq R \wedge id \leq S \wedge R^\circ \cdot S = \perp \quad (5.235)$$

□

³³ It turns out that universal property $X = [R^\circ, S^\circ]^\circ \Leftrightarrow i_1^\circ \cdot X = R \wedge i_2^\circ \cdot X = S$ holds, as is easy to derive from (5.114). So $[R^\circ, S^\circ]^\circ$ is the *categorical* product for relations:

$$A \rightarrow (B + C) \xrightleftharpoons{\cong} (A \rightarrow B) \times (A \rightarrow C)$$

That is, among relations, the product is obtained as the converse dual of the coproduct. This is called a *biproduct* [33].

Exercise 5.54. The Peano algebra $\mathbb{N}_0 \xleftarrow{\text{in}} 1 + \mathbb{N}_0 = [0, \text{succ}]$ is an isomorphism³⁴, and therefore injective. Check what (5.235) means in this case.

□

Exercise 5.55. An SQL-like relational operator is projection,

$$\pi_{g,f}R \stackrel{\text{def}}{=} g \cdot R \cdot f^\circ \quad \begin{array}{ccc} B & \xleftarrow{R} & A \\ g \downarrow & & \downarrow f \\ C & \xleftarrow{\pi_{g,f}R} & D \end{array} \quad (5.236)$$

whose set-theoretic meaning is³⁵

$$\pi_{g,f}R = \{(g\,b, f\,a) \mid b \in B \wedge a \in A \wedge b\,R\,a\} \quad (5.237)$$

Functions f and g are often referred to as **attributes** of R . Derive (5.237) from (5.236).

□

Exercise 5.56. A relation R is said to satisfy functional dependency (FD) $g \rightarrow f$, written $g \xrightarrow{R} f$ wherever projection $\pi_{f,g}R$ (5.236) is simple.

1. Recalling (5.230), prove the equivalence:

$$g \xrightarrow{R} f \quad \equiv \quad f \leq g \cdot R^\circ \quad (5.238)$$

2. Show that $g \xrightarrow{R} f$ trivially holds wherever g is injective and R is simple, for all (suitably typed) f .

3. Prove the composition rule of FDs:

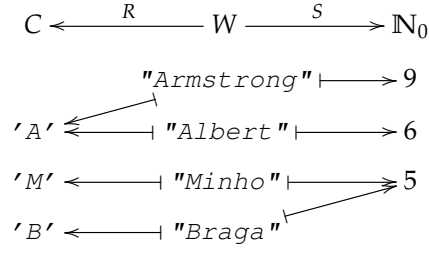
$$h \xleftarrow{S \cdot R} g \quad \Leftarrow \quad h \xleftarrow{S} f \quad \wedge \quad f \xleftarrow{R} g \quad (5.239)$$

□

³⁴ Recall section 3.1.

³⁵ Note that any relation $R: B \leftarrow A$ defines the set of pairs $\{(b, a) \mid b\,R\,a\}$. Predicate $b\,R\,a$ describes R *intensionally*. The set $\{(b, a) \mid b\,R\,a\}$ is the *extension* of R .

Exercise 5.57. Let R and S be the two relations depicted as follows:



Check the assertions:

1. $R \leq S$
2. $S \leq R$
3. Both hold
4. None holds.

□

Exercise 5.58. As follow up to exercise 5.9,

- specify the relation R between students and teachers such that $t R s$ means: t is the mentor of s and also teaches one of her/his courses.
- Specify the property: mentors of students necessarily are among their teachers.

□

THE DEFINITION PREORDER The injectivity preorder works perfectly for functions, which are entire relations. For non-entire R it behaves in a mixed way, measuring not only injectivity but also definition (entireness). It is useful to order relations with respect to how defined they are:

$$R \preceq S \equiv \delta R \subseteq \delta S \quad (5.240)$$

From $\top = \ker !$ one draws another version of (5.240), $R \preceq S \equiv ! \cdot R \subseteq ! \cdot S$. The following Galois connections

$$R \cup S \preceq T \equiv R \preceq T \wedge S \preceq T \quad (5.241)$$

$$R \cdot f^\circ \preceq S \equiv R \preceq S \cdot f \quad (5.242)$$

are easy to prove. Recalling (5.223), (5.240) can also be written

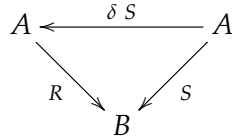
$$\delta R \subseteq \delta S \equiv R \subseteq \top \cdot S \quad (5.243)$$

THE REFINEMENT ORDER Standard programming theory relies on a notion of program *refinement*. As a rule, the starting point in any software design is a so-called *specification*, which indicates the expected behaviour of the program to be developed with no indication of *how* outputs are computed from the inputs. So, “vagueness” is a chief ingredient of a good specification, giving freedom to the programmer to choose a particular algorithmic solution.

Relation algebra captures this by ordering relations with respect to the degree in which they are closer to implementations:

$$S \vdash R \equiv \delta S \subseteq \delta R \wedge R \cdot \delta S \subseteq S \quad (5.244)$$

cf.



$S \vdash R$ is read: “ S is refined by R ”. In the limit situation, R is a function f , and then

$$S \vdash f \Leftrightarrow \delta S \subseteq f^\circ \cdot S \quad (5.245)$$

by shunting (5.46). This is a limit in the sense that f can be neither more defined nor more deterministic.

As maxima of the refinement ordering, functions are regarded as implementations “*par excellence*”. Note how (5.245) captures *implicit specification* S being refined by some function f — recall section 5.3. Back to points and thanks to (5.17) we obtain, in classical “VDM-speak”:

$$\forall a. \text{pre-}S(a) \Rightarrow \text{post-}S(f\ a, a)$$

In case S is entire, (5.245) simplifies to $S \vdash f \Leftrightarrow f \subseteq S$. As example of this particular case we go back to section 5.3 and prove that *abs*, explicitly defined by $\text{abs } i = \text{if } i < 0 \text{ then } -i \text{ else } i$, meets the implicit specification given there, here encoded by $S = \frac{\text{true}}{\text{geq0}} \cap (id \cup \text{sym})$ where $\text{geq0 } x = x \geq 0$ and $\text{sym } x = -x$. The explicit version below uses a McCarthy conditional, for $\text{lt0 } x = x < 0$. By exercise 5.51 term $id \cup \text{sym}$ in S can be ignored:

$$\begin{aligned} & \text{lt0} \rightarrow \text{sym}, id \subseteq \frac{\text{true}}{\text{geq0}} \\ \equiv & \quad \{ \text{shunting (5.46)} \} \\ & \text{geq0} \cdot (\text{lt0} \rightarrow \text{sym}, id) \subseteq \text{true} \\ \equiv & \quad \{ \text{fusion (2.71)} \} \\ & \text{lt0} \rightarrow \text{geq0} \cdot \text{sym}, \text{geq0} \subseteq \text{true} \\ \equiv & \quad \{ -x \geq 0 \Leftrightarrow x \leq 0 = \text{leq0 } x \} \\ & \text{lt0} \rightarrow \text{leq0}, \text{geq0} \subseteq \text{true} \end{aligned}$$

$$\begin{aligned}
&\equiv \{ x < 0 \Rightarrow x \leq 0 \text{ and } \neg (x < 0) \Leftrightarrow x \geq 0 \} \\
&\quad lt0 \rightarrow true, true \subseteq true \\
&\equiv \{ p \rightarrow f, f = f \text{ (exercise 5.51)} \} \\
&\quad true \\
&\square
\end{aligned}$$

Finally note that an equivalent way of stating (5.244) without using the domain operator is:

$$S \vdash R \equiv T \cdot S \cap T \cdot R \cap (R \cup S) = R \quad (5.246)$$

Exercise 5.59. Prove (5.246).

□

5.24 BACK TO FUNCTIONS

In this chapter we have argued that one needs *relations* in order to reason about *functions*. The inverse perspective — that relations can be represented as functions — also makes sense and it is, in many places, the approach that is followed.

Indeed, relations can be *transposed* back to functions without losing information. There are two transposes of interest. One is complete in the sense that it allows us to see *any* relation as a function. The other is specific, in the sense that it only applies to (the very important class of) simple relations (vulg. *partial* functions).

POWER TRANSPOSE Let $A \xrightarrow{R} B$ be a relation and define the function

$$\begin{aligned}
\Lambda R &: A \rightarrow \mathcal{P} B \\
\Lambda R a &= \{ b \mid b R a \}
\end{aligned}$$

which is such that:

$$\Lambda R = f \equiv \langle \forall b, a :: b R a \Leftrightarrow b \in f a \rangle \quad (5.247)$$

That is:

$$f = \Lambda R \Leftrightarrow \in \cdot f = R \quad (5.248)$$

cf.

$$\begin{array}{ccc}
A \rightarrow \mathcal{P} B & \xrightarrow{(\in \cdot)} & A \rightarrow B \\
& \cong & \\
& \xleftarrow{\Lambda} &
\end{array}$$

In words: any *relation* can be faithfully represented by set-valued *function*.

Moving the variables of (5.170) outwards by use of (5.17), we obtain the following *power transpose cancellation rule*:³⁶

$$\frac{\Lambda S}{\Lambda R} = \frac{S}{R} \quad (5.249)$$

Read from right to left, this shows a way of converting arbitrary symmetric divisions into function divisions.

“MAYBE” TRANSPOSE Let $A \xrightarrow{S} B$ be a *simple* relation. Define the function

$$\Gamma S : A \rightarrow B + 1$$

such that:

$$\Gamma S = f \Leftrightarrow \langle \forall b, a :: b S a \Leftrightarrow (i_1 b) = f a \rangle$$

That is:

$$f = \Gamma S \Leftrightarrow S = i_1^\circ \cdot f \quad (5.250)$$

cf.

$$\begin{array}{ccc} A \rightarrow B + 1 & \xrightarrow{(i_1^\circ \cdot)} & A \rightarrow B \\ & \cong & \\ & \xleftarrow{\Gamma} & \end{array}$$

In words: simple *relations* can always be represented by “Maybe”, or “pointer”-valued *functions*. Recall section 4.1, where the *Maybe* monad was used to “totalize” partial functions. Isomorphism (5.250) explains why such a totalization makes sense. For finite relations, and assuming these represented extensionally as lists of pairs, the function *lookup* :: *Eq* *a* \Rightarrow *a* \rightarrow [(*a*, *b*)] \rightarrow *Maybe* *b* in Haskell implements the “Maybe”-transpose.

5.25 BIBLIOGRAPHY NOTES

Chronologically, relational notation emerged — earlier than predicate logic itself — in the work of Augustus De Morgan (1806-71) on binary relations [35]. Later, Peirce (1839-1914) invented quantifier notation to explain De Morgan’s algebra of relations (see eg. [35] for details). De Morgan’s pioneering work was ill fated: the language³⁷ invented to explain his calculus of relations became eventually more popular than the calculus itself. Alfred Tarski (1901-83), who had a life-long struggle with quantified notation [14, 21], revived relation algebra. Together

³⁶ This rule is nothing but another way of stating exercise 4.48 proposed in [10]. Note that ΛR is always a function.

³⁷ Meanwhile named FOL, first order logic.

with Steve Givant he wrote a book (published posthumously) on *set theory without variables* [58].

Meanwhile, category theory was born, stressing the role of *arrows* and *diagrams* and on the arrow language of diagrams, which is inherently *pointfree*. The category of sets and functions immediately provided a basis for pointfree functional reasoning, but this was by and large ignored by John Backus (1924-2007) in his FP algebra of programs [7] which is APL-flavoured. (But there is far more in it than such a flavour, of course!) Anyway, Backus' landmark FP paper was among the first to show how relevant such reasoning style is to computing.

A bridge between the two pointfree schools, the relational and the categorial, was eventually established by Freyd and Ščedrov [17] in their proposal of the concept of an *allegory*. This gave birth to *typed* relation algebra and relation (semi-commutative) diagrams like those adopted in the current book for *relational thinking*. The pointfree algebra of programming (AoP) as it is understood today, stems directly from [17]. Its has reached higher education thanks to textbook [10] written by Bird and Moor.

In his book on *relational mathematics* [57], Gunther Schmidt makes extensive use of matrix displays, notation, concepts and operations in relation algebra. Winter [61] generalizes relation algebra to so-called Goguen categories.

In the early 1990s, the Groningen-Eindhoven MPC group led by Backhouse [1, 4] contributed decisively to the AoP by structuring relation algebra in terms of Galois connections. This elegant approach has been very influential in the way (typed) relation algebra was perceived afterwards, for instance in the way relation shrinking was introduced in the algebra [43, 53]. Galois connections are also the “Swiss knife” of [43].

Most of the current chapter was inspired by [4].

THEOREMS FOR FREE — BY CALCULATION

6.1 INTRODUCTION

As already stressed in previous chapters, type polymorphism remains one of the most useful and interesting ingredients of functional programming. For example, the two functions

$$\begin{aligned} \text{countBits} &: \mathbb{B}^* \rightarrow \mathbb{N}_0 \\ \text{countBits } [] &= 0 \\ \text{countBits } (b : bs) &= 1 + \text{countBits } bs \end{aligned}$$

and

$$\begin{aligned} \text{countNats} &: \mathbb{N}_0^* \rightarrow \mathbb{N}_0 \\ \text{countNats } [] &= 0 \\ \text{countNats } (b : bs) &= 1 + \text{countNats } bs \end{aligned}$$

are both subsumed by a single, *generic* (that is, parametric) program:

$$\begin{aligned} \text{count} &: (\forall A) A^* \rightarrow \mathbb{N}_0 \\ \text{count } [] &= 0 \\ \text{count } (a : as) &= 1 + \text{count } as \end{aligned}$$

Written as a catamorphism

$$(\text{in}_{\mathbb{N}_0} \cdot (id + \pi_2))$$

and thus even dispensing with a name, it becomes clear why this function is generic: nothing in

$$\text{in}_{\mathbb{N}_0} \cdot (id + \pi_2)$$

is susceptible to the *type* of the elements that are being counted up!

This form of polymorphism, known as *parametric polymorphism*, is attractive because

- one writes less code (*specific* solution = *generic* solution + *customization*);
- it is intellectually rewarding, as it brings elegance and economy in programming;

- and, last but not least¹,

“(...) from the type of a polymorphic function we can derive a theorem that it satisfies. (...) How useful are the theorems so generated? Only time and experience will tell (...)”

Recall that section 2.12 already addresses these theorems, also called *natural properties*. However, the full spread of naturality is not explored there. In particular, it does not address higher-order (exponential) types.

It turns out that the “free theorems” involving such types are easy to derive in relation algebra. The current chapter is devoted to such a generic derivation and includes a number of examples showing how vast the application of *free theorems* is.

6.2 POLYMORPHIC TYPE SIGNATURES

In any typed functional language, when declaring a polymorphic function one is bound to use the same generic format,

$$f : t$$

known as the function’s *signature*: f is the name of the function and t is a functional type written according to the following “grammar” of types:

$$t ::= t' \rightarrow t''$$

$$t ::= F(t_1, \dots, t_n) \quad F \text{ is a type constructor}$$

$$t ::= v \quad \text{a type variable, source of polymorphism.}$$

What does it mean for $f : t$ to be *parametrically* polymorphic? We shall see shortly that what matters in this respect is the formal structure of type t . Let

- V be the set of type variables involved in type expression t ;
- $\{R_v\}_{v \in V}$ be a V -indexed family of relations (f_v in case R_v is a function);
- R_t be a relation defined inductively as follows:

$$R_{t:=v} = R_v \tag{6.1}$$

$$R_{t:=F(t_1, \dots, t_n)} = F(R_{t_1}, \dots, R_{t_n}) \tag{6.2}$$

$$R_{t:=t' \rightarrow t''} = R_{t'} \rightarrow R_{t''} \tag{6.3}$$

Two questions arise: what does F in the right handside of (6.2) mean? What kind of relation is $R_{t'} \rightarrow R_{t''}$ in (6.3)?

First of all, and to answer the first question, we need the concept of *relator*, which extends that of a *functor* (introduced in section 3.8) to relations.

¹ Quoting *Theorems for free!*, by Philip Wadler [60].

6.3 RELATORS

A functor G is said to be a *relator* wherever, given a relation R from A to B , $G R$ extends R to G -structures: it is a relation from $G A$ to $G B$

$$\begin{array}{ccc} A & \xrightarrow{\quad} & G A \\ R \downarrow & & \downarrow G R \\ B & \xrightarrow{\quad} & G B \end{array} \quad (6.4)$$

which obeys the properties of a functor,

$$G id = id \quad (6.5)$$

$$G(R \cdot S) = (G R) \cdot (G S) \quad (6.6)$$

— recall (3.55) and (3.56) — plus the properties:

$$R \subseteq S \Rightarrow G R \subseteq G S \quad (6.7)$$

$$G(R^\circ) = (G R)^\circ \quad (6.8)$$

That is, a relator is a functor that is monotonic and commutes with converse. For instance, the “Maybe” functor $G X = 1 + X$ is an example of relator:

$$\begin{array}{ccc} A & \xrightarrow{\quad} & G A = 1 + A \\ R \downarrow & & \downarrow G R = id + R \\ B & \xrightarrow{\quad} & G B = 1 + B \end{array}$$

It is monotonic since $G R = id + R$ only involves monotonic operators and commutes with converse via (5.123). Let us unfold $G R = id + R$:

$$\begin{aligned} & y(id + R)x \\ \equiv & \quad \{ \text{unfolding the sum, cf. } id + R = [i_1 \cdot id, i_2 \cdot R] \text{ (5.119)} \} \\ & y(i_1 \cdot i_1^\circ \cup i_2 \cdot R \cdot i_2^\circ)x \\ \equiv & \quad \{ \text{relational union (5.57); image} \} \\ & y(\text{img } i_1)x \vee y(i_2 \cdot R \cdot i_2^\circ)x \\ \equiv & \quad \{ \text{let } NIL \text{ denote the sole inhabitant of the singleton type} \} \\ & y = x = i_1 NIL \vee \langle \exists b, a : y = i_2 b \wedge x = i_2 a : b R a \rangle \end{aligned}$$

In words: two “pointer-values” x and y are $G R$ -related iff they are both null or they are both defined and hold R -related data.

Finite lists also form a relator, $G X = X^*$. Given $B \xleftarrow{R} A$, relator $B^* \xleftarrow{R^*} A^*$ is the relation

$$\begin{aligned} s'(R^*)s & \Leftrightarrow \text{length } s' = \text{length } s \wedge \\ & \langle \forall i : 0 \leq i < \text{length } s : (s' !! i) R (s !! i) \rangle \end{aligned} \quad (6.9)$$

Exercise 6.1. Check properties (6.7) and (6.8) for the list relator defined above.

□

6.4 A RELATION ON FUNCTIONS

The next step needed to postulate free theorems requires a formal understanding of the arrow operator written on the right handside of (6.3).

This is achieved by defining the so-called “Reynolds arrow” relational operator, which establishes a relation on two functions f and g parametric on two other arbitrary relations R and S :

$$f(R \leftarrow S)g \equiv f \cdot S \subseteq R \cdot g \quad \begin{array}{ccc} A & \xleftarrow{S} & B \\ f \downarrow & \subseteq & \downarrow g \\ C & \xleftarrow{R} & D \end{array} \quad (6.10)$$

The typing rule is:

$$\frac{A \xleftarrow{S} B \quad C \xleftarrow{R} D}{C^A \xleftarrow{R \leftarrow S} D^B}$$

This is a powerful operator that satisfies many properties, for instance:

$$id \leftarrow id = id \quad (6.11)$$

$$(R \leftarrow S)^\circ = R^\circ \leftarrow S^\circ \quad (6.12)$$

$$R \leftarrow S \subseteq V \leftarrow U \Leftrightarrow R \subseteq V \wedge U \subseteq S \quad (6.13)$$

$$(R \leftarrow V) \cdot (S \leftarrow U) \subseteq (R \cdot S) \leftarrow (V \cdot U) \quad (6.14)$$

$$(f \leftarrow g^\circ)h = f \cdot h \cdot g \quad (6.15)$$

$$k(f \leftarrow g)h \equiv k \cdot g = f \cdot h \quad (6.16)$$

From property (6.13) we learn that the combinator is monotonic on the left hand side — and thus facts

$$S \leftarrow R \subseteq (S \cup V) \leftarrow R \quad (6.17)$$

$$\top \leftarrow S = \top \quad (6.18)$$

hold² — and anti-monotonic on the right hand side — and thus property

$$R \leftarrow \perp = \top \quad (6.19)$$

and the two distributive laws which follow:

$$S \leftarrow (R_1 \cup R_2) = (S \leftarrow R_1) \cap (S \leftarrow R_2) \quad (6.20)$$

$$(S_1 \cap S_2) \leftarrow R = (S_1 \leftarrow R) \cap (S_2 \leftarrow R) \quad (6.21)$$

It should be stressed that (6.14) expresses *fusion* only, not *fission*.

² Cf. $f \cdot S \cdot g^\circ \subseteq \top \Leftrightarrow \text{TRUE}$ concerning (6.18).

SUPREMA AND INFIMA Suppose relation R in (6.10) is a complete partial order \leq , that is, it has suprema and infima. What kind of relationship is established between two functions f and g such that

$$f ((\leq) \leftarrow S) g$$

holds? We reason:

$$\begin{aligned} & f ((\leq) \leftarrow S) g \\ \equiv & \{ (6.10) \} \\ & f \cdot S \subseteq (\leq) \cdot g \\ \equiv & \{ \text{shunting (5.46)} \} \\ & S \subseteq f^\circ \cdot (\leq) \cdot g \\ \equiv & \{ \text{go pointwise — (5.17), etc} \} \\ & \langle \forall a, b : a S b : f a \leq g b \rangle \\ \equiv & \{ \text{introduce supremum, for all } b \text{ (see below)} \} \\ & g b = \langle \bigvee a : a S b : f a \rangle \end{aligned}$$

The last step can be checked by unfolding the tautology $g b \leq g b$:

$$\begin{aligned} & \langle \forall b :: g b \leq g b \rangle \\ \equiv & \{ \text{unfold } g b \} \\ & \langle \forall b :: \langle \bigvee a : a S b : f a \rangle \leq g b \rangle \\ \equiv & \{ \text{universal law of suprema} \} \\ & \langle \forall b :: \langle \forall a : a S b : f a \leq g b \rangle \rangle \\ \equiv & \{ \text{quantifier calculus} \} \\ & \langle \forall a, b : a S b : f a \leq g b \rangle \\ & \square \end{aligned}$$

In summary:³

$$f ((\leq) \leftarrow S) g \quad \equiv \quad g b = \langle \bigvee a : a S b : f a \rangle \quad (6.22)$$

In words: $g b$ is the largest of all $(f a)$ such that $a S b$ holds.

Pattern $(\leq) \leftarrow \dots$ turns up quite often in relation algebra. Consider, for instance, a Galois connection $\alpha \vdash \gamma$ (5.132), that is,

$$\begin{aligned} & \alpha^\circ \cdot (\sqsubseteq) = (\leq) \cdot \gamma \\ \equiv & \{ \text{ping pong} \} \\ & \alpha^\circ \cdot (\sqsubseteq) \subseteq (\leq) \cdot \gamma \wedge \gamma^\circ \cdot (\geq) \subseteq (\supseteq) \cdot \alpha \end{aligned}$$

Following the same strategy as just above, we obtain pointwise definitions for the two adjoints of the connection:

$$\gamma x = \langle \bigvee y : \alpha y \sqsubseteq x : y \rangle \quad (6.23)$$

$$\alpha y = \langle \bigwedge x : y \leq \gamma x : x \rangle \quad (6.24)$$

³ Similarly, introducing infimum, for all a : $f a = \langle \bigwedge b : a S b : g b \rangle$.

6.5 FREE THEOREM OF TYPE t

We are now ready to establish the *free theorem* (FT) of type t , which is the following remarkably simple result:⁴

Given any function $\theta : t$, and V as above, then

$$\theta R_t \theta$$

holds, for any relational instantiation of type variables in V .

□

Note that this theorem

- is a result about t ;
- holds *independently* of the actual definition of θ .

So, it holds about any polymorphic function of type t .

6.6 EXAMPLES

Let us see the simplest of all examples, where the target function is the identity:

$$\theta = id : a \leftarrow a$$

We first calculate $R_{t=a \leftarrow a}$:

$$\begin{aligned} & R_{a \leftarrow a} \\ \equiv & \quad \{ \text{rule } R_{t=t' \leftarrow t''} = R_{t'} \leftarrow R_{t''} \} \\ & R_a \leftarrow R_a \end{aligned}$$

Then we derive the free theorem itself (R_a is abbreviated to R):

$$\begin{aligned} & id(R \leftarrow R)id \\ \equiv & \quad \{ (6.10) \} \\ & id \cdot R \subseteq R \cdot id \end{aligned}$$

In case R is a function f , the FT theorem boils down to *id's natural property*, $id \cdot f = f \cdot id$ — recall (2.10) — that can be read alternatively as stating that *id* is the *unit* of composition.

As a second example, consider $\theta = reverse : a^* \leftarrow a^*$, and first calculate $R_{t=a^* \leftarrow a^*}$:

$$\begin{aligned} & R_{a^* \leftarrow a^*} \\ \equiv & \quad \{ \text{rule } R_{t=t' \leftarrow t''} = R_{t'} \leftarrow R_{t''} \} \\ & R_{a^*} \leftarrow R_{a^*} \\ \equiv & \quad \{ \text{rule } R_{t=F(t_1, \dots, t_n)} = F(R_{t_1}, \dots, R_{t_n}) \} \\ & R_{a^*}^* \leftarrow R_{a^*}^* \end{aligned}$$

⁴ This result is due to J. Reynolds [56], advertised by P. Wadler [60] and re-written by Backhouse [2] in the pointfree style adopted in this book.

where $s R^* s'$ is given by (6.9). Next we calculate the FT itself (R_a abbreviated to R):

$$\begin{aligned} & \text{reverse}(R^* \leftarrow R^*) \text{reverse} \\ \equiv & \quad \{ \text{definition } f(R \leftarrow S)g \equiv f \cdot S \subseteq R \cdot g \} \\ & \text{reverse} \cdot R^* \subseteq R^* \cdot \text{reverse} \end{aligned}$$

In case R is a function r , this FT theorem boils down to *reverse's natural property*,

$$\text{reverse} \cdot r^* = r^* \cdot \text{reverse}$$

that is, $\text{reverse} [r \ a \mid a \leftarrow l] = [r \ b \mid b \leftarrow \text{reverse } l]$. For the general case, we obtain:

$$\begin{aligned} & \text{reverse} \cdot R^* \subseteq R^* \cdot \text{reverse} \\ \equiv & \quad \{ \text{shunting rule (5.46)} \} \\ & R^* \subseteq \text{reverse}^\circ \cdot R^* \cdot \text{reverse} \\ \equiv & \quad \{ \text{going pointwise (5.19, 5.17)} \} \\ & \langle \forall s, r :: s R^* r \Rightarrow (\text{reverse } s) R^* (\text{reverse } r) \rangle \end{aligned}$$

An instance of this pointwise version of *reverse-FT* will state that, for example, *reverse* will respect element-wise orderings ($R := <$):⁵

$$\begin{aligned} & \text{length } s = \text{length } r \wedge \langle \forall i : i \in \text{inds } s : (s !! i) < (r !! i) \rangle \\ & \quad \Downarrow \\ & \text{length}(\text{reverse } s) = \text{length}(\text{reverse } r) \\ & \quad \wedge \\ & \langle \forall j : j \in \text{inds } s : (\text{reverse } s !! j) < (\text{reverse } r !! j) \rangle \end{aligned}$$

(Guess other instances.)

As a third example, also involving finite lists, let us calculate the FT of

$$\text{sort} : a^* \leftarrow a^* \leftarrow (\text{Bool} \leftarrow (a \times a))$$

where the first parameter stands for the chosen ordering relation, expressed by a binary predicate:

$$\begin{aligned} & \text{sort}(R_{(a^* \leftarrow a^*) \leftarrow (\text{Bool} \leftarrow (a \times a))}) \text{sort} \\ \equiv & \quad \{ (6.2, 6.1, 6.3); \text{abbreviate } R_a := R \} \\ & \text{sort}((R^* \leftarrow R^*) \leftarrow (R_{\text{Bool}} \leftarrow (R \times R))) \text{sort} \\ \equiv & \quad \{ R_{t:=\text{Bool}} = \text{id} \text{ (constant relator)} \text{ — cf. exercise 6.11} \} \\ & \text{sort}((R^* \leftarrow R^*) \leftarrow (\text{id} \leftarrow (R \times R))) \text{sort} \\ \equiv & \quad \{ (6.10) \} \end{aligned}$$

⁵ Let $\text{inds } s$ denote the set $\{0, \dots, \text{length } s - 1\}$.

$$\begin{aligned}
& \text{sort} \cdot (\text{id} \leftarrow (R \times R)) \subseteq (R^* \leftarrow R^*) \cdot \text{sort} \\
\equiv & \quad \{ \text{shunting (5.46)} \} \\
& (\text{id} \leftarrow (R \times R)) \subseteq \text{sort}^\circ \cdot (R^* \leftarrow R^*) \cdot \text{sort} \\
\equiv & \quad \{ \text{introduce variables } f \text{ and } g \text{ (5.19, 5.17)} \} \\
& f(\text{id} \leftarrow (R \times R))g \Rightarrow (\text{sort } f)(R^* \leftarrow R^*)(\text{sort } g) \\
\equiv & \quad \{ (6.10) \text{ twice} \} \\
& f \cdot (R \times R) \subseteq g \Rightarrow (\text{sort } f) \cdot R^* \subseteq R^* \cdot (\text{sort } g)
\end{aligned}$$

Case $R := r$:

$$\begin{aligned}
& f \cdot (r \times r) = g \Rightarrow (\text{sort } f) \cdot r^* = r^* \cdot (\text{sort } g) \\
\equiv & \quad \{ \text{introduce variables} \} \\
& \left\langle \begin{array}{c} \forall a, b :: \\ f(r \ a, r \ b) = g(a, b) \end{array} \right\rangle \Rightarrow \left\langle \begin{array}{c} \forall l :: \\ (\text{sort } f)(r^* \ l) = r^*(\text{sort } g \ l) \end{array} \right\rangle
\end{aligned}$$

Denoting predicates f, g by infix orderings \leq, \preceq :

$$\left\langle \begin{array}{c} \forall a, b :: \\ r \ a \leq r \ b \equiv a \preceq b \end{array} \right\rangle \Rightarrow \left\langle \begin{array}{c} \forall l :: \\ \text{sort } (\leq)(r^* \ l) = r^*(\text{sort } (\preceq) \ l) \end{array} \right\rangle$$

That is, for r monotonic and injective,

$$\text{sort } (\leq) [r \ a \mid a \leftarrow l]$$

is always the same list as

$$[r \ a \mid a \leftarrow \text{sort } (\preceq) \ l]$$

Exercise 6.2. Let C be a nonempty data domain and let $c \in C$. Let \underline{c} be the “everywhere c ” function $\underline{c} : A \rightarrow C$ (2.12). Show that the free theorem of \underline{c} reduces to

$$\langle \forall R :: R \subseteq \top \rangle \tag{6.25}$$

□

Exercise 6.3. Calculate the free theorem associated with the projections

$$A \xleftarrow{\pi_1} A \times B \xrightarrow{\pi_2} B$$

and instantiate it to (a) functions; (b) coreflexives. Introduce variables and derive the corresponding pointwise expressions.

□

Exercise 6.4. As follow-up to exercise 6.2, consider higher order function $(\underline{_}) : a \rightarrow b \rightarrow a$ such that, given any x of type a , produces the constant function \underline{x} . Show that the equalities

$$\underline{f} \ x = f \cdot \underline{x} \quad (6.26)$$

$$\underline{x} \cdot f = \underline{x} \quad (6.27)$$

$$\underline{x}^\circ \cdot \underline{x} = \top \quad (6.28)$$

arise as corollaries of the free theorem of $(\underline{_})$.⁶

□

Exercise 6.5. The following is a well-known Haskell function

$$\text{filter} :: \forall a \cdot (a \rightarrow \mathbb{B}) \rightarrow [a] \rightarrow [a]$$

Calculate the free theorem associated with its type

$$\text{filter} : a^* \leftarrow a^* \leftarrow (\mathbb{B} \leftarrow a)$$

and instantiate it to the case where all relations are functions.

□

Exercise 6.6. In many sorting problems, data are sorted according to a given ranking function which computes each datum's numeric rank (eg. students marks, credits, etc). In this context one may parameterize sorting with an extra parameter f ranking data into a fixed numeric datatype, eg. the integers: $\text{serial} : (a \rightarrow \mathbb{N}_0) \rightarrow a^* \rightarrow a^*$. Calculate the FT of serial .

□

Exercise 6.7. Consider the following function from Haskell's Prelude:

$$\begin{aligned} \text{findIndices} &:: (a \rightarrow \mathbb{B}) \rightarrow [a] \rightarrow [\mathbb{Z}] \\ \text{findIndices } p \ xs &= [i \mid (x, i) \leftarrow \text{zip } xs \ [0..], p \ x] \end{aligned}$$

which yields the indices of elements in a sequence xs which satisfy p .

For instance, $\text{findIndices } (<0) \ [1, -2, 3, 0, -5] = [1, 4]$. Calculate the FT of this function.

□

⁶ Note that (6.27) is property (2.14) assumed in chapter 2.

Exercise 6.8. Wherever two equally typed functions f, g are such that $f a \leq g a$, for all a , we say that f is pointwise at most g and write $f \dot{\leq} g$,

$$f \dot{\leq} g = f \subseteq (\leq) \cdot g \quad \text{cf. diagram}$$

$$\begin{array}{ccc} & A & \\ f \swarrow & \subseteq & \searrow g \\ B & \longleftarrow & B \\ & \leq & \end{array}$$

recall (5.93). Show that implication

$$f \dot{\leq} g \Rightarrow (\text{map } f) \dot{\leq}^* (\text{map } g) \quad (6.29)$$

follows from the FT of the function $\text{map} : (a \rightarrow b) \rightarrow a^* \rightarrow b^*$.

□

Exercise 6.9. Infer the FT of the following function, written in Haskell syntax,

while :: $(a \rightarrow \mathbb{B}) \rightarrow (a \rightarrow a) \rightarrow (a \rightarrow b) \rightarrow a \rightarrow b$
while $p f g x = \text{if } \neg (p x) \text{ then } g x \text{ else while } p f g (f x)$

which implements a generic *while*-loop. Derive its corollary for functions.

□

6.7 CATAMORPHISM LAWS AS FREE THEOREMS

Recall from section 3.13 the concept of a catamorphism over a parametric type $\mathsf{T} a$:

$$\begin{array}{ccc} \mathsf{T} a & \xleftarrow{\text{in}_{\mathsf{T} a}} & \mathsf{B}(a, \mathsf{T} a) \\ \downarrow \llbracket g \rrbracket & & \downarrow \mathsf{B}(\text{id}, \llbracket g \rrbracket) \\ b & \xleftarrow{g} & \mathsf{B}(a, b) \end{array}$$

So $\llbracket - \rrbracket$ has generic type

$$\llbracket - \rrbracket : b \leftarrow \mathsf{T} a \leftarrow (b \leftarrow \mathsf{B}(a, b))$$

where $\mathsf{T} a \cong \mathsf{B}(a, \mathsf{T} a)$. Then the free theorem of $\llbracket - \rrbracket$ is

$$\llbracket - \rrbracket \cdot (R_b \leftarrow \mathsf{B}(R_a, R_b)) \subseteq (R_b \leftarrow \mathsf{F} R_a) \cdot \llbracket - \rrbracket$$

This unfolds into $(R_a, R_b$ abbreviated to R, S):

$$\llbracket - \rrbracket \cdot (S \leftarrow \mathsf{B}(R, S)) \subseteq (S \leftarrow \mathsf{T} R) \cdot \llbracket - \rrbracket$$

$$\begin{aligned}
&\equiv \{ \text{shunting (5.46)} \} \\
&\quad (S \leftarrow B(R, S)) \subseteq \llbracket _ \rrbracket^\circ (S \leftarrow T R) \cdot \llbracket _ \rrbracket \\
&\equiv \{ \text{introduce variables } f \text{ and } g \text{ (5.19, 5.17)} \} \\
&\quad f(S \leftarrow B(R, S))g \Rightarrow \llbracket f \rrbracket (S \leftarrow T R) \llbracket g \rrbracket \\
&\equiv \{ \text{definition } f(R \leftarrow S)g \equiv f \cdot S \subseteq R \cdot g \} \\
&\quad f \cdot B(R, S) \subseteq S \cdot g \Rightarrow \llbracket f \rrbracket \cdot T R \subseteq S \cdot \llbracket g \rrbracket
\end{aligned}$$

From the calculated free theorem of the catamorphism combinator,

$$f \cdot B(R, S) \subseteq S \cdot g \Rightarrow \llbracket f \rrbracket \cdot T R \subseteq S \cdot \llbracket g \rrbracket$$

we can infer:

- $\llbracket _ \rrbracket$ -fusion $(R, S := id, s)$:

$$f \cdot B(id, s) = s \cdot g \Rightarrow \llbracket f \rrbracket = s \cdot \llbracket g \rrbracket$$

— recall (3.71), for $F f = B(id, f)$;

- $\llbracket _ \rrbracket$ -absorption $(R, S := r, id)$:

$$f \cdot B(r, id) = g \Rightarrow \llbracket f \rrbracket \cdot T r = \llbracket g \rrbracket$$

whereby, substituting $g := f \cdot B(r, id)$:

$$\llbracket f \rrbracket \cdot T r = \llbracket f \cdot B(r, id) \rrbracket$$

— recall (3.77).

Exercise 6.10. Let

$$iprod = \llbracket [1, (\times)] \rrbracket$$

be the function that multiplies all natural numbers in a given list, and even be the predicate which tests natural numbers for evenness. Finally, let

$$exists = \llbracket [\underline{\text{FALSE}}, (\vee)] \rrbracket$$

be the function that implements existential quantification over a list of Booleans. From (6.30) infer

$$even \cdot iprod = exists \cdot even^*$$

meaning that the product $n_1 \times n_2 \times \dots \times n_m$ is even if and only if some n_i is so.

□

Exercise 6.11. Show that the identity relator Id , which is such that $Id R = R$ and the constant relator K (for a given data type K) which is such that $K R = id_K$

are indeed relators.

□

Exercise 6.12. Show that product

$$\begin{array}{ccc}
 A & C & \cdots \cdots \cdots G(A, C) = A \times C \\
 \downarrow R & \downarrow S & \downarrow G(R, S) = R \times S \\
 B & D & \cdots \cdots \cdots G(B, D) = B \times D
 \end{array}$$

is a (binary) relator.

□

6.8 BIBLIOGRAPHY NOTES

The free theorem of a polymorphic function is a result due to computer scientist John Reynolds [56]. It became popular under the “theorems for free” heading coined by Phil Wadler [60]. The original pointwise setting of this result was re-written in the pointfree style in [2] thanks to the *relation on functions* combinator (6.10) first introduced by Roland Backhouse in [3].

More recently, Janis Voigtlaender devoted a whole research project to free theorems, showing their usefulness in several areas of computer science [39]. One outcome of this project was an automatic generator of free theorems for types written in Haskell syntax. This is (was?) available from Janis Voigtlaender’s home page:

<http://www-ps.iai.uni-bonn.de/ft>

The relators used in the calculational style followed in this book are implemented in this automatic generator by so-called structural functor *lifting*.

 CONTRACT-ORIENTED PROGRAMMING

The chapters of the first part of this book rely on a type-polymorphic notion of computation, captured by the omnipresent use of the arrow notation

$$B \xleftarrow{f} A$$

where A and B are *types*.

The generalization from functions to relations carried out in the previous two chapters has preserved the same principle — all relational combinators are typed in the same way. There is thus an implicit assumption of *static type checking* in the overall approach — types are checked at “compile time”. Expressions which don’t type are automatically excluded.

However, examples such as the Alcuin puzzle show that this is insufficient. Why? Because the types involved are most often “too large”: the whole purpose of the puzzle is to consider only the inhabitants of type $Bank^{being}$ — functions that describe all possible configurations in the puzzle — that satisfy the “starvation property”, recall (5.76). Moreover, the *carry* – operation (5.208) *should* preserve this property — something we didn’t at all check in the previous chapter!

Let us generalize the situation in this puzzle to that of a function $f : A \rightarrow A$ and a predicate $p : A \rightarrow \mathbb{B}$ that should be preserved by f . Predicates such as p have become known as *invariants* by software theorists. The preservation requirement is captured by:

$$\langle \forall a : p\ a : p\ (f\ a) \rangle$$

Note how the type A is now divided in two parts — a “good one”, $\{a \mid a \in A \wedge p\ a\}$ and a “bad one”, $\{a \mid a \in A \wedge \neg (p\ a)\}$. By identifying p as an invariant, the programmer is *obliged* to ensure a “good” output $f\ a$ wherever a “good” input is passed to f . For “bad” inputs nothing is requested.

The situation above can be generalized to some $f : A \rightarrow B$ where B is subject to some invariant $q : B \rightarrow \mathbb{B}$. So f is *obliged* to ensure “good” outputs satisfying q . It may well be the case that the only way for f to ensure “good” outputs is to restrict its inputs by some precondition $p : A \rightarrow \mathbb{B}$. Thus the proof obligation above generalizes to:

$$\langle \forall a : p\ a : q\ (f\ a) \rangle \tag{7.1}$$

One might tentatively try and express this requirement by writing

$$p \xrightarrow{f} q$$

where predicates p and q take the place of the original types A and B , respectively. This is what we shall do, calling assertion $p \xrightarrow{f} q$ a *contract*. Note how we are back to the function-as-a-contract view of section 2.1 but in a wider setting:

f commits itself to producing a “good” B -value (wrt. q) provided it is supplied with a “suitable” A -value (wrt. p).

The main difference compared to section 2.1 is that the well-typing of $p \xrightarrow{f} q$ cannot be mechanically ascertained at “compile time” — it has to be validated by a formal proof — the proof obligation (7.1) mentioned above. This kind of type checking is often referred to as “extended type checking”.

In real life software design data type invariants can be arbitrarily complex — think of all legal restrictions imposed on the organized societies of today! The increasing “softwarization” of our times forces us to think that, as in the regular functioning of such organized *societies*, programs should interact with each other via *formal contracts* establishing what they rely upon or guarantee among themselves. This is the only way to ensure *safety* and *security* essential to reliable, mechanized operations.

This chapter will use relation algebra to describe such contracts and develop a simple theory about them, enabling compositionality as before. Relations (including functions) will play a double role — they will not only describe computations but also the data structures involved in such computations, in a unified and elegant way.

7.1 CONTRACTS

It should be routine work for the reader to check that

$$f \cdot \Phi_p \subseteq \Phi_q \cdot f \tag{7.2}$$

means exactly the same as (7.1) above. In software design terminology, this is known as a (functional) *contract*, and we shall write

$$p \xrightarrow{f} q \tag{7.3}$$

to denote it — a notation that generalizes the type $A \rightarrow B$ of f , as already observed. Thanks to (5.207), (7.2) can also be written:

$$f \cdot \Phi_p \subseteq \Phi_q \cdot \top \tag{7.4}$$

Predicates p and q in contract $p \xrightarrow{f} q$ shall be referred to as the contract's *precondition* and *postcondition*, respectively. Contracts compose sequentially, see the following exercise.

Exercise 7.1. Show that $q \xleftarrow{g \cdot f} p$ holds provided $r \xleftarrow{f} p$ and $q \xleftarrow{g} r$ hold.

□

WEAKEST PRE-CONDITIONS Note that more than one (*pre*) condition p may ensure (*post*) condition q on the outputs of f . Indeed, contract $false \xrightarrow{f} q$ always holds, but it is useless — pre-condition *false* is “*unacceptably strong*”.

Clearly, the weaker p the better. The question is, then: is there a *weakest* such p ? We calculate:

$$\begin{aligned}
 & f \cdot \Phi_p \subseteq \Phi_q \cdot f \\
 \equiv & \quad \{ \text{recall (5.207)} \} \\
 & f \cdot \Phi_p \subseteq \Phi_q \cdot \top \\
 \equiv & \quad \{ \text{shunting (5.46); (5.205)} \} \\
 & \Phi_p \subseteq f^\circ \cdot \frac{true}{q} \\
 \equiv & \quad \{ (5.52) \} \\
 & \Phi_p \subseteq \frac{true}{q \cdot f} \\
 \equiv & \quad \{ \Phi_p \subseteq id ; (5.58) \} \\
 & \Phi_p \subseteq id \cap \frac{true}{q \cdot f} \\
 \equiv & \quad \{ (5.198) \} \\
 & \Phi_p \subseteq \Phi_{(q \cdot f)}
 \end{aligned}$$

We conclude that $q \cdot f$ is such a *weakest* pre-condition. Notation $wp(f, q) = q \cdot f$ is often used to denote a *weakest* pre-condition (WP). This is the weakest constraint on inputs for outputs by f to fall within q . The special situation of a weakest precondition is nicely captured by the universal property:

$$f \cdot \Phi_p = \Phi_q \cdot f \quad \equiv \quad p = q \cdot f \quad (7.5)$$

where $p = wp(f, q)$ could be written instead of $p = q \cdot f$, as seen above. Property (7.5) enables a “logic-free” calculation of weakest pre-conditions, as we shall soon see: given f and post-condition q , there

always exists a unique (weakest) precondition p such that $\Phi_q \cdot f$ can be replaced by $f \cdot \Phi_p$. Moreover:

$$\frac{f}{f} \cdot \Phi_p = \Phi_p \cdot \frac{f}{f} \iff p \leq f \quad (7.6)$$

where \leq denotes the injectivity preorder (5.230) on functions.¹

Exercise 7.2. Calculate the weakest pre-condition $wp(f, q)$ for the following function / post-condition pairs:

- $f \ x = x^2 + 1, q \ y = y \leq 10$ (in \mathbb{R})
- $f = \mathbb{N}_0 \xrightarrow{\text{succ}} \mathbb{N}_0, q = \text{even}$
- $f \ x = x^2 + 1, q \ y = y \leq 0$ (in \mathbb{R})

□

INVARIANTS In case $p = q$ in a contract (7.3), that is, in case of $q \xrightarrow{f} q$ holding, we say that q is an *invariant* of f , meaning that the “truth value” of q remains unchanged by execution of f . More generally, invariant q is *preserved* by function f provided contract $p \xrightarrow{f} q$ holds and $p \Rightarrow q$, that is, $\Phi_p \subseteq \Phi_q$.

Some pre-conditions are weaker than others wrt. invariant preservation. We shall say that w is the *weakest* pre-condition for f to preserve invariant q wherever $wp(f, q) = w \wedge q$, where $\Phi_{p \wedge q} = \Phi_p \cdot \Phi_q$.

Recalling the Alcuin puzzle, let us define the *starvation* invariant as a predicate on the state of the puzzle, passing the *where* function as a parameter w :

$$\text{starving } w = w \cdot \text{CanEat} \subseteq w \cdot \underline{\text{Farmer}}$$

Then the *contract*

$$\text{starving} \xrightarrow{\text{carry } b} \text{starving}$$

would mean that the function *carry b* — that should transfer the beings in b to the other bank of the river — always preserves the invariant:

$$wp(\text{carry } b, \text{starving}) = \text{starving}.$$

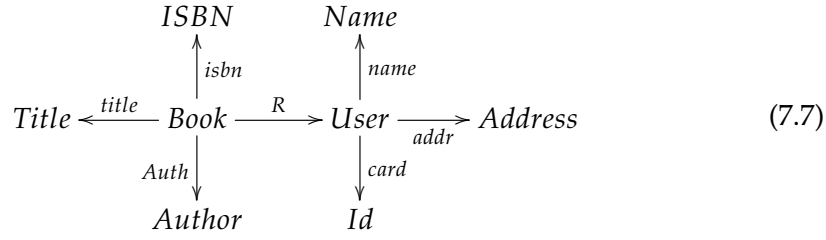
Things are not that easy, however: there is a need for a *pre-condition* ensuring that b includes the farmer together with a good choice of the being to carry!

Let us see some simpler examples first.

¹ The interested reader will find the proofs of (7.5) and (7.6) in reference [52].

7.2 LIBRARY LOAN EXAMPLE

Consider the following relational data model of a library involving books and users that can borrow its books:



All arrows denote attributes (functions) but two — *Auth* and *R*. The former is a relation because a book can have more than one author.² The latter is the most interesting relation of the model, $u R b$ meaning “book b currently on loan to library user u ”. Quite a few invariants are required in this model, for instance:

- the same book cannot be not on loan to more than one user;
- no book exists with no authors;
- no two different users have the same card *Id*;
- books with the same *ISBN* should have the same title and the same authors.

Such properties (invariants) are easy to encode:

- no book on loan to more than one user:

$$Book \xrightarrow{R} User \text{ is simple}$$

- no book without an author:

$$Book \xrightarrow{Auth} Author \text{ is entire}$$

- no two users with the same card *Id*:

$$User \xrightarrow{card} Id \text{ is injective}$$

- *ISBN* is a key attribute:

$$ISBN \xrightarrow{title \cdot isbn^\circ} Title \text{ and } ISBN \xrightarrow{\Lambda Auth \cdot isbn^\circ} P Author \text{ are simple relations.}$$

Since all other arrows are functions, they are simple and entire.

Let us now spell out such invariants in terms of relational assertions (note the role of the injectivity preorder):

- no book on loan to more than one user:

² Its power transpose (5.247) — $\Lambda Auth : Book \rightarrow P Author$ — gives the set of authors of a book.

$$id \leq R^\circ$$

equivalent to $\text{img } R \subseteq id$;

- no book without an author:

$$id \subseteq \ker Auth$$

- no two users with the same card Id:

$$id \leq card$$

equivalent to $\ker card \subseteq id$.

- ISBN is a *key* attribute:

$$title \leq isbn \wedge \Lambda Auth \leq isbn$$

equivalent to $\frac{isbn}{isbn} \subseteq \frac{title}{title}$ and $\frac{isbn}{isbn} \subseteq \frac{Auth}{Auth}$, respectively.³

Below we focus on the first invariant, *no book on loan to more than one user*. To bring life to our model, let us think of two operations on $User \xleftarrow{R} Book$, one that *returns* books to the library and another that *records* new borrowings:

$$(\text{return } S) R = R - S \quad (7.8)$$

$$(\text{borrow } S) R = S \cup R \quad (7.9)$$

Note that parameter S is of type $User \xleftarrow{R} Book$, indicating which users borrow/return which books. Clearly, these operations only change the *books-on-loan* relation R , which is conditioned by invariant

$$\text{inv } R = \text{img } R \subseteq id \quad (7.10)$$

The question is, then: are the following “types”

$$\text{inv} \xleftarrow{\text{return } S} \text{inv} \quad (7.11)$$

$$\text{inv} \xleftarrow{\text{borrow } S} \text{inv} \quad (7.12)$$

valid? Let us check (7.11):

$$\begin{aligned} & \text{inv} (\text{return } S R) \\ \equiv & \quad \{ \text{inline definitions} \} \\ & \text{img } (R - S) \subseteq id \\ \Leftarrow & \quad \{ \text{since img is monotonic} \} \\ & \text{img } R \subseteq id \\ \equiv & \quad \{ \text{definition} \} \\ & \text{inv } R \\ & \square \end{aligned}$$

³ Note the use of (5.170) in the second case.

So, for all R , $inv\ R \Rightarrow inv\ (\text{return } S\ R)$ holds — invariant inv is preserved.

At this point note that (7.11) was checked only as a *warming-up* exercise — we don't actually need to worry about it! Why?

As $R - S$ is smaller than R (exercise 5.41) and “*smaller than injective is injective*” (5.82), it is immediate that inv (7.10) is preserved.

To see this better, we unfold and draw definition (7.10) in the form of a diagram:

$$inv\ R = \begin{array}{ccc} Book & \xleftarrow{R^\circ} & User \\ R \downarrow & \subseteq & \downarrow id \\ User & \xleftarrow{id} & User \end{array}$$

As R occurs only in the lower-path of the diagram, it can always get smaller.

This “rule of thumb” does not work for *borrow* S because, in general, $R \subseteq \text{borrow } S\ R$. This time R gets bigger, not smaller, and we do have to check the contract:

$$\begin{aligned} & inv\ (\text{borrow } S\ R) \\ \equiv & \quad \{ \text{inline definitions} \} \\ & \text{img } (S \cup R) \subseteq id \\ \equiv & \quad \{ \text{exercise 5.15} \} \\ & \text{img } R \subseteq id \wedge \text{img } S \subseteq id \wedge S \cdot R^\circ \subseteq id \\ \equiv & \quad \{ \text{definition of } inv \} \\ & inv\ R \wedge \underbrace{\text{img } S \subseteq id \wedge S \cdot R^\circ \subseteq id}_{wp\ (\text{borrow } S, inv)} \end{aligned}$$

Thus the complete definition of the *borrow* operation becomes, in the notation of section 5.3:

$$\begin{aligned} & \text{Borrow } (S, R : Book \rightarrow User)\ R' : Book \rightarrow User \\ & \text{pre } S \cdot S^\circ \subseteq id \wedge S \cdot R^\circ \subseteq id \\ & \text{post } R' = R \cup S \end{aligned}$$

Why have we written *Borrow* instead of *borrow* as before? This is because *borrow* has become a *simple* relation

$$\text{Borrow} = \text{borrow} \cdot \Phi_{\text{pre}}$$

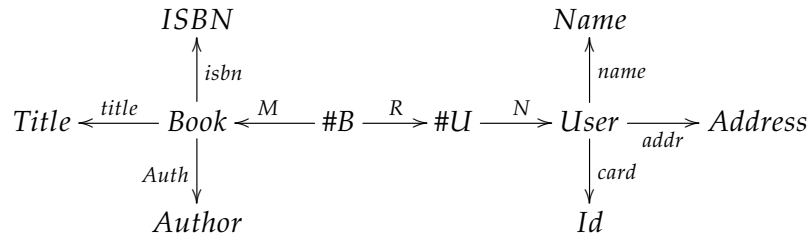
It is no longer a function since its (weakest) precondition is not the predicate *true*. (Recall that lowercase identifiers are reserved to functions only.) This precondition was to be expected, as spelt out by rendering $S \cdot R^\circ \subseteq id$ in pointwise notation: for all users u, u' ,

$$\langle \exists b : u S b : u' R b \rangle \Rightarrow u = u'$$

should hold. So, after the operation takes place, the result state $R' = R \cup S$ won't have the same book on loan twice to different users. (Of course, the same must happen about S itself, which is the same predicate for $R = S$.) Interestingly, the weakest precondition is not ruling out the situation in which $u S b$ and $u R b$ hold, for some book b and user u . Not only this does not harm the model but also it corresponds to a kind of renewal of a previous borrowing.

EVOLUTION The library loan model (7.7) given above is not realistic in the following sense — it only “gives life” to the borrowing relation R . In a sense, it assumes that all books have been bought and all users are registered.

How do we improve the model so that new books can be acquired and new users can join the library? Does this evolution entail a complete revision of (7.7)? Not at all. What we have to do is to *add* two new relations, say M and N , the first recording the books currently available in the library and the second the users currently registered for loaning:



Two new datatypes have been added: $\#U$ (unique identifier of each user) and $\#B$ (key identifying each book). Relations M and N have to be simple. The operations defined thus far stay the same, provided $\#B$ replaces $Book$ and $\#U$ replaces $User$ — advantages of a polymorphic notation. New operations can be added for

- acquiring new books — will change relation M only;
- registering new users — will change relation N only;
- cancelling users' registrations — will change relation N only.

There is, however, something that has not been considered: think of a starting state where $M = \perp$ and $N = \perp$, that is, the library has no users, no books yet. Then necessarily $R = \perp$. In general, users cannot borrow books that don't exist,

$$\delta R \subseteq \delta M$$

and not-registered users cannot borrow books at all:

$$\rho R \subseteq \delta N$$

Invariants of this kind capture so-called *referential integrity* constraints. They can be written with less symbols, cf.

$$R \subseteq T \cdot M$$

and

$$R \subseteq N^\circ \cdot T$$

respectively. Using the “thumb” rules as above, it is clear that, with respect to *referential integrity*:

- returning books is no problem, because R is only on the lower side of both inclusions;
- *borrowing* books calls for new contracts — R is on the lower side and it increases!
- registering new users and buying new books are no problem, because M and N are on the upper side only;
- unregistering users calls for a contract because N is on the upper side and decreases — users must return all books before unregistering!

7.3 MOBILE PHONE EXAMPLE

In this example we go back to the *store* operation on a mobile phone list of calls specified by (5.2). Of the three invariants we select (b), the one requiring no duplicate calls in the list. Recall, in Haskell, the function $(!!) :: [a] \rightarrow \mathbb{Z} \rightarrow a$. This tells how a finite list s is converted into a partial function $(s!!)$ of type $\mathbb{Z} \rightarrow a$. In fact, the partiality extends to the negative numbers⁴ and so we should regard $(s!!)$ as a *simple* relation⁵ even if restricted to the type $a \leftarrow \mathbb{N}_0$, as we shall do below.

The no-duplicates requirement requests $(s!!)$ to be injective: in case $s!!i$ and $s!!j$ are defined, $i \neq j \Rightarrow s!!i \neq s!!j$. Let $L = (s!!)$. Then we can re-specify the operations of *store* in terms of L , as follows:⁶

$$\begin{aligned} \text{inv } L &= \text{id} \leq L \\ \text{filter } (c \neq) L &= L - \underline{c} \\ c : L &= [\underline{c}, L] \cdot \text{in}^\circ \end{aligned}$$

where $\text{in} = [\underline{0}, \text{succ}]$ — the Peano algebra which builds up natural numbers.⁷ By (5.121) the definition of $c : L$ can also be written $\underline{c} \cdot \underline{0}^\circ \cup$

⁴ Try $[2, 3, 3]!!(-1)$, for instance.

⁵ Partial functions are *simple* relations, as we know.

⁶ Knowing that take 10 will always yield its input or a smaller list, and that *smaller than injective is injective* (5.82), we only need to focus on $(c:) \cdot \text{filter } (c \neq)$.

⁷ Recall section 3.1.

$L \cdot \text{succ}^\circ$, explicitly telling that c is placed in position 0 while L is shifted one position up to make room for the new element. We calculate:

$$\begin{aligned}
& \text{inv } (c : (\text{filter } (c \neq) L)) \\
\equiv & \quad \{ \text{inv } L = \text{id} \leq L, \text{ using the injectivity preorder } \} \\
& \text{id} \leq c : (\text{filter } (c \neq) L) \\
\equiv & \quad \{ \text{in-line definitions } \} \\
& \text{id} \leq [\underline{c}, L - \underline{c}] \cdot \text{in}^\circ \\
\equiv & \quad \{ \text{Galois connection (5.233)} \} \\
& \text{in} \leq [\underline{c}, L - \underline{c}] \\
\equiv & \quad \{ (5.235) ; \text{in is as injective as id} \} \\
& \text{id} \leq \underline{c} \wedge \text{id} \leq L - \underline{c} \wedge \underline{c}^\circ \cdot (L - \underline{c}) \subseteq \perp \\
\Leftarrow & \quad \{ \text{constant functions are injective; } L \subseteq \top \} \\
& \text{id} \leq L - \underline{c} \wedge \underline{c}^\circ \cdot (\top - \underline{c}) \subseteq \perp \\
\Leftarrow & \quad \{ \text{smaller than injective is injective ; } \underline{c}^\circ \cdot (\top - \underline{c}) = \perp \text{ (5.155)} \} \\
& \text{id} \leq L \\
& \square
\end{aligned}$$

Having given two examples of contract checking in two quite different domains, let us prepare for checking that of the Alcuin puzzle. By exercise 5.20 we already know that any of the starting states $w = \underline{Left}$ or $w = \underline{Right}$ satisfy the invariant:

$$\text{starving } w = w \cdot \text{CanEat} \subseteq w \cdot \underline{Farmer}.$$

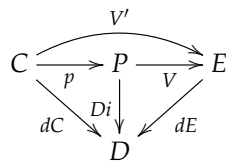
The only operation defined is

$$\text{carry who where} = (\in \text{who}) \rightarrow \text{cross} \cdot \text{where}, \text{ where}$$

Clearly, calculating the weakest precondition for this operation to preserve *starving* is expected to be far more complex than in the previous examples, since *where* is everywhere in the invariant. Can this be made simpler?

The answer is positive provided we understand a technique to be adopted, called *abstract interpretation*. So we postpone the topic of this paragraph to section 7.5, where abstract interpretation will be introduced. In between, we shall study a number of rules that can be used to address contracts in a structured way.

Exercise 7.3. Consider the voting system described by the relations of the diagram below,



where electors can vote in political parties or nominally in members of such parties. In detail: (a) $p \ c$ denotes the party of candidate c ; (b) $dC \ c$ denotes the district of candidate c ; (c) $dE \ e$ denotes the district of elector e ; (d) $d \ Di \ p$ records that party p has a list of candidates in district d ; (e) $e \ V \ p$ indicates that elector e voted in party p ; (f) $e \ V' \ c$ indicates that elector e voted nominally in candidate c .

There are several invariants to take into account in this model, namely:

$$\text{inv1 } (V, V') = V : E \leftarrow P \text{ and } V' : E \leftarrow C \text{ are injective} \quad (7.13)$$

$$\text{inv2 } (V, V') = V^\circ \cdot V' = \perp \quad (7.14)$$

since an elector cannot vote in more than one candidate or party;

$$\text{inv3 } (V, V') = dE \cdot [V, V'] \subseteq [Di, dC] \quad (7.15)$$

since each elector is registered in one district and can only vote in candidates of that district.

When the elections take place, relations p , dC , dE and Di are static, since all lists and candidates are fixed before people can vote. Once it is over, the scrutinity of the votes is carried out function

$$\text{batch } (V, V', X) = \dots$$

where $X : E \rightarrow (P + C)$ is a batch of votes to be loaded into the system.

Complete the definition of batch and discharge the proof obligations of the contracts that this function must satisfy.

□

7.4 A CALCULUS OF FUNCTIONAL CONTRACTS

The number and complexity of invariants in real life problems invites us to develop *divide & conquer* rules alleviating the proof obligations that have to be discharged wherever contracts are needed. All such rules have definition (7.2) as starting point. Let us see, for instance, what happens wherever the input predicate in (7.3) is a disjunction:

$$\begin{aligned}
 & \Phi_q \xleftarrow{f} \Phi_{p_1} \cup \Phi_{p_2} \\
 \equiv & \quad \{ (7.2) \} \\
 & f \cdot (\Phi_{p_1} \cup \Phi_{p_2}) \subseteq \Phi_q \cdot f \\
 \equiv & \quad \{ \text{distribution of } (f \cdot) \text{ by } \cup \text{ (5.60)} \} \\
 & f \cdot \Phi_{p_1} \cup f \cdot \Phi_{p_2} \subseteq \Phi_q \cdot f \\
 \equiv & \quad \{ \cup\text{-universal (5.59)} \} \\
 & f \cdot \Phi_{p_1} \subseteq \Phi_q \cdot f \wedge f \cdot \Phi_{p_2} \subseteq \Phi_q \cdot f \\
 \equiv & \quad \{ (7.2) \text{ twice} \} \\
 & \Phi_q \xleftarrow{f} \Phi_{p_1} \wedge \Phi_q \xleftarrow{f} \Phi_{p_2}
 \end{aligned}$$

Recall that the disjunction $p \vee q$ of two predicates is such that $\Phi_{p \vee q} = \Phi_p \cup \Phi_q$ holds. So we can write the result above in the simpler notation (7.3) as the contract decomposition rule:

$$q \xleftarrow{f} p \vee r \quad \equiv \quad q \xleftarrow{f} p \wedge q \xleftarrow{f} r \quad (7.16)$$

The dual rule,

$$\Phi_q \cdot \Phi_r \xleftarrow{f} \Phi_p \quad \equiv \quad \Phi_q \xleftarrow{f} \Phi_p \wedge \Phi_r \xleftarrow{f} \Phi_p$$

is calculated in the same way and written

$$q \wedge r \xleftarrow{f} p \quad \equiv \quad q \xleftarrow{f} p \wedge r \xleftarrow{f} p \quad (7.17)$$

in the same notation, since $\Phi_{p \wedge q} = \Phi_p \cap \Phi_q$. The fact that contracts compose sequentially (exercise 7.1) enables the corresponding decomposition, once a suitable middle predicate r is found:

$$q \xleftarrow{g \cdot h} p \quad \Leftarrow \quad q \xleftarrow{g} r \wedge r \xleftarrow{h} p \quad (7.18)$$

This follows straight from (7.3, 7.2), as does the obvious rule concerning identity

$$q \xleftarrow{id} p \quad \equiv \quad q \Leftarrow p \quad (7.19)$$

since $p \Rightarrow q \Leftrightarrow \Phi_p \subseteq \Phi_q$. The expected

$$p \xleftarrow{id} p$$

immediately follows from (7.19).

Now suppose that we have contracts $q \xleftarrow{f} p$ and $r \xleftarrow{g} p$. What kind of contract can we infer for $\langle f, g \rangle$? We calculate:

$$\begin{aligned} & \Phi_q \xleftarrow{f} \Phi_p \quad \wedge \quad \Phi_r \xleftarrow{g} \Phi_p \\ \equiv & \quad \{ (7.3, 7.2) \text{ twice} \} \\ & f \cdot \Phi_p \subseteq \Phi_q \cdot f \quad \wedge \quad g \cdot \Phi_p \subseteq \Phi_r \cdot g \\ \equiv & \quad \{ \text{cancellations (2.22)} \} \\ & \pi_1 \cdot \langle f, g \rangle \cdot \Phi_p \subseteq \Phi_q \cdot f \quad \wedge \quad \pi_2 \cdot \langle f, g \rangle \cdot \Phi_p \subseteq \Phi_r \cdot g \\ \equiv & \quad \{ \text{universal property (5.103)} \} \\ & \langle f, g \rangle \cdot \Phi_p \subseteq \langle \Phi_q \cdot f, \Phi_r \cdot g \rangle \\ \equiv & \quad \{ \text{absorption (5.106)} \} \\ & \langle f, g \rangle \cdot \Phi_p \subseteq (\Phi_q \times \Phi_r) \cdot \langle f, g \rangle \\ \equiv & \quad \{ (7.3, 7.2) \} \\ & \Phi_q \times \Phi_r \xleftarrow{\langle f, g \rangle} \Phi_p \end{aligned}$$

Defining $p \boxtimes q$ such that $\Phi_{p \boxtimes q} = \Phi_p \times \Phi_q$ we obtain the contract decomposition rule:

$$q \boxtimes r \xleftarrow{\langle f, g \rangle} p \quad \equiv \quad q \xleftarrow{f} p \wedge r \xleftarrow{g} p \quad (7.20)$$

which justifies the existence of arrow $\langle f, g \rangle$ in the diagram

$$\begin{array}{ccccc} & & q \boxtimes r & & \\ \pi_1 \swarrow & & & \searrow \pi_2 & \\ q & & q \boxtimes r & & r \\ & \swarrow f & \uparrow \langle f, g \rangle & \searrow g & \\ & p & & & \end{array} \quad (7.21)$$

where predicates (coreflexives) are promoted to objects (nodes in diagrams).

Exercise 7.4. Check the contracts $q \xleftarrow{\pi_1} q \boxtimes r$ and $q \boxtimes r \xrightarrow{\pi_2} r$ of diagram (7.21).

□

Let us finally see how to handle conditional expressions of the form *if* ($c \ x$) *then* ($f \ x$) *else* ($g \ x$) which, by (5.213), transform into

$$c \rightarrow f, g = f \cdot \Phi_c \cup g \cdot \Phi_{\neg c} \quad (7.22)$$

In this case, (7.4) offers a better standpoint for calculation than (7.2), as the reader may check in calculating the following rule for conditionals:

$$\Phi_q \xleftarrow{c \rightarrow f, g} \Phi_p \quad \equiv \quad \begin{cases} \Phi_q \xleftarrow{f} \Phi_p \cdot \Phi_c \\ \Phi_q \xleftarrow{g} \Phi_p \cdot \Phi_{\neg c} \end{cases} \quad (7.23)$$

This is because it is hard to handle $c \rightarrow f, g$ on the upper side, \top being more convenient.

Further contract rules can be calculated on the same basis, either elaborating on the predicate structure or on the combinator structure. However, all the cases above involve functions only and the semantics of computations are, in general, relations. So our strategy is to generalize definition (7.2) from functions to arbitrary relations.

RELATIONAL CONTRACTS Note that $S = R \cdot \Phi_p$ means

$$b \ S \ a \Leftrightarrow p \ a \wedge b \ R \ a$$

— that is, S is R pre-conditioned by p . Dually, $\Phi_q \cdot R$ is the largest part of R which yields outputs satisfying q — R post-conditioned by q . By writing

$$R \cdot \Phi_p \subseteq \Phi_q \cdot R \quad (7.24)$$

— which is equivalent to

$$R \cdot \Phi_p \subseteq \Phi_q \cdot \top \quad (7.25)$$

by (5.207) and even equivalent to

$$\Phi_p \subseteq R \setminus (\Phi_q \cdot \top) \quad (7.26)$$

by (5.159) — we express a very important fact about R regarded as a (possibly non-deterministic, undefined) program R : condition p on the inputs is *sufficient* for condition q to hold on the outputs:

$$\langle \forall a : p a : \langle \forall b : b R a : q b \rangle \rangle$$

Thus we generalize functional contracts (7.2) to arbitrary relations,

$$p \xrightarrow{R} q \equiv R \cdot \Phi_p \subseteq \Phi_q \cdot R \quad (7.27)$$

a definition equivalent to

$$p \xrightarrow{R} q \equiv R \cdot \Phi_p \subseteq \Phi_q \cdot \top \quad (7.28)$$

as seen above.

7.5 ABSTRACT INTERPRETATION

The proofs involved in verifying contracts may be hard to perform due to the intricacies of real-life sized software specifications, which may involve hundreds of invariants of arbitrary complexity. Such situations can only be tackled with the support of a theorem prover, and in many situations even this is not enough to accomplish the task. This problem has made software theorists to think of strategies helping designers to simplify their proofs. One such strategy is *abstract interpretation*.

It is often the case that the proof of a given contract does not require the whole model because the contract is only concerned with a particular *view* of the whole thing. As a very simple example, think of a model that is made of two independent parts $A \times B$ and of an invariant that constrains part A only. Then one may safely ignore B in the proofs. This is equivalent to applying projection $\pi_1 : A \times B \rightarrow A$ (2.21) to the original model. Note that π_1 is an *abstraction*, since it is a surjective function (recall figure 5.3).

In general, software models are not as “separable” as $A \times B$ is, but abstraction functions exist that yield much simpler models where the proofs can be made easier. Different abstractions help in different proofs — a kind of “on demand” *abstraction* making a model more *abstract* with respect to the *specific* property one wishes to check. In general, techniques of this kind are known as *abstract interpretation* techniques and play a major role in *program analysis*, for instance. To explain abstract interpretation we need to introduce the notion of a *relational type*.

RELATIONS AS TYPES A function h is said to have *relation type* $R \rightarrow S$, written $R \xrightarrow{h} S$ if

$$h \cdot R \subseteq S \cdot h \quad \begin{array}{ccc} B & \xleftarrow{R} & B \\ h \downarrow & & \downarrow h \\ A & \xleftarrow{S} & A \end{array} \quad (7.29)$$

holds. Note that (7.29) could be written $h (S \leftarrow R) h$ in the notation of (6.10). In case $h : B \rightarrow A$ is surjective, i.e. h is an *abstraction function*, we also say that $A \xleftarrow{S} A$ is an *abstract simulation* of $B \xleftarrow{R} B$ through h .

A special case of relational type defines so-called *invariant functions*. A function of relation type $R \xrightarrow{h} id$ is said to be *R-invariant*, in the sense that

$$\langle \forall b, a : b R a : h b = h a \rangle \quad (7.30)$$

holds. When h is *R-invariant*, observations by h are not affected by *R*-transitions. In pointfree notation, an *R-invariant* function h is always such that:

$$R \subseteq \frac{h}{h} \quad (7.31)$$

For instance, a binary operation θ is *commutative* iff θ is *swap-invariant*, that is

$$\text{swap} \subseteq \frac{\theta}{\theta} \quad (7.32)$$

holds.

Exercise 7.5. What does (7.29) mean in case R and S are partial orders?

□

Exercise 7.6. Show that relational types compose, that is $Q \xleftarrow{k} S$ and $S \xleftarrow{h} R$ entail $Q \xleftarrow{k \cdot h} R$.

□

Exercise 7.7. Show that an alternative way of stating (7.27) is

$$p \xrightarrow{R} q \equiv R \cdot \Phi_p \subseteq \Phi_q \cdot \top \quad (7.33)$$

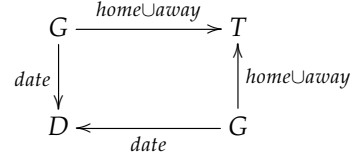
□

Exercise 7.8. Recalling exercise 5.12, let the following relation specify that two dates are at least one week apart in time:

$$d \text{ Ok } d' \Leftrightarrow |d - d'| > 1 \text{ week}$$

Looking at the type diagram below, say in your own words the meaning of the invariant specified by the relational type (7.29) statement below, on the left:

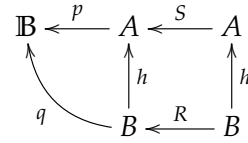
$$\ker (\text{home} \cup \text{away}) - \text{id} \xrightarrow{\text{date}} \text{Ok}$$



□

ABSTRACT INTERPRETATION Suppose that one wishes to show that $q : B \rightarrow \mathbb{B}$ is an invariant of some operation $B \xrightarrow{R} B$, i.e. that $q \xrightarrow{R} q$ holds and you know that $q = p \cdot h$, for some $h : B \rightarrow A$, as shown in the diagram. Then one can factor the proof in two steps:

- show that there is an abstract *simulation* S such that $R \xrightarrow{h} S$;
- prove $p \xrightarrow{S} p$, that is, that p is an (abstract) *invariant* of (abstract) S .



This strategy is captured by the following calculation:

$$\begin{aligned}
 & R \cdot \Phi_q \subseteq \Phi_q \cdot \top \\
 \equiv & \{ q = p \cdot h \} \\
 & R \cdot \Phi_{(p \cdot h)} \subseteq \Phi_{(p \cdot h)} \cdot \top \\
 \equiv & \{ (5.205) \text{ etc } \} \\
 & R \cdot \Phi_{(p \cdot h)} \subseteq h^\circ \cdot \Phi_p \cdot \top \\
 \equiv & \{ \text{shunting} \} \\
 & h \cdot R \cdot \Phi_{(p \cdot h)} \subseteq \Phi_p \cdot \top \\
 \Leftarrow & \{ R \xrightarrow{h} S \} \\
 & S \cdot h \cdot \Phi_{(p \cdot h)} \subseteq \Phi_p \cdot \top \\
 \Leftarrow & \{ \Phi_{(p \cdot h)} \subseteq h^\circ \cdot \Phi_p \cdot h (5.210) \} \\
 & S \cdot h \cdot h^\circ \cdot \Phi_p \cdot h \subseteq \Phi_p \cdot \top
 \end{aligned}$$

$$\begin{aligned} & \Leftarrow \{ \top = \top \cdot h \text{ (cancel } h); \text{img } h \subseteq id \} \\ & S \cdot \Phi_p \subseteq \Phi_p \cdot \top \\ & \square \end{aligned}$$

Abstract interpretation techniques usually assume that h is an adjoint of a Galois connection. Our first examples below do not assume this, for an easy start.

7.6 SAFETY AND LIVENESS PROPERTIES

Before showing examples of abstract interpretation, let us be more specific about what was meant by “some operation $B \xrightarrow{R} B$ ” above. In section 4.9 a monad was studied called the *state monad*. This monad is inhabited by state-transitions encoding state-based automata known as *Mealy machines*.

With relations one may be more relaxed on how to characterize state automata. In general, functional models generalize to so called *state-based* relational models in which there is

- a set Σ of *states*
- a subset $I \subseteq \Sigma$ of *initial* states
- a *step* relation $\Sigma \xrightarrow{R} \Sigma$ which expresses transition of states.

We define:

- $R^0 = id$ — no action or transition takes place
- $R^{i+1} = R \cdot R^i$ — all “paths” made of $i + 1$ R -transitions
- $R^* = \bigcup_{i \geq 0} R^i$ — the set of all possible R -paths.

We represent the set I of initial states by the coreflexive $\Sigma \xrightarrow{\Phi(\in I)} \Sigma$, simplified to $\Sigma \xrightarrow{I} \Sigma$ to avoid symbol cluttering.

Given $\Sigma \xrightarrow{R, I} \Sigma$ (i.e. a nondeterministic automaton, model) there are two kinds of property that one may wish to prove — *safety* and *liveness* properties. *Safety* properties are of the form $R^* \cdot I \subseteq S$, that is,

$$\langle \forall n : n \geq 0 : R^n \cdot I \subseteq S \rangle \quad (7.34)$$

for some safety relation $S : \Sigma \rightarrow \Sigma$, meaning: *All paths in the model originating from its initial states are bounded by S*. In the particular case $S = \frac{true}{p}$ ⁸

$$\langle \forall n : n \geq 0 : R^n \cdot I \subseteq \frac{true}{p} \rangle \quad (7.35)$$

⁸ Recall that $\frac{true}{p} = \Phi_p \cdot \top$ (5.205).

meaning that formula p holds for every state reachable by R from an initial state. Invariant preservation is an example of a safety property: if starting from a “good” state, the automaton only visits “good” (valid) states.

In contrast to safety properties, the so-called *liveness* properties are of the form

$$\langle \exists n : n \geq 0 : Q \subseteq R^n \cdot I \rangle \quad (7.36)$$

for some *target* relation $Q : \Sigma \rightarrow \Sigma$, meaning: *the target relation Q is eventually realizable, after n steps starting from an initial state*. In the particular case $Q = \frac{\text{true}}{p}$ we have

$$\langle \exists n : n \geq 0 : \frac{\text{true}}{p} \subseteq R^n \cdot I \rangle \quad (7.37)$$

meaning that, for a sufficiently large n , formula p will eventually hold.

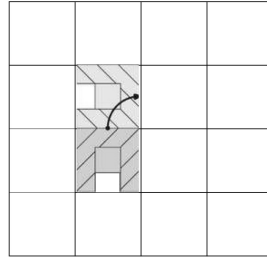
7.7 EXAMPLES

The Alcuin puzzle is an example of a problem that is characterized by a liveness and safety property:

- From initial state *where* = Left, state *where* = Right is eventually reachable — a *liveness* property.
- Initial state *where* = Left is valid and no step of the automaton leads to invalid *where* states — a *safety* property.

The first difficulty in ensuring properties such as (7.35) e (7.37) is the quantification on the number of path steps. In the case of (7.37) one can try and find a particular path using a *model checker*. In both cases, the complexity /size of the *state space* may offer some impedance to proving / model checking. Below we show how to circumvent such difficulties by use of *abstract interpretation*.

THE HEAVY ARMCHAIR PROBLEM Let us show a simple, but effective example of abstract interpretation applied to a well-known problem — the *heavy armchair* problem.⁹ Consider the following picture:



⁹ Credits: this version of the problem and the pictures shown are taken from [6].

We wish to move the armchair to an adjacent square, horizontally or vertically. However, because the armchair is too heavy, it can only be rotated over one of its four legs, as shown in the picture.

The standard model for this problem is a pair (p, o) where $p = (y, x)$ captures the square where the armchair is positioned and o is one of the complex numbers $\{i, -i, 1, -1\}$ indicating the orientation of the armchair (that is, it can face N,S,E,W). Let the following the step-relation be proposed,

$$R = P \times Q$$

where P captures the *adjacency* of two squares and Q captures 90° rotations. A *rotation* multiplies an orientation o by $\pm i$, depending on choosing a clockwise ($-i$) or anti-clockwise (i) rotation. Altogether:

$$\begin{aligned} ((y', x'), d') R ((y, x), d) &\Leftrightarrow \\ \left\{ \begin{array}{l} y' = y \pm 1 \wedge x' = x \vee y' = y \wedge x' = x \pm 1 \\ d' = (\pm i) d \end{array} \right. \end{aligned}$$

We want to check the *liveness* property:

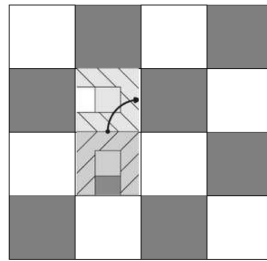
$$\text{For some } n, ((y, x + 1), d) R^n ((y, x), d) \text{ holds.} \quad (7.38)$$

That is, we wish to move the armchair to the adjacent square on its right, keeping the armchair's orientation. This is exactly what the pointfree version of (7.38) tells:

$$\langle \exists n :: (id \times (1+)) \times id \subseteq R^n \rangle$$

In other words: *there is a path with n steps that realizes the function $move = (id \times (1+)) \times id$.*

Note that the state of this problem is arbitrarily large. (The squared area is unbounded.) Moreover, the specification of the problem is non-deterministic. (For each state, there are four possible successor states.) We resort to *abstract interpretation* to obtain a bounded, deterministic (*functional*) model: the floor is coloured as a chess board and the armchair behaviour is abstracted by function $h = col \times dir$ which tells the *colour* of the square where the armchair is and the *direction* of its current orientation:



Since there are two colours (black, white) and two directions (horizontal, vertical), both can be modelled by Booleans. Then the action of

moving to any adjacent square abstracts to *color* negation and any 90° rotation abstracts to *direction* negation:

$$P \xrightarrow{col} (\neg) \quad (7.39)$$

$$Q \xrightarrow{dir} (\neg) \quad (7.40)$$

In detail:

$$col(y, x) = even(y + x)$$

$$dir\ x = x \in \{1, -1\}$$

For instance, $col(0, 0) = \text{True}$ (black in the picture), $col(1, 1) = \text{True}$, $col(1, 2) = \text{False}$ and so on; $dir\ 1 = \text{True}$ (horizontal orientation), $dir(-i) = \text{False}$, and so on. Checking (7.40):

$$\begin{aligned} & dir((\pm i)\ x) \\ = & \{ \ dir\ x = x \in \{1, -1\} \} \\ & (\pm i)\ x \in \{1, -1\} \\ = & \{ \ multiply\ by\ (\pm i)\ \text{within}\ \{1, i, -1, -i\} \} \\ & x \in \{-i, i\} \\ = & \{ \ the\ remainder\ of\ \{-i, i\}\ \text{is}\ \{1, -1\} \} \\ & \neg(x \in \{1, -1\}) \\ = & \{ \ dir\ x = x \in \{1, -1\} \} \\ & \neg(dir\ x) \end{aligned}$$

□

Checking (7.39):

$$\begin{aligned} & (\neg) \xleftarrow{col} P \\ \equiv & \{ \ (7.29)\ \text{for functions} \} \\ & col \cdot P \subseteq \neg \cdot col \\ \equiv & \{ \ shunting ; go\ pointwise \} \\ & (y', x')\ P(y, x) \Rightarrow even(y' + x') = \neg even(y + x) \\ \equiv & \{ \ unfold \} \\ & \begin{cases} y' = y \pm 1 \wedge x' = x \Rightarrow even(y' + x') = \neg even(y + x) \\ y' = y \wedge x' = x \pm 1 \Rightarrow even(y' + x') = \neg even(y + x) \end{cases} \\ \equiv & \{ \ substitutions ; trivia \} \\ & \begin{cases} even(y \pm 1) = \neg even\ y \\ even(x \pm 1) = \neg even\ x \end{cases} \\ \equiv & \{ \ trivia \} \\ & true \end{aligned}$$

□

Altogether:

$$R \xrightarrow{\text{col} \times \text{dir}} (\neg \times \neg)$$

That is, step relation R is simulated by $s = \neg \times \neg$, i.e. the function

$$s(c, d) = (\neg c, \neg d)$$

over a state space with 4 possibilities only: wherever the armchair turns over one of its legs, whatever this is, it changes *both* the colour of the square where it is, and its direction.

At this level, we note that *observation* function

$$f(c, d) = c \oplus d \tag{7.41}$$

is *s*-invariant (7.30), that is

$$f \cdot s = f \tag{7.42}$$

since $\neg c \oplus \neg d = c \oplus d$ holds. By induction on n , $f \cdot s^n = f$ holds too.

Expressed under this abstraction, (7.38) is rephrased into: *there is a number of steps n such that $s^n(c, d) = (\neg c, d)$ holds.* Let us check this abstract version of the original property, assuming variable n existentially quantified:

$$\begin{aligned} & s^n(c, d) = (\neg c, d) \\ \Rightarrow & \quad \{ \text{Leibniz} \} \\ & f(s^n(c, d)) = f(\neg c, d) \\ \equiv & \quad \{ f \text{ is } s\text{-invariant} \} \\ & f(c, d) = f(\neg c, d) \\ \equiv & \quad \{ (7.41) \} \\ & c \oplus d = \neg c \oplus d \\ \equiv & \quad \{ 1 \oplus d = \neg d \text{ and } 0 \oplus d = d \} \\ & d = \neg d \\ \equiv & \quad \{ \text{trivia} \} \\ & \text{false} \end{aligned}$$

Thus, for all paths of arbitrary length n , $s^n(c, d) \neq (\neg c, d)$. We conclude that the proposed liveness property does not at all hold!

ALCUIN PUZZLE EXAMPLE Abstract interpretation applies nicely to this problem, thanks to its symmetries. On the one hand, one does not need to work over the 16 functions in $\text{Bank}^{\text{Being}}$, since starting from

the left margin or from the right margin is irrelevant. Another symmetry can be found in type *Being*, suggesting the following abstraction of beings into three classes:

$$f : \text{Being} \rightarrow \{\alpha, \beta, \gamma\}$$

$$f = \begin{pmatrix} \text{Goose} \longrightarrow \alpha \\ \text{Fox} \longrightarrow \beta \\ \text{Beans} \nearrow \beta \\ \text{Farmer} \longrightarrow \gamma \end{pmatrix}$$

The abstraction consists in unifying , the maximum and minimum elements of the “food chain”. In fact, the simultaneous presence of one α and one β is enough for defining the invariant — which *specific* being eats the other is irrelevant detail. This double abstraction is captured by

$$\begin{array}{ccc} \text{Bank} & \xleftarrow{w} & \text{Being} \\ \text{Left} \uparrow & & f \downarrow \\ 1 & \xleftarrow[V]{} & \{\alpha, \beta, \gamma\} \end{array} \quad V = \underline{\text{Left}}^\circ \cdot w \cdot f^\circ$$

where the choice of *Left* as reference bank is arbitrary. Thus function w is abstracted by the row *vector* relation V ¹⁰ such that:

$$_ V x = \langle \exists b : x = f b : w b = \text{Left} \rangle$$

Vector V tells whether at least one being of class x can be found in the reference bank. Noting that there could be more than one β there, we refine the abstraction a bit so that the number of beings of each class is counted.¹¹ This leads to the following *state-abstraction* (higher order) function h based on f :

$$\begin{aligned} h : (\text{Being} \rightarrow \text{Bank}) &\rightarrow \{\alpha, \beta, \gamma\} \rightarrow \{0, 1, 2\} \\ h w x &= \langle \sum b : x = f b \wedge w b = \text{Left} : 1 \rangle \end{aligned}$$

For instance,

$$\begin{aligned} h \underline{\text{Left}} &= 121 \\ h \underline{\text{Right}} &= 000 \end{aligned}$$

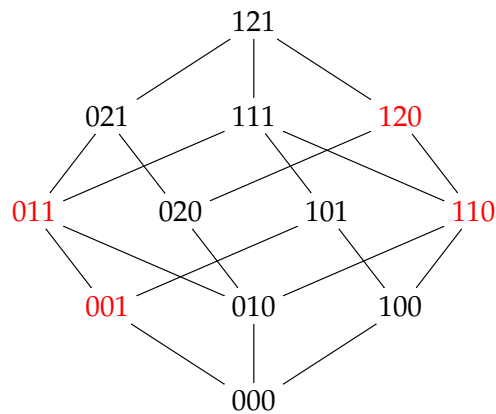
abbreviating by vector xyz the mapping $\{\alpha \mapsto x, \beta \mapsto y, \gamma \mapsto z\}$.¹² To obtain the other bank just compute: $\bar{x} = 121 - x$. Note that there are

¹⁰ A fragment of $! : \{\alpha, \beta, \gamma\} \rightarrow 1$, recall section 5.5.

¹¹ This suggests that linear algebra would be a good alternative to relation algebra here!

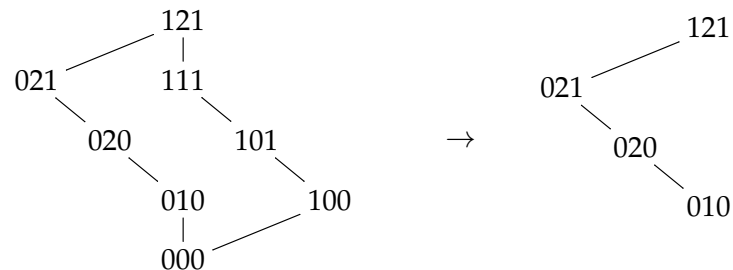
¹² This version of the model is inspired in [6].

$2 \times 3 \times 2 = 12$ possible state vectors, 4 of which are invalid (these are marked in red):

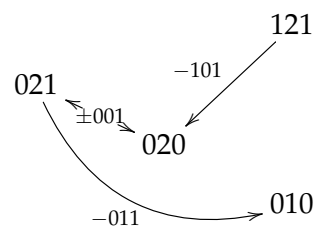


The ordering implicit in the lattice above is pointwise (\leq). This is complemented by $\bar{x} = 121 - x$, which gives the information of the other bank.

The 8 valid states can be further abstracted to only 4 of them,



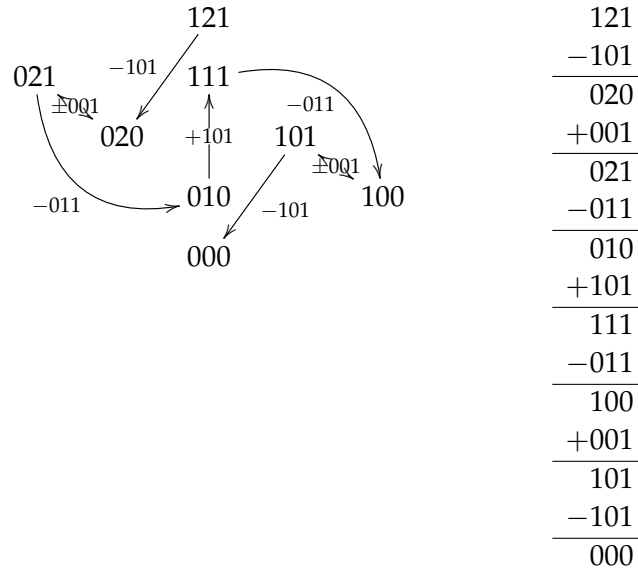
since, due to complementation (cf. the Left-Right margin symmetry), we only need to reach state 010. Then we reverse the path through the complements. In this setting, the automaton is deterministic, captured by the abstract automaton:



Termination is ensured by disabling toggling between states 021 and 020:

$$\begin{array}{r}
 121 \\
 -101 \\
 \hline
 020 \\
 +001 \\
 \hline
 021 \\
 -011 \\
 \hline
 010
 \end{array}$$

We then take the complemented path $111 \rightarrow 100 \rightarrow 101 \rightarrow 000$. So the abstract solution for the Alcuin puzzle is, finally:



At this point note that, according to the principles of abstract interpretation stated above, quite a few steps are pending in this exercise: abstract the *starving* invariant to the vector level, find an abstract simulation of *carry*, and so on and so forth. But — why bother doing all that? There no other operation in the problem, so the abstraction found is, in a sense, universal: we should have started from the vector model and not from the *Being* \rightarrow *Bank* model, which is not *sufficiently* abstract.

The current scientific basis of programming enables the calculation of programs, following the scientific method. So, programming is lesser and lesser an *art*. Where is creativity gone to? To the *art* of abstract modelling and elegant proving — this is where it can be found nowadays.

Exercise 7.9. Verification of code involves calculations of real numbers and is often done on the basis of an abstract interpretation called *sign analysis*:

$$\begin{aligned} \text{sign} : \mathbb{R} &\rightarrow \{-, 0, +\} \\ \text{sign } 0 &= 0 \\ \text{sign } x &= \text{if } x > 0 \text{ then } + \text{ else } - \end{aligned}$$

Suppose there is evidence that the operation $\theta : \{-, 0, +\}^2 \rightarrow \{-, 0, +\}$ defined by

θ	-	0	+
-	+	0	-
0	0	0	0
+	-	0	+

(7.43)

is the abstract simulation induced by *sign* of a given concrete operation $f : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$, that is, that

$$\theta \cdot (\text{sign} \times \text{sign}) = \text{sign} \cdot f \quad (7.44)$$

holds. It is easy to see, by inspection of (7.43), that θ is a commutative operation, recalling (7.32).

- Show that $\text{sign} \cdot f$ is necessarily commutative as well. (Hint: the free theorem of swap can be useful here.)
- Does the previous question guarantee that the specific operation f is also commutative? Answer informally.

□

7.8 “FREE CONTRACTS”

In design by contract, many functional *contracts* arise naturally as corollaries of *free theorems*. This has the advantage of saving us from proving such contracts explicitly.

The following exercises provide ample evidence of this.

Exercise 7.10. The type of functional composition (\cdot) is

$$(b \rightarrow c) \rightarrow (a \rightarrow b) \rightarrow a \rightarrow c$$

Show that contract composition (7.18) is a corollary of the free theorem (FT) of this type.

□

Exercise 7.11. Show that contract $q^* \xleftarrow{\text{map } f} p^*$ holds provided contract $q \xleftarrow{f} p$ holds.

□

Exercise 7.12. Suppose a functional programmer wishes to prove the following property of lists:

$$\langle \forall a, s : (p\ a) \wedge \langle \forall a' : a' \in \text{elems } s : p\ a' \rangle : \langle \forall a'' : a'' \in \text{elems } (a : s) : p\ a'' \rangle \rangle$$

Show that this property is a contract arising (for free) from the polymorphic type of the `cons` operation $(:)$ on lists.

□

7.9 REASONING BY APPROXIMATION

Currently in preparation

7.10 BIBLIOGRAPHY NOTES

To be completed

PROGRAMS AS RELATIONAL HYLOMORPHISMS

not given in the current version of this textbook

CALCULATIONAL PROGRAM REFINEMENT

not given in the current version of this textbook

Part III

CALCULATING WITH MATRICES

10

TOWARDS A LINEAR ALGEBRA OF PROGRAMMING

This part of the book will build upon references [50, 45, 54]. Another chapter will address the application of typed linear algebra to analytical data processing, cf. e.g. [49].



BACKGROUND — EINDHOVEN QUANTIFIER CALCULUS

This appendix is a quick reference summary of section 4.3 of reference [4].

A.1 NOTATION

The Eindhoven quantifier calculus adopts the following notation standards:

- $\langle \forall x : R : T \rangle$ means: “for *all* x in the range R , term T holds”, where R and T are logical expressions involving x .
- $\langle \exists x : R : T \rangle$ means: “for *some* x in the range R , term T holds”.

A.2 RULES

The main rules of the Eindhoven quantifier calculus are listed below:

Trading:

$$\langle \forall k : R \wedge S : T \rangle = \langle \forall k : R : S \Rightarrow T \rangle \quad (\text{A.1})$$

$$\langle \exists k : R \wedge S : T \rangle = \langle \exists k : R : S \wedge T \rangle \quad (\text{A.2})$$

de Morgan:

$$\neg \langle \forall k : R : T \rangle = \langle \exists k : R : \neg T \rangle \quad (\text{A.3})$$

$$\neg \langle \exists k : R : T \rangle = \langle \forall k : R : \neg T \rangle \quad (\text{A.4})$$

One-point:

$$\langle \forall k : k = e : T \rangle = T[k := e] \quad (\text{A.5})$$

$$\langle \exists k : k = e : T \rangle = T[k := e] \quad (\text{A.6})$$

Nesting:

$$\langle \forall a, b : R \wedge S : T \rangle = \langle \forall a : R : \langle \forall b : S : T \rangle \rangle \quad (\text{A.7})$$

$$\langle \exists a, b : R \wedge S : T \rangle = \langle \exists a : R : \langle \exists b : S : T \rangle \rangle \quad (\text{A.8})$$

Rearranging- \forall :

$$\langle \forall k : R \vee S : T \rangle = \langle \forall k : R : T \rangle \wedge \langle \forall k : S : T \rangle \quad (\text{A.9})$$

$$\langle \forall k : R : T \wedge S \rangle = \langle \forall k : R : T \rangle \wedge \langle \forall k : R : S \rangle \quad (\text{A.10})$$

Rearranging- \exists :

$$\langle \exists k : R : T \vee S \rangle = \langle \exists k : R : T \rangle \vee \langle \exists k : R : S \rangle \quad (\text{A.11})$$

$$\langle \exists k : R \vee S : T \rangle = \langle \exists k : R : T \rangle \vee \langle \exists k : S : T \rangle \quad (\text{A.12})$$

Splitting:

$$\langle \forall j : R : \langle \forall k : S : T \rangle \rangle = \langle \forall k : \langle \exists j : R : S \rangle : T \rangle \quad (\text{A.13})$$

$$\langle \exists j : R : \langle \exists k : S : T \rangle \rangle = \langle \exists k : \langle \exists j : R : S \rangle : T \rangle \quad (\text{A.14})$$

B

HASKELL SUPPORT LIBRARY

This library, written in the Haskell functional programming language, is still evolving.

```
infix 5 ×  
infix 4 +
```

Products

$$\begin{aligned}\langle \cdot, \cdot \rangle &:: (a \rightarrow b) \rightarrow (a \rightarrow c) \rightarrow a \rightarrow (b, c) \\ \langle f, g \rangle x &= (f\ x, g\ x) \\ (\times) &:: (a \rightarrow b) \rightarrow (c \rightarrow d) \rightarrow (a, c) \rightarrow (b, d) \\ f \times g &= \langle f \cdot \pi_1, g \cdot \pi_2 \rangle\end{aligned}$$

The 0-adic split is the unique function of its type

$$\begin{aligned}(!) &:: a \rightarrow () \\ (!) &= \underline{\quad}\end{aligned}$$

Renamings:

$$\begin{aligned}\pi_1 &= \text{fst} \\ \pi_2 &= \text{snd}\end{aligned}$$

Coproduct

Renamings:

$$\begin{aligned}i_1 &= \text{Left} \\ i_2 &= \text{Right}\end{aligned}$$

Either is predefined:

$$\begin{aligned}(+) &:: (a \rightarrow b) \rightarrow (c \rightarrow d) \rightarrow a + c \rightarrow b + d \\ f + g &= [i_1 \cdot f, i_2 \cdot g]\end{aligned}$$

McCarthy's conditional:

$$p \rightarrow f, g = [f, g] \cdot p?$$

Exponentiation

Curry is predefined.

$$\begin{aligned} ap &:: (a \rightarrow b, a) \rightarrow b \\ ap &= (\widehat{\$}) \end{aligned}$$

Functor:

$$\begin{aligned} \cdot &:: (b \rightarrow c) \rightarrow (a \rightarrow b) \rightarrow a \rightarrow c \\ f\cdot &= \overline{f \cdot ap} \end{aligned}$$

Pair-to-predicate isomorphism (2.99):

$$\begin{aligned} p2p &:: (b, b) \rightarrow \mathbb{B} \rightarrow b \\ p2p\ p\ b &= \text{if } b \text{ then } (\pi_2\ p) \text{ else } (\pi_1\ p) \end{aligned}$$

The exponentiation functor is $(a \rightarrow)$ predefined:

$$\begin{aligned} &\textbf{instance Functor } ((\rightarrow)\ s) \textbf{ where} \\ &fmap\ f\ g = f \cdot g \end{aligned}$$

Guards

$$\begin{aligned} \cdot? &:: (a \rightarrow \mathbb{B}) \rightarrow a \rightarrow a + a \\ p? \ x &= \text{if } p\ x \text{ then } i_1\ x \text{ else } i_2\ x \end{aligned}$$

Others

$\underline{\cdot} :: a \rightarrow b \rightarrow a$ such that $\underline{a}\ x = a$ is predefined.

Natural isomorphisms

$$\begin{aligned} \text{swap} &:: (a, b) \rightarrow (b, a) \\ \text{swap} &= \langle \pi_2, \pi_1 \rangle \\ \text{assocr} &:: ((a, b), c) \rightarrow (a, (b, c)) \\ \text{assocr} &= \langle \pi_1 \cdot \pi_1, \text{snd} \times id \rangle \\ \text{assocl} &:: (a, (b, c)) \rightarrow ((a, b), c) \\ \text{assocl} &= \langle id \times \pi_1, \pi_2 \cdot \pi_2 \rangle \\ \text{undistr} &:: a, b + a, c \rightarrow (a, b + c) \\ \text{undistr} &= [id \times i_1, id \times i_2] \\ \text{undistl} &:: b, c + a, c \rightarrow (b + a, c) \\ \text{undistl} &= [i_1 \times id, i_2 \times id] \\ \text{coswap} &:: a + b \rightarrow b + a \\ \text{coswap} &= [i_2, i_1] \\ \text{coassocr} &:: a + b + c \rightarrow a + b + c \\ \text{coassocr} &= [id + i_1, i_2 \cdot i_2] \\ \text{coassocl} &:: b + a + c \rightarrow b + a + c \end{aligned}$$

```

coassocl = [i1 · i1, i2 + id]
distl :: (c + a, b) → c, b + a, b
distl =  $\widehat{[i_1, i_2]}$ 
distr :: (b, c + a) → b, c + b, a
distr = (swap + swap) · distl · swap
flatr :: (a, (b, c)) → (a, b, c)
flatr (a, (b, c)) = (a, b, c)
flatl :: ((a, b), c) → (a, b, c)
flatl ((b, c), d) = (b, c, d)
br = ⟨id, !⟩
bl = swap · br

```

Class bifunctor

```

class BiFunctor f where
  bmap :: (a → b) → (c → d) → (f a c → f b d)
instance BiFunctor (· + ·) where
  bmap f g = f + g
instance BiFunctor (· × ·) where
  bmap f g = f × g

```

Monads

Kleisli monadic composition:

```

infix 4 •
(•) :: Monad a ⇒ (b → a c) → (d → a b) → d → a c
(f • g) a = (g a) >>= f

```

Multiplication, also known as join:

```

mult :: (Monad m) ⇒ m (m b) → m b
mult = (>>= id)

```

Monadic binding:

```

ap' :: (Monad m) ⇒ (a → m b, m a) → m b
ap' = flip  $\widehat{(>>=)}$ 

```

List monad:

```

singl :: a → [a]
singl = return

```

Strong monads:

```

class (Functor f, Monad f) ⇒ Strong f where
  rstr :: (f a, b) → f (a, b)
  rstr (x, b) = do a ← x; return (a, b)

```

```

lstr :: (b, f a) → f (b, a)
lstr (b, x) = do a ← x; return (b, a)
instance Strong IO
instance Strong []
instance Strong Maybe

```

Double strength:

```

dstr :: Strong m ⇒ (m a, m b) → m (a, b)
dstr = rstr • lstr

```

Exercise 4.8.13 in Jacobs' "Introduction to Coalgebra" [28]:

```

splitm :: Strong F ⇒ F (a → b) → a → F b
splitm = fmap ap · rstr

```

Monad transformers:

```

class (Monad m, Monad (t m)) ⇒ MT t m where    -- monad transformer class
    lift :: m a → t m a

```

Nested lifting:

```

dlift :: (MT t (t1 m), MT t1 m) ⇒ m a → t (t1 m) a
dlift = lift · lift

```

Basic functions, abbreviations

```

zero = 0
one = 1
nil = []
cons = ∘
add = +
mul = *
conc = ++
inMaybe :: +a → Maybe a
inMaybe = [Nothing, Just]

```

More advanced

```

class (Functor f) ⇒ Unzipable f where
    unzp :: f (a, b) → (f a, f b)
    unzp = ⟨fmap π1, fmap π2⟩
class Functor g ⇒ DistL g where
    λ :: Monad m ⇒ g (m a) → m (g a)
instance DistL [] where λ = sequence
instance DistL Maybe where

```

```

λ Nothing = return Nothing
λ (Just a) = mp Just a where mp f = (return · f) • id

```

Convert Monad into Applicative:

```

aap :: Monad m => m (a -> b) -> m a -> m b
aap mf mx = do {f <- mf; x <- mx; return (f x)}

```

BIBLIOGRAPHY

- [1] C. Aarts, R.C. Backhouse, P. Hoogendijk, E. Voermans, and J. van der Woude. A relational theory of datatypes, December 1992. Available from www.cs.nott.ac.uk/~rcb.
- [2] K. Backhouse and R.C. Backhouse. Safety of abstract interpretations for free, via logical relations and Galois connections. *SCP*, 15(1–2):153–196, 2004.
- [3] R.C. Backhouse. On a relation on functions. In *Beauty is our business: a birthday salute to Edsger W. Dijkstra*, pages 7–18, New York, NY, USA, 1990. Springer-Verlag.
- [4] R.C. Backhouse. *Mathematics of Program Construction*. Univ. of Nottingham, 2004. Draft of book in preparation. 608 pages.
- [5] R.C. Backhouse and M.M. Fokkinga. The associativity of equivalence and the Towers of Hanoi problem. *Information Processing letters*, 77(2–4):71–76, 2001.
- [6] Roland Backhouse. *Algorithmic Problem Solving*. Wiley Publishing, 1st edition, 2011.
- [7] J. Backus. Can programming be liberated from the von Neumann style? a functional style and its algebra of programs. *CACM*, 21(8):613–639, August 1978.
- [8] L.S. Barbosa. *Components as Coalgebras*. University of Minho, December 2001. Ph. D. thesis.
- [9] R. Bird. Introduction to Functional Programming. Series in Computer Science. Prentice-Hall International, 2nd edition, 1998. C.A.R. Hoare, series editor.
- [10] R. Bird and O. de Moor. *Algebra of Programming*. Series in Computer Science. Prentice-Hall, 1997.
- [11] R.M. Burstall and J. Darlington. A transformation system for developing recursive programs. *JACM*, 24(1):44–67, January 1977.
- [12] V. Cerf. Where is the science in computer science? *CACM*, 55(10):5, October 2012.
- [13] M. Erwig and S. Kollmannsberger. Functional pearls: Probabilistic functional programming in Haskell. *J. Funct. Program.*, 16:21–34, January 2006.
- [14] S. Feferman. Tarski’s influence on computer science. *Logical Methods in Computer Science*, 2:1–1–13, 2006.
- [15] R.W. Floyd. Assigning meanings to programs. In J.T. Schwartz, editor, *Mathematical Aspects of Computer Science*, volume 19, pages 19–32. American Mathematical Society, 1967. Proc. Symposia in Applied Mathematics.
- [16] M.M. Fokkinga. *Law and Order in Algorithmics*. PhD thesis, University of Twente, Dept INF, Enschede, The Netherlands, 1992.

- [17] P.J. Freyd and A. Scedrov. *Categories, Allegories*, volume 39 of *Mathematical Library*. North-Holland, 1990.
- [18] J. Gibbons. Kernels, in a nut shell. *JLAMP*, 85(5, Part 2):921–930, 2016.
- [19] J. Gibbons and R. Hinze. Just do it: simple monadic equational reasoning. In *Proceedings of the 16th ACM SIGPLAN international conference on Functional programming*, ICFP’11, pages 2–14, New York, NY, USA, 2011. ACM.
- [20] Jeremy Gibbons, Graham Hutton, and Thorsten Altenkirch. When is a function a fold or an unfold?, 2001. WGP, July 2001 (slides).
- [21] S. Givant. The calculus of relations as a foundation for mathematics. *J. Autom. Reasoning*, 37(4):277–322, 2006.
- [22] A.S. Green, P.L. Lumsdaine, N.J. Ross, P. Selinger, and B. Valiron. An introduction to quantum programming in Quipper. *CoRR*, cs.PL(arXiv:1304.5485v1), 2013.
- [23] Ralf Hinze. Adjoint folds and unfolds — an extended study. *Science of Computer Programming*, 78(11):2108–2159, 2013.
- [24] Ralf Hinze. Adjoint folds and unfolds — an extended study. *Science of Computer Programming*, 78(11):2108–2159, 2013.
- [25] Ralf Hinze, Nicolas Wu, and Jeremy Gibbons. Conjugate hylomorphisms – or: The mother of all structured recursion schemes. In *Proceedings of the 42Nd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, POPL ’15, pages 527–538, New York, NY, USA, 2015. ACM.
- [26] P. Hudak. *The Haskell School of Expression - Learning Functional Programming Through Multimedia*. Cambridge University Press, 1st edition, 2000. ISBN 0-521-64408-9.
- [27] Graham Hutton and Erik Meijer. Monadic parsing in Haskell. *Journal of Functional Programming*, 8(4), 1993.
- [28] B. Jacobs. *Introduction to Coalgebra. Towards Mathematics of States and Observations*. Cambridge University Press, 2016.
- [29] P. Jansson and J. Jeuring. Polylib — a library of polytypic functions. In *Workshop on Generic Programming (WGP’98), Marstrand, Sweden*, 1998.
- [30] J. Jeuring and P. Jansson. Polytypic programming. In *Advanced Functional Programming*, number 1129 in LNCS, pages 68–114. Springer-Verlag, 1996.
- [31] S.L. Peyton Jones. *Haskell 98 Language and Libraries*. Cambridge University Press, Cambridge, UK, 2003. Also published as a Special Issue of the Journal of Functional Programming, 13(1) Jan. 2003.
- [32] R. Lämmel and J. Visser. A Strafunski Application Letter. In V. Dahl and P.L. Wadler, editors, *Proc. of Practical Aspects of Declarative Programming (PADL’03)*, volume 2562 of LNCS, pages 357–375. Springer-Verlag, January 2003.
- [33] H.D. Macedo and J.N. Oliveira. Typing linear algebra: A biproduct-oriented approach. *SCP*, 78(11):2160–2191, 2013.

- [34] S. MacLane. *Categories for the Working Mathematician*. Springer-Verlag, 1971.
- [35] R.D. Maddux. The origin of relation algebras in the development and axiomatization of the calculus of relations. *Studia Logica*, 50:421–455, 1991.
- [36] G. Malcolm. Data structures and program transformation. *Science of Computer Programming*, 14:255–279, 1990.
- [37] E.G. Manes and M.A. Arbib. *Algebraic Approaches to Program Semantics*. Texts and Monographs in Computer Science. Springer-Verlag, 1986. D. Gries, series editor.
- [38] J. McCarthy. Towards a mathematical science of computation. In C.M. Popplewell, editor, *Proc. of IFIP 62*, pages 21–28, Amsterdam-London, 1963. North-Holland Pub. Company.
- [39] S. Mehner, D. Seidel, L. Straßburger, and J. Voigtländer. Parametricity and proving free theorems for functional-logic languages. *PPDP'14*, pages 19–30, New York, NY, USA, 2014. ACM.
- [40] E. Meijer and G. Hutton. Bananas in space: Extending fold and unfold to exponential types. In S. Peyton Jones, editor, *Proceedings of Functional Programming Languages and Computer Architecture (FPCA95)*, 1995.
- [41] E. Moggi. Computational lambda-calculus and monads. In *Proceedings 4th Annual IEEE Symp. on Logic in Computer Science, LICS'89, Pacific Grove, CA, USA, 5–8 June 1989*, pages 14–23. IEEE Computer Society Press, Washington, DC, 1989.
- [42] S.-C. Mu, Z. Hu, and M. Takeichi. An injective language for reversible computation. In *MPC 2004*, pages 289–313, 2004.
- [43] S.-C. Mu and J.N. Oliveira. Programming from Galois connections. *JLAP*, 81(6):680–704, 2012.
- [44] S.C. Mu and R. Bird. Quantum functional programming, 2001. 2nd Asian Workshop on Programming Languages and Systems, KAIST, Daejeon, Korea, December 17-18, 2001.
- [45] D. Murta and J.N. Oliveira. A study of risk-aware program transformation. *SCP*, 110:51–77, 2015.
- [46] P. Naur and B. Randell, editors. *Software Engineering: Report on a conference sponsored by the NATO SCIENCE COMMITTEE, Garmisch, Germany, 7th to 11th October 1968*. Scientific Affairs Division, NATO, 1969.
- [47] P. Nunes. *Libro de Algebra en Arithmetica y Geometria*. Original edition by Arnoldo Birckman (Anvers), 1567.
- [48] A. Oettinger. The hardware-software complementarity. *Commun. ACM*, 10:604–606, October 1967.
- [49] J. N. Oliveira and H. D. Macedo. The data cube as a typed linear algebra operator. In *Proc. of the 16th Int. Symposium on Database Programming Languages, DBPL '17*, pages 6:1–6:11, New York, NY, USA, 2017. ACM.
- [50] J.N. Oliveira. Towards a linear algebra of programming. *FAoC*, 24(4-6):433–458, 2012.

- [51] J.N. Oliveira. Lecture notes on relational methods in software design, 2015. Available from ResearchGate:
https://www.researchgate.net/profile/Jose_Oliveira34.
- [52] J.N. Oliveira. Programming from metaphorisms. *Journal of Logical and Algebraic Methods in Programming*, 94:15–44, January 2018.
- [53] J.N. Oliveira and M.A. Ferreira. Alloy meets the algebra of programming: A case study. *IEEE Trans. Soft. Eng.*, 39(3):305–326, 2013.
- [54] J.N. Oliveira and V.C. Miraldo. “Keep definition, change category” — a practical approach to state-based system calculi. *JLAMP*, 85(4):449–474, 2016.
- [55] M.S. Paterson and C.E. Hewitt. Comparative schematology. In *Project MAC Conference on Concurrent Systems and Parallel Computation*, pages 119–127, August 1970.
- [56] J.C. Reynolds. Types, abstraction and parametric polymorphism. *Information Processing 83*, pages 513–523, 1983.
- [57] G. Schmidt. *Relational Mathematics*. Number 132 in Encyclopedia of Mathematics and its Applications. Cambridge University Press, November 2010.
- [58] A. Tarski and S. Givant. *A Formalization of Set Theory without Variables*. American Mathematical Society, 1987. AMS Colloquium Publications, volume 41, Providence, Rhode Island.
- [59] G. Villavicencio and J.N. Oliveira. *Reverse Program Calculation Supported by Code Slicing*. In *Proceedings of the Eighth Working Conference on Reverse Engineering (WCRE 2001) 2-5 October 2001, Stuttgart, Germany*, pages 35–46. IEEE Computer Society, 2001.
- [60] P.L. Wadler. Theorems for free! In *4th International Symposium on Functional Programming Languages and Computer Architecture*, pages 347–359, London, Sep. 1989. ACM.
- [61] M. Winter. Arrow categories. *Fuzzy Sets and Systems*, 160(20):2893–2909, 2009.