

J.N. Oliveira

**PROGRAM DESIGN BY  
CALCULATION**

(DRAFT)

University of Minho  
*(in preparation)*



J.N. Oliveira

**PROGRAM DESIGN BY  
CALCULATION**

(DRAFT)

University of Minho  
*(in preparation)*



# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Why Program Design by Calculation? . . . . .	5
1.2	Why this Book? . . . . .	6
1.3	Book Structure . . . . .	7
1.4	How to Read This Book . . . . .	7
<b>I</b>	<b>Calculating with Functions</b>	<b>9</b>
<b>2</b>	<b>An Introduction to Pointfree Programming</b>	<b>11</b>
2.1	Introducing functions and types . . . . .	12
2.2	Functional application . . . . .	13
2.3	Functional equality and composition . . . . .	13
2.4	Identity functions . . . . .	16
2.5	Constant functions . . . . .	16
2.6	Monics and epics . . . . .	17
2.7	Isos . . . . .	19
2.8	Gluing functions which do not compose — products . . . . .	20
2.9	Gluing functions which do not compose — coproducts . . . . .	27
2.10	Mixing products and coproducts . . . . .	30
2.11	Natural properties . . . . .	33
2.12	Universal properties . . . . .	34
2.13	Guards and McCarthy’s conditional . . . . .	36
2.14	Gluing functions which do not compose — exponentials . . . . .	38
2.15	Elementary datatypes . . . . .	43
2.16	Finitary products and coproducts . . . . .	45
2.17	Initial and terminal datatypes . . . . .	47
2.18	Sums and products in HASKELL . . . . .	48

2.19 Exercises . . . . .	51
2.20 Bibliography notes . . . . .	53
<b>3 Recursion in the Pointfree Style</b>	<b>55</b>
3.1 Motivation . . . . .	55
3.2 Introducing inductive datatypes . . . . .	61
3.3 Observing an inductive datatype . . . . .	66
3.4 Synthesizing an inductive datatype . . . . .	70
3.5 Introducing (list) catas, anas and hylos . . . . .	71
3.6 Inductive types more generally . . . . .	76
3.7 Functors . . . . .	77
3.8 Polynomial functors . . . . .	79
3.9 Polynomial inductive types . . . . .	81
3.10 F-algebras and F-homomorphisms . . . . .	82
3.11 F-catamorphisms . . . . .	83
3.12 Parameterization and type functors . . . . .	86
3.13 A catalogue of standard polynomial inductive types . . . . .	90
3.14 Functors and type functors in HASKELL . . . . .	92
3.15 The mutual-recursion law . . . . .	93
3.16 “Banana-split”: a corollary of the mutual-recursion law . . . . .	101
3.17 Inductive datatype isomorphism . . . . .	102
3.18 Bibliography notes . . . . .	102
<b>4 Why Monads Matter</b>	<b>105</b>
4.1 Partial functions . . . . .	105
4.2 Putting partial functions together . . . . .	106
4.3 Lists . . . . .	108
4.4 Monads . . . . .	109
4.4.1 Properties involving (Kleisli) composition . . . . .	111
4.5 Monadic application (binding) . . . . .	112
4.6 Sequencing and the <code>do</code> -notation . . . . .	113
4.7 Generators and comprehensions . . . . .	114
4.8 Monads in HASKELL . . . . .	115
4.8.1 Monadic I/O . . . . .	116
4.9 The state monad . . . . .	117
4.10 Bibliography notes . . . . .	124

<b>II</b>	<b>Moving Away From (Pure) Functions</b>	<b>125</b>
<b>5</b>	<b>Quasi-inductive datatypes</b>	<b>127</b>
5.1	Introducing Non-inductive Datatypes . . . . .	127
5.2	Structural Induction . . . . .	133
5.3	Well-founded coalgebras and induction . . . . .	134
5.4	Datatype invariants and proof obligations . . . . .	136
5.5	Binary relations and finite mappings . . . . .	138
5.6	Finite mapping induction principle . . . . .	140
5.7	An overview of the powerset and finite mapping algebras . . . . .	141
5.8	Exercises . . . . .	141
5.9	Bibliography notes . . . . .	143
<b>III</b>	<b>Calculating with Relations</b>	<b>145</b>
<b>6</b>	<b>Introduction to Relational Calculation</b>	<b>147</b>
6.1	Functions are not enough . . . . .	147
6.2	Relational composition and converse . . . . .	151
6.3	Relational equality . . . . .	152
6.4	Specifications as “properties” . . . . .	154
6.5	Relational approach . . . . .	154
6.6	Pre/post specification style . . . . .	155
6.7	From predicates to relations . . . . .	156
6.8	Basic relational combinators . . . . .	158
6.9	Converse . . . . .	160
6.10	Meet . . . . .	161
6.11	Pointwise vs pointfree notation . . . . .	162
6.12	Orders and their taxonomy . . . . .	164
6.13	Derived combinators . . . . .	166
6.14	Entireness and simplicity . . . . .	168
6.15	Surjectiveness and injectiveness . . . . .	168
6.16	Binary relation taxonomy . . . . .	169
6.17	Reasoning about functions . . . . .	170
6.18	Galois connections . . . . .	174
6.19	Converse in a Galois connection . . . . .	177
6.20	Functions in a Galois connection . . . . .	177
6.21	Relational division . . . . .	178

6.22	Meet and join . . . . .	182
6.23	Relational <i>split</i> . . . . .	183
6.24	Relational <i>either</i> . . . . .	184
6.25	Meaning of VDM-SL specs . . . . .	188
6.26	Relational semantics of VDM-SL . . . . .	190
6.27	Relational McCarthy conditional . . . . .	191
6.28	Reasoning about VDM-SL . . . . .	192
6.29	Bibliography notes . . . . .	196
<b>7</b>	<b>An Introduction to Relational Hylomorphisms</b>	<b>197</b>
7.1	“How” does one specify? . . . . .	197
7.2	Divide-and-conquer (formally) . . . . .	198
7.3	Relators . . . . .	199
7.4	Properties of relators . . . . .	200
7.5	Equations and fixpoints . . . . .	202
7.6	Solving (Fixpoint) Equations . . . . .	203
7.7	Solving (Fixpoint) Equations III . . . . .	204
7.8	Solving relational equations . . . . .	205
7.9	Laws of the Fixpoint Calculus . . . . .	206
7.10	$\mu$ -fusion theorem . . . . .	209
7.11	Hylo(cata)-fusion . . . . .	211
7.12	Hylo(ana)-fusion . . . . .	212
7.13	Examples: VDM collective types . . . . .	213
7.14	Relational cata(ana)morphisms . . . . .	214
7.15	Inductive coreflexives . . . . .	215
7.16	Hylos as unique solutions . . . . .	218
7.17	Accessibility and membership . . . . .	219
7.18	Hylo-factorization Theorem . . . . .	222
7.19	Virtual data-structuring . . . . .	223
7.20	Final note on inductive relation $\prec$ . . . . .	224
7.21	Bibliography notes . . . . .	224
<b>8</b>	<b>Theorems for Free</b>	<b>225</b>
8.1	Parametric polymorphism: why? . . . . .	225
8.2	Free theorem of type $t$ . . . . .	226
8.3	Reynolds arrow operator . . . . .	227
8.4	Pointwise version of FT . . . . .	229
8.5	Second example: FT of $( -)$ . . . . .	230



8.6	Bibliography notes . . . . .	231
<b>IV Data Refinement by Calculation</b>		<b>233</b>
<b>9</b>	<b>On Data Representation</b>	<b>235</b>
9.1	Introduction . . . . .	235
9.2	Data refinement . . . . .	237
9.3	Representation relations . . . . .	238
9.4	Right invertibility . . . . .	239
9.5	Refinement inequations . . . . .	239
9.6	Relational representation . . . . .	241
9.7	Functional representation . . . . .	242
9.8	Concrete invariants . . . . .	242
9.9	A fundamental iso abstraction . . . . .	243
9.10	Pointfree $untot = (i_1^\circ)$ . . . . .	244
9.11	Properties of $\leq$ . . . . .	245
9.12	Structural data refinement . . . . .	246
9.13	The finite map bifunctor . . . . .	251
9.14	Transposing relations . . . . .	254
9.15	Transposing finite relations . . . . .	255
9.16	Recursive data refinement . . . . .	256
9.17	Recursion “removal” . . . . .	257
9.18	Closure and wellfoundedness . . . . .	260
9.19	Object oriented Data Implementation . . . . .	261
9.20	Multiple inheritance . . . . .	261
9.21	Bibliography notes . . . . .	262
<b>10</b>	<b>On Algorithmic Refinement</b>	<b>265</b>
10.1	Implicit/explicit refinement . . . . .	265
10.2	Handling refinement equations . . . . .	268
10.3	Solving refinement equations . . . . .	268
10.4	Properties of $\vdash$ . . . . .	269
10.5	Stepwise refinement . . . . .	270
10.6	Refinement is a partial order . . . . .	271
10.7	Main refinement strategies . . . . .	272
10.8	Data refinement in full . . . . .	273
10.9	Analysis of refinement equation . . . . .	274

10.10	Solving refinement equations . . . . .	275
10.11	Functional solutions . . . . .	275
10.12	Calculation of <code>while</code> / <code>for</code> loops . . . . .	277
10.13	Bibliography notes . . . . .	277
<b>A</b>	<b>Haskell Support Library in Haskell</b>	<b>279</b>
A.0.1	Set.hs . . . . .	279
<b>B</b>	<b>Solutions to Selected Exercises</b>	<b>281</b>

# List of Exercises

Exercise 2.1 . . . . .	17
Exercise 2.2 . . . . .	19
Exercise 2.3 . . . . .	26
Exercise 2.4 . . . . .	26
Exercise 2.5 . . . . .	30
Exercise 2.6 . . . . .	30
Exercise 2.7 . . . . .	32
Exercise 2.8 . . . . .	32
Exercise 2.9 . . . . .	32
Exercise 2.10 . . . . .	34
Exercise 2.11 . . . . .	34
Exercise 2.12 . . . . .	36
Exercise 2.13 . . . . .	37
Exercise 2.14 . . . . .	37
Exercise 2.15 . . . . .	38
Exercise 2.16 . . . . .	42
Exercise 2.17 . . . . .	45
Exercise 2.18 . . . . .	46
Exercise 2.19 . . . . .	48
Exercise 2.20 . . . . .	51
Exercise 2.21 . . . . .	52
Exercise 2.22 . . . . .	52
Exercise 2.23 . . . . .	52
Exercise 2.24 . . . . .	52
Exercise 2.25 . . . . .	53
Exercise 3.1 . . . . .	66
Exercise 3.2 . . . . .	75
Exercise 3.3 . . . . .	76

Exercise 3.4 . . . . .	76
Exercise 3.5 . . . . .	81
Exercise 3.6 . . . . .	91
Exercise 3.7 . . . . .	92
Exercise 3.8 . . . . .	92
Exercise 3.9 . . . . .	93
Exercise 3.10 . . . . .	93
Exercise 3.11 . . . . .	98
Exercise 3.12 . . . . .	98
Exercise 3.13 . . . . .	99
Exercise 3.14 . . . . .	99
Exercise 3.15 . . . . .	102
Exercise 4.1 . . . . .	109
Exercise 4.2 . . . . .	109
Exercise 4.3 . . . . .	112
Exercise 4.4 . . . . .	117
Exercise 4.5 . . . . .	117
Exercise 5.1 . . . . .	132
Exercise 5.2 . . . . .	132
Exercise 5.3 . . . . .	134
Exercise 5.4 . . . . .	136
Exercise 5.5 . . . . .	138
Exercise 5.6 . . . . .	141
Exercise 5.7 . . . . .	142
Exercise 5.8 . . . . .	142
Exercise 5.9 . . . . .	142

# Preamble

This textbook in preparation has arisen from the author's research and teaching experience. Its main aim is to provide software practitioners with 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 so as to enable algebraic reasoning.

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. The terminology seems to date from the Garmisch NATO conference in 1968, from whose report [NR69] 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. However, how "scientific" do such foundations turn out to be, now that four decades have since elapsed?

Ten years ago, Richard Bird and Oege de Moore published a textbook [BdM97] in whose preface C.A.R. Hoare writes:

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

The formulæ and equations mentioned in this quotation are those of the 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. [ABH<sup>+</sup>92] and [Bac04], among many others. Unfortunately, both of these references 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 [BdM97] too difficult for the students to follow, mainly because of its 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 [BdM97] to 4th and 5th year students (today at master level).

This draft book is by and large incomplete, most chapters being still in *slide form*. Such chapters are omitted from the current print-out.

Braga, University of Minho, December 2008

José N. Oliveira

# **Part I**

## **Calculating with Functions**





## Chapter 2

# 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.

This chapter provides a light-weight introduction to the theory of functional programming. Its emphasis is on explaining how to construct new functions out of other functions using a minimal set of predefined functional combinators. This leads to a programming style which is *point free* in the sense that function descriptions dispense with variables (definition *points*).

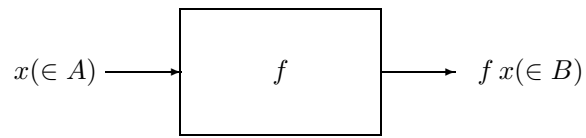
Many technical issues are deliberately ignored and deferred to later chapters. Most programming examples will be provided in the HASKELL functional programming language. Appendix A includes the listings of some HASKELL modules which complement the HUGS *Standard Prelude* (which is based very closely on the *Standard Prelude* for HASKELL 1.4.) and help to “animate” the main concepts introduced in this chapter.

## 2.1 Introducing functions and types

The definition of a function

$$f : A \longrightarrow B \tag{2.1}$$

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:



From another viewpoint,  $f$  can be regarded as a kind of “contract”: it 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,

$$\begin{aligned} \text{succ} & : \quad \mathbf{N} \longrightarrow \mathbf{N} \\ \text{succ } n & \stackrel{\text{def}}{=} n + 1 \end{aligned}$$

expresses the computation rule of the *successor* function — the function  $\text{succ}$  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  $\mathbf{N} \longrightarrow \mathbf{N}$  in the case of  $\text{succ}$  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  $\text{succ}$  and  $\text{sq}$  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 well-formedness.

It has become standard to use arrows to denote function signatures or function types, recall (2.1). In this book the following variants will be used interchangeably to denote the fact that function  $f$  accepts arguments of type  $A$  and produces results of type  $B$ :  $f :$

$B \longleftarrow A, f : A \longrightarrow B, B \xleftarrow{f} A$  or  $A \xrightarrow{f} B$ . This corresponds to writing `f :: a -> 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 which denote sets of values, very 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 “ $\cdot$ ” 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.

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 drawing (2.4)?

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  $12.cu.\bar{p}.18.ce.\bar{p}.27.co.\bar{p}.17$  for what is nowadays written as  $12x^3 + 18x^2 + 27x + 17$ . We may find such *syncopated* notation odd, but should not forget that at its time it was replacing even more obscure 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 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  $\mathbb{R} \longleftarrow \mathbb{R}$  functions are the same function. Formally, we will say that two functions  $f, g : B \longleftarrow A$  are equal if they agree at pointwise-level, that is

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

where  $=_B$  denotes equality at  $B$ -level.

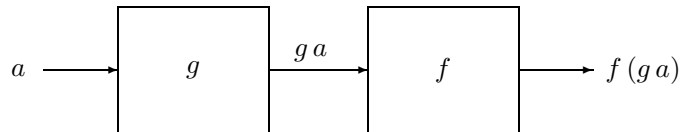
Concerning application, the pointfree style replaces it by the more generic concept of functional *composition* suggested by function-arrow chaining: wherever two functions are such that the target type of one of them, say  $B \xleftarrow{g} A$  is the same as the source type of the other, say  $C \xleftarrow{f} B$ , 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:

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

This situation is pictured by the following arrow-diagram

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

or by block-diagram



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

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

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)<sup>1</sup>. This is because of one of its most relevant properties,

$$(f \cdot g) \cdot h = f \cdot (g \cdot h) \tag{2.8}$$

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, 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  $f_1 \cdots f_n$ .

<sup>1</sup>It even has a place in script languages such as UNIX’s, 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  $A$ :

$$\begin{aligned} id_A & : A \longleftarrow A \\ id_A a & \stackrel{\text{def}}{=} a \end{aligned} \tag{2.9}$$

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  $id_{\mathbb{N}}$ , *etc.*. 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,

$$f \cdot id = id \cdot f = f \tag{2.10}$$

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* of composition. In a diagram, (2.10) looks like this:

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

Note the graphical analogy of diagrams (2.4) and (2.11). Diagrams of this kind are very common and express important 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  $A$ :

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

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):

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

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.** *The HUGS Standard Prelude provides for constant functions: you write `const c` for  $\underline{c}$ . Check that HUGS assigns the same type to expressions  $\underline{c} \cdot \text{const } c$  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 the 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

$$B \xleftarrow{f} A \begin{array}{l} \xleftarrow{h} \\ \xleftarrow{k} \end{array} C$$

( $f$  is “cancellable on the left”).

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

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

( $f$  is “cancellable on the right”).

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 2.2.** 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  $f^\circ$  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,

Weekday =

$\{Sunday, Monday, Tuesday, Wednesday, Thursday, Friday, Saturday\}$

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,

$$\begin{aligned} f : \text{Weekday} &\longrightarrow 7 \\ f \text{ Monday} &= 1 \\ f \text{ Tuesday} &= 2 \\ f \text{ Wednesday} &= 3 \\ f \text{ Thursday} &= 4 \\ f \text{ Friday} &= 5 \\ f \text{ Saturday} &= 6 \\ f \text{ Sunday} &= 7 \end{aligned}$$

is iso (guess  $f^\circ$ ). 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

$$\text{Weekday} \cong 7 \tag{2.17}$$

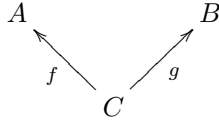
Finally, note that all classes of functions referred to so far — constants, identities, epics, monics and isos — are closed under composition, that is, the composition of two constants is a constant, the composition of two epics is epic, *etc.*

## 2.8 Gluing functions which do not compose — products

Function composition has been presented above as the 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  $f : A \longleftarrow C$  and  $g : B \longleftarrow C$  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

## 2.8. GLUING FUNCTIONS WHICH DO NOT COMPOSE — PRODUCTS 21

about gluing  $f$  and  $g$  together is to draw a diagram expressing this fact, something like



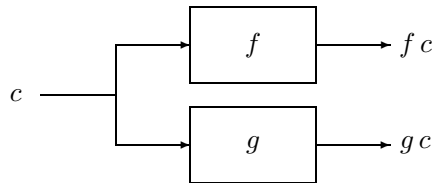
Because  $f$  and  $g$  share the same domain, their outputs can be paired, that is, we may write ordered pair  $(f\ c, g\ c)$  for each  $c \in C$ . Such pairs belong to the Cartesian product of  $A$  and  $B$ , that is, to the set

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

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:

$$\begin{aligned} \langle f, g \rangle & : C \longrightarrow A \times B \\ \langle f, g \rangle c & \stackrel{\text{def}}{=} (f\ c, g\ c) \end{aligned} \quad (2.18)$$

Function 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*  $\pi_1$  or  $\pi_2$ ,

$$A \xleftarrow{\pi_1} A \times B \xrightarrow{\pi_2} B \quad (2.19)$$

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:

$$\pi_1 \cdot \langle f, g \rangle = f \quad \text{and} \quad \pi_2 \cdot \langle f, g \rangle = g \quad (2.20)$$

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

$$\begin{array}{ccccc}
 A & \xleftarrow{\pi_1} & A \times B & \xrightarrow{\pi_2} & B \\
 & \searrow f & \uparrow \langle f, g \rangle & \nearrow g & \\
 & & C & & 
 \end{array} \tag{2.21}$$

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. But a way to “approach” the domains of  $f$  and  $g$ ,  $C$  and  $D$  respectively, is to regard them as targets of the projections  $\pi_1$  and  $\pi_2$  of  $C \times D$ :

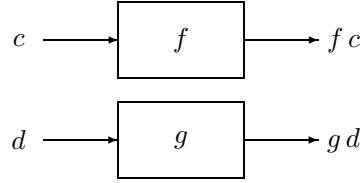
$$\begin{array}{ccccc}
 A & \xleftarrow{\pi_1} & A \times B & \xrightarrow{\pi_2} & B \\
 \uparrow f & & & & \uparrow g \\
 C & \xleftarrow{\pi_1} & C \times D & \xrightarrow{\pi_2} & D
 \end{array}$$

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

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

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:

2.8. GLUING FUNCTIONS WHICH DO NOT COMPOSE — PRODUCTS 23



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

$$f \times g \stackrel{\text{def}}{=} \langle f \cdot \pi_1, g \cdot \pi_2 \rangle \tag{2.22}$$

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

$$\frac{A \xleftarrow{f} C \quad B \xleftarrow{g} D}{A \times B \xleftarrow{f \times g} C \times D} \tag{2.23}$$

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:

$$\begin{array}{c}
 A \xleftarrow{\pi_1} A \times B \xrightarrow{\pi_2} B \\
 \swarrow \scriptstyle g \quad \searrow \scriptstyle h \\
 \quad \scriptstyle \langle g, h \rangle \\
 \quad \quad \quad \uparrow \\
 \quad \quad \quad C \\
 \swarrow \scriptstyle g \cdot f \quad \searrow \scriptstyle h \cdot f \\
 \quad \quad \quad \uparrow \scriptstyle f \\
 \quad \quad \quad D
 \end{array}
 \quad \langle g, h \rangle \cdot f = \langle g \cdot f, h \cdot f \rangle \tag{2.24}$$

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  $f = i \times j$ :

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

$$\begin{aligned}
 & \langle (i \cdot \pi_1) \cdot \langle g, h \rangle, (j \cdot \pi_2) \cdot \langle g, h \rangle \rangle \\
 = & \quad \{ \text{by (2.8)} \} \\
 & \langle i \cdot (\pi_1 \cdot \langle g, h \rangle), j \cdot (\pi_2 \cdot \langle g, h \rangle) \rangle \\
 = & \quad \{ \text{by } \times\text{-cancellation (2.20)} \} \\
 & \langle i \cdot g, j \cdot h \rangle
 \end{aligned}$$

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:

$$\begin{array}{ccc}
 A & \xleftarrow{\pi_1} & A \times B & \xrightarrow{\pi_2} & B \\
 \uparrow i & & \uparrow i \times j & & \uparrow j \\
 D & \xleftarrow{\pi_1} & D \times E & \xrightarrow{\pi_2} & E \\
 \swarrow g & & \uparrow \langle g, h \rangle & & \searrow h \\
 & & C & & 
 \end{array} \quad (i \times j) \cdot \langle g, h \rangle = \langle i \cdot g, j \cdot h \rangle \quad (2.25)$$

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

$$i \cdot \pi_1 = \pi_1 \cdot (i \times j) \quad (2.26)$$

$$j \cdot \pi_2 = \pi_2 \cdot (i \times j) \quad (2.27)$$

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:

$$(g \cdot h) \times (i \cdot j) = (g \times i) \cdot (h \times j) \quad (2.28)$$

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:

$$id_A \times id_B = id_{A \times B} \quad (2.29)$$

These two properties will be identified as the *functorial properties* of product. This choice of terminology will be explained later on.

Let us finally analyse the particular situation in which a *split* is built involving projections  $\pi_1$  and  $\pi_2$  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:

$$\begin{array}{ccc}
 A & \xleftarrow{\pi_1} & A \times B & \xrightarrow{\pi_2} & B \\
 \swarrow \pi_1 & & \uparrow id_{A \times B} & & \searrow \pi_2 \\
 & & A \times B & & 
 \end{array} \quad \langle \pi_1, \pi_2 \rangle = id_{A \times B} \quad (2.30)$$

## 2.8. GLUING FUNCTIONS WHICH DO NOT COMPOSE — PRODUCTS 25

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

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

which looks very much the same if submitted to a  $180^\circ$  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, as can be checked easily as follows:

$$\begin{aligned}
 & \textit{swap} \cdot \textit{swap} \\
 = & \quad \{ \text{by definition } \textit{swap} \stackrel{\text{def}}{=} \langle \pi_2, \pi_1 \rangle \} \\
 & \langle \pi_2, \pi_1 \rangle \cdot \textit{swap} \\
 = & \quad \{ \text{by } \times\text{-fusion (2.24)} \} \\
 & \langle \pi_2 \cdot \textit{swap}, \pi_1 \cdot \textit{swap} \rangle \\
 = & \quad \{ \text{definition of } \textit{swap} \text{ twice} \} \\
 & \langle \pi_2 \cdot \langle \pi_2, \pi_1 \rangle, \pi_1 \cdot \langle \pi_2, \pi_1 \rangle \rangle \\
 = & \quad \{ \text{by } \times\text{-cancellation (2.20)} \} \\
 & \langle \pi_1, \pi_2 \rangle \\
 = & \quad \{ \text{by } \times\text{-reflexion (2.30)} \} \\
 & \textit{id}
 \end{aligned}$$

Therefore, *swap* is iso and establishes the following isomorphism

$$A \times B \cong B \times A \tag{2.31}$$

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 predefined datatypes, `fst` to denote  $\pi_1$  and `snd` to denote  $\pi_2$ . 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 (2.31) 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

<pre> record   F: A;   S: record     F: B;     S: C;   end end; </pre>	is abstractly the same as	<pre> record   F: record     F: A;     S: B;   end;   S: C; end; </pre>
--	---------------------------	---

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

$$A \times (B \times C) \cong (A \times B) \times C \quad (2.32)$$

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

$$\text{assocr} \stackrel{\text{def}}{=} \langle \pi_1 \cdot \pi_1, \langle \pi_2 \cdot \pi_1, \pi_2 \rangle \rangle \quad (2.33)$$

Section A.0.1 in the appendix lists an extension to the HUGS *Standard Prelude*, called `Set.hs`, which makes isomorphisms such as *swap* and *assocr* 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 2.3.** Show that *assocr* is iso by conjecturing its inverse *assocl* and proving that functional equality  $\text{assocr} \cdot \text{assocl} = \text{id}$  holds.

□

---

**Exercise 2.4.** Use (2.22) to prove properties (2.28) and (2.29).

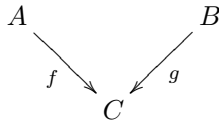
□

---



## 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  $f : C \longleftarrow A$  and  $g : C \longleftarrow B$  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*:

$$A \xrightarrow{i_1} A + B \xleftarrow{i_2} B$$

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. The tagging functions  $i_1$  and  $i_2$  are called *injections*:

$$i_1 a = (t_1, a) \quad , \quad i_2 b = (t_2, b) \tag{2.34}$$

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 and consistently. 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 HUGS *Standard Prelude*, which is based very closely on the *Standard Prelude* for HASKELL 1.4., 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:

$$\begin{aligned}
 [f, g] & : A + B \longrightarrow C \\
 [f, g] x & \stackrel{\text{def}}{=} \begin{cases} x = i_1 a \Rightarrow f a \\ x = i_2 b \Rightarrow g b \end{cases}
 \end{aligned} \tag{2.35}$$

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

$$\begin{array}{ccccc}
 A & \xrightarrow{i_1} & A + B & \xleftarrow{i_2} & B \\
 & \searrow f & \downarrow [f, g] & \swarrow g & \\
 & & C & & 
 \end{array} \tag{2.36}$$

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 a 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  $f \times g$ :

$$f + g \stackrel{\text{def}}{=} [i_1 \cdot f, i_2 \cdot g] \tag{2.37}$$

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

2.9. GLUING FUNCTIONS WHICH DO NOT COMPOSE — COPRODUCTS29

**+cancellation :**

$$\begin{array}{ccc}
 A & \xrightarrow{i_1} & A + B & \xleftarrow{i_2} & B \\
 & \searrow g & \downarrow [g,h] & \swarrow h & \\
 & & C & & 
 \end{array}
 \quad [g, h] \cdot i_1 = g, [g, h] \cdot i_2 = h \quad (2.38)$$

**+reflexion :**

$$\begin{array}{ccc}
 A & \xrightarrow{i_1} & A + B & \xleftarrow{i_2} & B \\
 & \searrow i_1 & \downarrow id_{A+B} & \swarrow i_2 & \\
 & & A + B & & 
 \end{array}
 \quad [i_1, i_2] = id_{A+B} \quad (2.39)$$

**+fusion :**

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

**+absorption :**

$$\begin{array}{ccc}
 A & \xrightarrow{i_1} & A + B & \xleftarrow{i_2} & B \\
 \downarrow i & & \downarrow i+j & & \downarrow j \\
 D & \xrightarrow{i_1} & D + E & \xleftarrow{i_2} & E \\
 & \searrow g & \downarrow [g,h] & \swarrow h & \\
 & & C & & 
 \end{array}
 \quad [g, h] \cdot (i + j) = [g \cdot i, h \cdot j] \quad (2.41)$$

**+functor :**

$$(g \cdot h) + (i \cdot j) = (g + i) \cdot (h + j) \quad (2.42)$$

**+functor-id :**

$$id_A + id_B = id_{A+B} \quad (2.43)$$

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

$$\frac{\begin{array}{c} C \xleftarrow{f} A \\ C \xleftarrow{g} B \end{array}}{C \xleftarrow{[f,g]} A + B} \quad \frac{\begin{array}{c} C \xleftarrow{f} A \\ D \xleftarrow{g} B \end{array}}{C + D \xleftarrow{f+g} A + B} \quad (2.44)$$

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

$$A + B \cong B + A \quad (2.45)$$

holds.

□

---

**Exercise 2.6.** *Dualize (2.33), that is, write the iso which witnesses fact*

$$A + (B + C) \cong (A + B) + C \quad (2.46)$$

*from right to left. Use the `either` syntax available from the HUGS 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:

$$[\langle f, g \rangle, \langle h, k \rangle] = \langle [f, h], [g, k] \rangle \quad (2.47)$$

It can be checked by type-inference that both the left-hand side and the right-hand side expressions of this equality have type  $B \times D \longleftarrow A + C$ , 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

$$A \times (B + C) \xleftarrow{\text{undistr}} (A \times B) + (A \times C) \quad (2.48)$$

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

$$\text{undistr} \stackrel{\text{def}}{=} [id \times i_1, id \times i_2] \quad (2.49)$$

and witnesses the fact that product distributes through coproduct:

$$A \times (B + C) \cong (A \times B) + (A \times C) \quad (2.50)$$

In this context, suppose that we know of three functions  $D \xleftarrow{f} A$ ,  $E \xleftarrow{g} B$  and  $F \xleftarrow{h} C$ . By (2.44) we infer  $E + F \xleftarrow{g+h} B + C$ . Then, by (2.23) we infer

$$D \times (E + F) \xleftarrow{f \times (g+h)} A \times (B + C) \quad (2.51)$$

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 (2.51).

This kind of abstraction works for every combination of products and coproducts. For instance, if we now abstract the right-hand side of (2.48) 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

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

which unfolds to

$$\begin{array}{ccc}
 A \times (B + C) & \xleftarrow{\text{undistr}} & (A \times B) + (A \times C) & (2.52) \\
 f \times (g+h) \downarrow & & \downarrow (f \times g) + (f \times h) & \\
 D \times (E + F) & & (D \times E) + (D \times F) &
 \end{array}$$

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

$$\begin{array}{ccc}
 & & ? \\
 [x,y] & \swarrow & \downarrow id+id \times f \\
 ? & \xleftarrow{\quad} & ? \\
 & [k,g] &
 \end{array}$$

and then solve the implicit equation for *x* and *y*.

□

---

**Exercise 2.9.** Repeat exercise 2.8 with respect to diagram

$$\begin{array}{ccc}
 ? & \xrightarrow{h+\langle i,j \rangle} & ? \\
 & \searrow & \downarrow id+id \times f \\
 x+y & & ?
 \end{array}$$

□

---

## 2.11 Natural properties

Let us resume discussion about *undistr* and the two other functions in diagram (2.52). What about using *undistr* itself to close this diagram, at the bottom? Note that definition (2.49) 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 (2.48) was rather arbitrary.) Therefore, we get:

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

which expresses a very important property of *undistr*:

$$(f \times (g + h)) \cdot \text{undistr} = \text{undistr} \cdot ((f \times g) + (f \times h)) \quad (2.53)$$

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:

$$\begin{array}{ccc} A & F A \xleftarrow{\phi} G A & (F f) \cdot \phi = \phi \cdot (G f) \\ f \downarrow & F f \downarrow \quad \downarrow G f & \\ B & F B \xleftarrow{\phi} G B & \end{array} \quad (2.54)$$

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

Diagram (2.54) 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, (2.28) or (2.42).

**Exercise 2.10.** Show that (2.26) and (2.27) are natural properties. Dualize these properties. *Hint:* recall diagram (2.41).

□

**Exercise 2.11.** Establish the natural properties of the swap (2.31) and assoc (2.33) isomorphisms.

□

## 2.12 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 come from? As a rule, for each combinator to be defined one has to define suitable constructions at “interface”-level <sup>2</sup>, 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 derived.

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

$$k = \langle f, g \rangle \Leftrightarrow \begin{cases} \pi_1 \cdot k = f \\ \pi_2 \cdot k = g \end{cases} \quad (2.55)$$

holds — recall diagram (2.21) — for all  $A \times B \xleftarrow{k} C$ . This equivalence states that  $\langle f, g \rangle$  is the *unique* arrow satisfying the property on the right. In fact, read (2.55) 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

<sup>2</sup>In the current context, *programs* “are” functions and *program-interfaces* “are” the datatypes involved in functional signatures.



$\times$ -cancellation (2.20). Reading (2.55) in the  $\Leftarrow$  direction we understand that, if some  $k$  satisfies such properties, then it “has to be” the same arrow as  $\langle f, g \rangle$ .

It is easy to see other properties of  $\langle f, g \rangle$  arising from (2.55). For instance, for  $k = id$  we get  $\times$ -reflexion (2.30),

$$\begin{aligned}
 id = \langle f, g \rangle &\Leftrightarrow \begin{cases} \pi_1 \cdot id = f \\ \pi_2 \cdot id = g \end{cases} \\
 \equiv &\quad \{ \text{by (2.10)} \} \\
 id = \langle f, g \rangle &\Leftrightarrow \begin{cases} \pi_1 = f \\ \pi_2 = g \end{cases} \\
 \equiv &\quad \{ \text{by substitution of } f \text{ and } g \} \\
 id &= \langle \pi_1, \pi_2 \rangle
 \end{aligned}$$

and for  $k = \langle i, j \rangle \cdot h$  we get  $\times$ -fusion (2.24):

$$\begin{aligned}
 \langle i, j \rangle \cdot h = \langle f, g \rangle &\Leftrightarrow \begin{cases} \pi_1 \cdot (\langle i, j \rangle \cdot h) = f \\ \pi_2 \cdot (\langle i, j \rangle \cdot h) = g \end{cases} \\
 \equiv &\quad \{ \text{composition is associative (2.8)} \} \\
 \langle i, j \rangle \cdot h = \langle f, g \rangle &\Leftrightarrow \begin{cases} (\pi_1 \cdot \langle i, j \rangle) \cdot h = f \\ (\pi_2 \cdot \langle i, j \rangle) \cdot h = g \end{cases} \\
 \equiv &\quad \{ \text{by } \times\text{-cancellation (just derived)} \} \\
 \langle i, j \rangle \cdot h = \langle f, g \rangle &\Leftrightarrow \begin{cases} i \cdot h = f \\ j \cdot h = g \end{cases} \\
 \equiv &\quad \{ \text{by substitution of } f \text{ and } g \} \\
 \langle i, j \rangle \cdot h &= \langle i \cdot h, j \cdot h \rangle
 \end{aligned}$$

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

$$\langle i, j \rangle = \langle f, g \rangle \Leftrightarrow \begin{cases} i = f \\ j = g \end{cases} \tag{2.56}$$

from universal property (2.55) — just let  $k = \langle i, j \rangle$ .

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

$$k = [ f, g ] \Leftrightarrow \begin{cases} k \cdot i_1 = f \\ k \cdot i_2 = g \end{cases} \tag{2.57}$$

from which structural equality of *eithers* can be inferred,

$$[i, j] = [f, g] \Leftrightarrow \begin{cases} i = f \\ j = g \end{cases} \quad (2.58)$$

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

**Exercise 2.12.** *Derive +-cancellation (2.38), +-reflexion (2.39) and +-fusion (2.40) from universal property (2.57). Then derive the exchange law (2.47) from the universal property of product (2.55) or coproduct (2.57).*

□

---

## 2.13 Guards and McCarthy's conditional

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

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

meaning

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

for some given predicate  $\text{Bool} \xleftarrow{p} A$ , some “then”-function  $B \xleftarrow{g} A$  and some “else”-function  $B \xleftarrow{h} A$ .  $\text{Bool}$  is the primitive datatype containing truth values  $\text{FALSE}$  and  $\text{TRUE}$ .

Can such expressions be written in the pointfree style? They can, provided we introduce the so-called “McCarthy conditional” functional form

$$p \rightarrow g, h$$

which is defined by

$$p \rightarrow g, h \stackrel{\text{def}}{=} [g, h] \cdot p? \quad (2.59)$$

a definition we can understand provided we know the meaning of the “ $p?$ ” construct.

We call  $A + A \xleftarrow{p?} A$  a *guard*, or better, the guard associated to a given predicate

$\text{Bool} \xleftarrow{p} A$ . Every predicate  $p$  gives birth to its own guard  $p?$  which, at point-level, is defined as follows:

$$(p?)a = \begin{cases} pa & \Rightarrow i_1 a \\ \neg(pa) & \Rightarrow i_2 a \end{cases} \quad (2.60)$$

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.

The following fact, which we will refer to as *McCarthy's conditional fusion law*, is a consequence of  $+$ -fusion (2.40):

$$f \cdot (p \rightarrow g, h) = p \rightarrow f \cdot g, f \cdot h \quad (2.61)$$

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  $\text{Bool} \xleftarrow{=} A \times A$  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  $\text{Eq}$ , which exports equality predicate  $(=)$ . This does not, however, guarantee the reflexive, transitive and symmetry properties, which need to be proved by dedicated mathematical arguments.

**Exercise 2.13.** *Prove that the following equality between two conditional expressions*

$$\begin{aligned} & k(\text{if } px \text{ then } fx \text{ else } hx, \text{if } px \text{ then } gx \text{ else } ix) \\ &= \text{if } px \text{ then } k(fx, gx) \text{ else } k(hx, ix) \end{aligned}$$

*holds by rewriting it in the pointfree style (using the McCarthy's conditional combinator) and applying the exchange law (2.47), among others.*

□

---

**Exercise 2.14.** *Prove law (2.61).*

□

---

**Exercise 2.15.** From (2.59) and property

$$p? \cdot f = (f + f) \cdot (p \cdot f)? \quad (2.62)$$

infer

$$(p \rightarrow f, g) \cdot h = (p \cdot h) \rightarrow (f \cdot h), (g \cdot h) \quad (2.63)$$

□

---

## 2.14 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  $B \xleftarrow{f} C \times A$  and  $D \xleftarrow{g} A$ . 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

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

from  $f$  which combines nicely with  $g$  via the *split* combinator:  $\langle f_c, g \rangle$  is well-defined and bears type  $B \times D \longleftarrow A$ . 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:

$$=_{c\ a} \stackrel{\text{def}}{=} a = c \quad (2.64)$$

## 2.14. GLUING FUNCTIONS WHICH DO NOT COMPOSE — EXPONENTIALS 39

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  $B \longleftarrow A$ . We will say that  $f_c \in B^A$  by introducing a new datatype construct which we will refer to as the *exponential*:

$$B^A \stackrel{\text{def}}{=} \{g \mid g : B \longleftarrow A\} \quad (2.65)$$

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  $g : B \longleftarrow A$  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

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

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  $B^A \longleftarrow C$  which expresses  $f$  as a kind of  $C$ -indexed family of functions of signature  $B \longleftarrow A$ . 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:

$$f(c, a) = (\bar{f} c) a \quad (2.66)$$

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,

$$k = \bar{f} \Leftrightarrow f = ap \cdot (k \times id) \quad (2.67)$$

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

**Exponentials cancellation :**

$$\begin{array}{ccc}
 B^A & B^A \times A \xrightarrow{ap} & B \\
 \bar{f} \uparrow & \bar{f} \times id \uparrow & \nearrow f \\
 C & C \times A &
 \end{array}
 \qquad
 f = ap \cdot (\bar{f} \times id)
 \qquad
 (2.68)$$

**Exponentials reflexion :**

$$\begin{array}{ccc}
 B^A & B^A \times A \xrightarrow{ap} & B \\
 id_{B^A} \uparrow & id_{B^A} \times id_A \uparrow & \nearrow ap \\
 B^A & B^A \times A &
 \end{array}
 \qquad
 \overline{ap} = id_{B^A}
 \qquad
 (2.69)$$

**Exponentials fusion :**

$$\begin{array}{ccc}
 B^A & B^A \times A \xrightarrow{ap} & B \\
 \bar{g} \uparrow & \bar{g} \times id \uparrow & \nearrow g \\
 C & C \times A & \nearrow g \cdot (f \times id) \\
 f \uparrow & f \times id \uparrow & \\
 D & D \times A &
 \end{array}
 \qquad
 \overline{g \cdot (f \times id)} = \bar{g} \cdot f
 \qquad
 (2.70)$$

Note that the cancellation law is nothing but fact (2.66) 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:

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

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

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

2.14. GLUING FUNCTIONS WHICH DO NOT COMPOSE — EXPONENTIALS 41

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 2.23). So  $f^A$  is the “compose with  $f$ ” functional combinator:

$$(f^A)g \stackrel{\text{def}}{=} f \cdot g \tag{2.71}$$

Now we are ready to understand the laws which follow:

**Exponentials absorption :**

$$\begin{array}{ccc}
 D^A & D^A \times A \xrightarrow{ap} & D \\
 \uparrow f^A & \uparrow f^A \times id & \uparrow f \\
 B^A & B^A \times A \xrightarrow{ap} & B \\
 \uparrow \bar{g} & \uparrow \bar{g} \times id & \nearrow g \\
 C & C \times A & 
 \end{array}
 \qquad
 \overline{f \cdot g} = f^A \cdot \bar{g} \tag{2.72}$$

**Exponentials-functor :**

$$(g \cdot h)^A = g^A \cdot h^A \tag{2.73}$$

**Exponentials-functor-id :**

$$id^A = id \tag{2.74}$$

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. Let us introduce the following operator

$$curry f \stackrel{\text{def}}{=} \bar{f} \tag{2.75}$$

which maps a function  $f$  to its transpose  $\bar{f}$ . This operator, which is very familiar to functional programmers, maps functions in some function space  $B^{C \times A}$  to functions in  $(B^A)^C$ . Its inverse (known as the  $\hat{\phantom{f}}$ -function) also exists. In the HUGS *Standard Prelude* we find them declared as follows:

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

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

From (2.75) it is obvious see that writing  $\overline{f}$  or *curry*  $f$  is a matter of taste, the latter being more in the tradition of functional programming. For instance, the fusion law (2.70) can be re-written as

$$\text{curry } (g \cdot (f \times \text{id})) = \text{curry } g \cdot f$$

and so on.

It is known from mathematics that *curry* and  $\hat{\phantom{x}}$  are isos witnessing the following isomorphism which is at the core of the theory of functional programming:

$$B^{C \times A} \cong (B^A)^C \quad (2.76)$$

Fact (2.76) clearly resembles a well known equality concerning numeric exponentials,  $b^{c \times a} = (b^a)^c$ . But other known facts about numeric exponentials, e.g.  $a^{b+c} = a^b \times a^c$  or  $(b \times c)^a = b^a \times c^a$  find their counterpart in functional exponentials. The counterpart of the former,

$$A^{B+C} \cong A^B \times A^C \quad (2.77)$$

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

$$(B \times C)^A \cong B^A \times C^A \quad (2.78)$$

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  $1^a = 1$ ? We need to know what 0 and 1 mean as datatypes. Such elementary datatypes are presented in the section which follows.

**Exercise 2.16.** Load module *Set.hs* (cf. section A.0.1) into the HUGS interpreter and check the types assigned to the following functional expressions:

```
curry ap
\f -> ap . ( f >< id)
uncurry . curry
```

Which of these is functionally equivalent to the *uncurry* function and why? Which of these are functionally equivalent to identity functions? Justify.

□

---



## 2.15 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  $\mathbf{R}$ ,  $\mathbf{Bool}$  and  $\mathbf{N}$  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  $\mathbf{N}_0$  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:

$$A + 0 \cong A \quad (2.79)$$

$$A \times 0 \cong 0 \quad (2.80)$$

### 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 relevant facts involving  $1$ :

$$A \times 1 \cong A \quad (2.81)$$

$$A^0 \cong 1 \quad (2.82)$$

We can read (2.81) informally as follows: if the second component of a record (“struct”) cannot change, then it is useless and can be ignored. Selector  $\pi_1$  is, in this context, an iso mapping the left-hand side of (2.81) to its right-hand side. Its inverse is  $\langle id, c \rangle$  where  $c$  is a particular choice of inhabitant for datatype  $1$ . Concerning (2.82),  $A^0$  denotes the set of all functions from the empty set to some  $A$ . What does (2.82) mean? It simply tells

us that there is only one 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

$$1^A \cong 1 \tag{2.83}$$

$$A^1 \cong A \tag{2.84}$$

which are the functional exponentiation counterparts of  $1^a = 1$  and  $a^1 = a$ . Fact (2.83) 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. Fact (2.84) 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  $A$ .

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.

Finally, what can we say about  $1 + A$ ? Every function  $B \xleftarrow{f} 1 + A$  observing this type is bound to be an *either*  $[\underline{b_0}, g]$  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 HUGS *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:  $f = [\underline{b_1}, \underline{b_2}]$ . 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  $\text{Bool}$  is expected.

**Exercise 2.17.** *Relate HASKELL expressions*

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

and

```
\f->(f True, f False)
```

to the following isomorphisms involving generic elementary type  $2$ :

$$2 \times A \cong A + A \quad (2.85)$$

$$A \times A \cong A^2 \quad (2.86)$$

Apply the exchange law (2.47) to the first expression above.

□

## 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 2.1.

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  $\pi_i$  are associated as the number  $n$  of factors involved. Of course, *splits* become  $n$ -ary as well

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

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

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, *eithers* become  $n$ -ary

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

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

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
$B^A$	array[A] of B	B ...[A]	Arrays
$1 + A$	$\hat{A}$	A *...	Pointers

Table 2.1: Abstract notation versus programming language data-structures.

### 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 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

$$\underbrace{A \times \dots \times A}_n \cong A^n \quad (2.87)$$

$$\underbrace{A + \dots + A}_n \cong n \times A \quad (2.88)$$

hold.

**Exercise 2.18.** *On the basis of table 2.1, encode  $\text{undistr}$  (2.49) 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 *eithers* 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 (2.21) and (2.36) 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 (2.83). 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 (2.82).

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 :

$$\begin{array}{ccc} ?_0 = id_0 & & \\ \curvearrowright & & \\ 0 & & \end{array} \qquad ?_0 = id_0 \qquad (2.89)$$

### Initial datatype fusion :

$$\begin{array}{ccc} 0 & & \\ ?_A \downarrow & \searrow ?_B & \\ A & \xrightarrow{f} & B \end{array} \qquad f \cdot ?_A = ?_B \qquad (2.90)$$

### Terminal datatype reflexion :

$$\begin{array}{ccc} !_1 = id_1 & & \\ \curvearrowleft & & \\ 1 & & \end{array} \qquad !_1 = id_1 \qquad (2.91)$$

**Terminal datatype fusion :**

$$\begin{array}{ccc}
 & 1 & \\
 & \swarrow \!_B & \\
 A & \xleftarrow{f} & B
 \end{array}
 \qquad
 !_A \cdot f = !_B
 \qquad
 (2.92)$$

**Exercise 2.19.** Particularize the exchange law (2.47) 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 Int | CC Int
```

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 HUGS *interpreter* about what it knows about `CostumerId`, the reply will contain the following information:

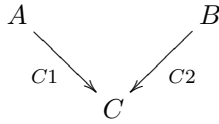
```
Main> :i CostumerId
-- type constructor
data CostumerId

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

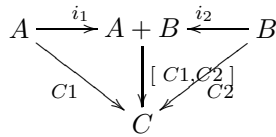
In general, let  $A$  and  $B$  be two known datatypes. Via declaration

$$\text{data } C = C1 A \mid C2 B
 \qquad
 (2.93)$$

one obtains from HUGS a new datatype  $C$  equipped with constructors  $C \xleftarrow{C_1} A$  and  $C \xleftarrow{C_2} B$ , in fact the only ones available for constructing values of  $C$ :

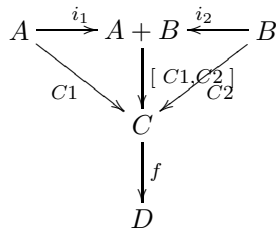


This diagram leads to an obvious instance of coproduct diagram (2.36),



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  $[C_1, C_2]$  has an inverse, say  $inv$ , which is such that  $inv \cdot [C_1, C_2] = id$ . 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  $+fusion$  (2.40), whereby we obtain

$$f \cdot [C_1, C_2] = [f \cdot C_1, f \cdot C_2]$$

Therefore, the observation will be fully described provided we explain how  $f$  behaves with respect to  $C_1$  — *cf.*  $f \cdot C_1$  — and with respect to  $C_2$  — *cf.*  $f \cdot C_2$ . This is what is behind the typical *inductive* structure of pointwise  $f$ , which will be made of two and only two clauses:

$$\begin{aligned} f &: C \longrightarrow D \\ f(C_1 a) &= \dots \\ f(C_2 b) &= \dots \end{aligned}$$

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

$$\begin{aligned}
 & inv \cdot [C1, C2] = id \\
 \equiv & \quad \{ \text{by } +\textit{-fusion} (2.40) \} \\
 & [inv \cdot C1, inv \cdot C2] = id \\
 \equiv & \quad \{ \text{by } +\textit{-reflexion} (2.39) \} \\
 & [inv \cdot C1, inv \cdot C2] = [i_1, i_2] \\
 \equiv & \quad \{ \textit{either structural equality} (2.58) \} \\
 & inv \cdot C1 = i_1 \wedge inv \cdot C2 = i_2
 \end{aligned}$$

Therefore:

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

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

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

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

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

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  $\pi_1, \pi_2$  “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  $y$ ):

$$\begin{array}{ccccc}
 A & \xleftarrow{\pi_1} & A \times A & \xrightarrow{\pi_2} & A \\
 & \searrow x & \downarrow \textit{Point} & \nearrow y & \\
 & & \textit{Point} & & 
 \end{array} \tag{2.94}$$



This rises an algebra ( $Point$ ) and a coalgebra ( $\langle x, y \rangle$ ) for datatype  $Point$ :

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

In HASKELL one writes

```
data Point a = Point { x :: a, y :: a }
```

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

```
Point :: a -> a -> Point a
```

Finally, what is the “pointer”-equivalent in HASKELL? This corresponds to  $A = 1$  in (2.93) and to the following HASKELL declaration:

```
data C = C1 () | 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 | 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 (2.84) helps in comparing these two alternative situations.

## 2.19 Exercises

**Exercise 2.20.** Let  $A$  and  $B$  be two disjoint datatypes, that is,  $A \cap B = \emptyset$  holds. Show that isomorphism

$$A \cup B \cong A + B \tag{2.95}$$

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

□

---

**Exercise 2.21.** 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  $(B + C) \times A \cong B \times A + C \times A$ :

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

$\xrightarrow{\quad\quad\quad} \quad\quad\quad \xrightarrow{distl}$

□

---

**Exercise 2.22.** In the context of exercise 2.21, prove

$$[g, h] \times f = [g \times f, h \times f] \cdot distr \tag{2.96}$$

knowing that

$$f \times [g, h] = [f \times g, f \times h] \cdot distr$$

holds.

□

---

**Exercise 2.23.** Show that  $\overline{(f \cdot ap)} g = f \cdot g$  holds, cf. (2.71).

□

---

**Exercise 2.24.** Let  $C \xrightarrow{const} C^A$  be the function of exercise 2.1, that is,  $const\ c = \underline{c}_A$ . Which fact is expressed by the following diagram featuring  $const$ ?

$$\begin{array}{ccc} C & \xrightarrow{const} & C^A \\ f \downarrow & & \downarrow f^A \\ B & \xrightarrow{const} & B^A \end{array}$$

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

□

---

**Exercise 2.25.** *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 [Bac78]. 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<sup>3</sup>. 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

---

<sup>3</sup>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.

[BdM97]. 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 [McC63], of Floyd's work on program meaning [Flo67] and of Paterson and Hewitt's *comparative schematology* [PH70]. Work of the so-called *program transformation* school was already very expressive in the mid 1970s, see for instance references [BD77].

The mathematics adequate for the effective integration of these related but independent lines of thought was provided by the categorial approach of Manes and Arbib compiled in a textbook [MA86] 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 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 [Bir98] 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.

# Chapter 3

## 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 2.1 — which records an analogy between abstract type notation and the corresponding data-structures available in common, imperative languages.

This analogy is precisely our point of departure for extending the abstract notation towards a most important field of programming: *recursion*.

### 3.1 Motivation

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

```
L =  $\hat{N}$ ;  
N = record  
    first: A;  
    next:  $\hat{N}$   
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 2.1. 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:

$$\begin{cases} L = 1 + N \\ N = A \times (1 + N) \end{cases} \quad (3.1)$$

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 (3.1) is thus equivalent to:

$$\begin{cases} L = 1 + A \times L \\ N = A \times (1 + N) \end{cases} \quad (3.2)$$

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,

$$x = 1 + \frac{x}{2} \quad ? \quad (3.3)$$

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

$$x = 1 + \frac{x}{2}$$

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

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  $b \neq 0$ :

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

$$\begin{aligned}
& 2 \times x - 2 \times \frac{x}{2} = 2 \\
\equiv & \quad \{ 2 \text{ cancels its inverse } \frac{1}{2} \} \\
& 2 \times x - 1 \times x = 2 \\
\equiv & \quad \{ \text{multiplication distributes over addition} \} \\
& (2 - 1) \times x = 2 \\
\equiv & \quad \{ 2 - 1 = 1 \text{ and } 1 \text{ is the unit of multiplication} \} \\
& x = 2
\end{aligned}$$

Back to (3.2), we would like to submit each of the equations, *e.g.*

$$L = 1 + A \times L \tag{3.4}$$

to a similar reasoning. Can we do it? The analogy which can be found between this equation and (3.3) goes beyond pattern similarity. From chapter 2 we know that many properties required in the reasoning above hold in the context of (3.4), provided the “=” sign is replaced by the “ $\cong$ ” sign, that of set-theoretical isomorphism. Recall that, for instance, + is associative (2.46), 0 is the unit of addition (2.79), 1 is the unit of multiplication (2.81), multiplication distributes over addition (2.50) *etc.* Moreover, the first step above assumed that addition is compatible (monotonic) with respect to equality,

$$\begin{array}{rcc}
a & = & b \\
c & = & d \\
\hline
a + c & = & b + d
\end{array}$$

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

$$\begin{array}{rcc}
A & \cong & B \\
C & \cong & D \\
\hline
A + C & \cong & B + D
\end{array}$$

— recall (2.44) 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 inverses<sup>1</sup>, and so we cannot “calculate by cancellation”. How do we circumvent this limitation?

Just think of how we would have gone about (3.3) 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...

$$\begin{aligned}
 x &= 1 + \frac{x}{2} \\
 \equiv & \quad \{ x \text{ by } 1 + \frac{x}{2} \text{ substitution followed by simplification} \} \\
 x &= 1 + \frac{1 + \frac{x}{2}}{2} = 1 + \frac{1}{2} + \frac{x}{4} \\
 \equiv & \quad \{ \text{the same as above} \} \\
 x &= 1 + \frac{1}{2} + \frac{1 + \frac{x}{2}}{4} = 1 + \frac{1}{2} + \frac{1}{4} + \frac{x}{8} \\
 \equiv & \quad \{ \text{over and over again, } n\text{-times} \} \\
 & \dots \\
 \equiv & \quad \{ \text{simplification} \} \\
 x &= \sum_{i=0}^n \frac{1}{2^i} + \frac{x}{2^{n+1}} \\
 \equiv & \quad \{ \text{sum of } n \text{ first terms of a geometric progression} \} \\
 x &= (2 - \frac{1}{2^n}) + \frac{x}{2^{n+1}} \\
 \equiv & \quad \{ \text{let } n \rightarrow \infty \} \\
 x &= (2 - 0) + 0 \\
 \equiv & \quad \{ \text{simplification} \} \\
 x &= 2
 \end{aligned}$$

Clearly, this is a much more complicated way of finding solution  $x = 2$  for equation (3.3). But we would have loved it in case it were the only known way, and this is precisely what happens with respect to (3.4). In this case we have:

$$\begin{aligned}
 L &= 1 + A \times L \\
 \equiv & \quad \{ \text{substitution of } 1 + A \times L \text{ for } L \}
 \end{aligned}$$

---

<sup>1</sup>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.

$$\begin{aligned}
& L = 1 + A \times (1 + A \times L) \\
\equiv & \quad \{ \text{distributive property (2.50)} \} \\
& L \cong 1 + A \times 1 + A \times (A \times L) \\
\equiv & \quad \{ \text{unit of product (2.81) and associativity of product (2.32)} \} \\
& L \cong 1 + A + (A \times A) \times L \\
\equiv & \quad \{ \text{by (2.82), (2.84) and (2.87)} \} \\
& L \cong A^0 + A^1 + A^2 \times L \\
\equiv & \quad \{ \text{another substitution as above and similar simplifications} \} \\
& L \cong A^0 + A^1 + A^2 + A^3 \times L \\
\equiv & \quad \{ \text{after } (n+1)\text{-many similar steps} \} \\
& L \cong \sum_{i=0}^n A^i + A^{n+1} \times L
\end{aligned}$$

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:

$$A^* \stackrel{\text{def}}{=} \bigcup_{i \geq 0} A^i \tag{3.5}$$

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 2.20), we finally have:

$$L \cong A^*$$

All in all, we have obtained  $A^*$  as a solution to equation (3.4). 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 (3.4) 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 in chapter 7, 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.

## 3.2 Introducing inductive datatypes

Datatype  $L$  as defined by (3.4) is said to be *recursive* because  $L$  “recurs” in the definition of  $L$  itself<sup>2</sup>. From the discussion above, it is clear that set-theoretical equality “=” in this equation should give place to set-theoretical isomorphism (“ $\cong$ ”):

$$L \cong 1 + A \times L \tag{3.6}$$

Which isomorphism  $L \xleftarrow{in} 1 + A \times L$  do we expect to witness (3.4)? This will depend on which particular solution to (3.4) 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 (3.4) by replacing unknown  $L$  with solution  $A^*$ , that is

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

---

<sup>2</sup>By analogy, we may regard (3.3) as a “recursive definition” of number 2.

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

$$\begin{array}{ccc}
 1 & \xrightarrow{i_1} & 1 + A \times A^* & \xleftarrow{i_2} & A \times A^* \\
 & \searrow & \downarrow in & \swarrow & \\
 & & A^* & & 
 \end{array}
 \quad (3.7)$$

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:

$$\begin{aligned}
 cons &: A \times A^* \longrightarrow A^* \\
 cons(a, [a_1, \dots, a_n]) &= [a, a_1, \dots, a_n]
 \end{aligned}
 \quad (3.8)$$

In a diagram:

$$\begin{array}{ccc}
 1 & \xrightarrow{i_1} & 1 + A \times A^* & \xleftarrow{i_2} & A \times A^* \\
 & \searrow & \downarrow [ [], cons ] & \swarrow & \\
 & & A^* & & 
 \end{array}
 \quad (3.9)$$

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

$$out \stackrel{\text{def}}{=} (! + \langle hd, tl \rangle) \cdot (= []?) \quad (3.10)$$

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

$$\begin{aligned}
 hd &: A^* \longrightarrow A \\
 hd[a_1, a_2, \dots, a_n] &= a_1
 \end{aligned}
 \quad (3.11)$$

$$\begin{aligned}
 tl &: A^* \longrightarrow A^* \\
 tl[a_1, a_2, \dots, a_n] &= [a_2, \dots, a_n]
 \end{aligned}
 \quad (3.12)$$

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

$$\begin{aligned}
 & in \cdot out = id \\
 \equiv & \quad \{ \text{definitions of } in \text{ and } out \} \\
 & [ [], cons ] \cdot (! + \langle hd, tl \rangle) \cdot (= []?) = id \\
 \equiv & \quad \{ +-absorption (2.41) \text{ and } (2.15) \}
 \end{aligned}$$

$$\begin{aligned}
& [\underline{[]}, \text{cons} \cdot \langle \text{hd}, \text{tl} \rangle] \cdot (=_{[]}?) = \text{id} \\
\equiv & \quad \{ \text{property of sequences: } \text{cons}(\text{hd } s, \text{tl } s) = s \} \\
& [\underline{[]}, \text{id}] \cdot (=_{[]}?) = \text{id} \\
\equiv & \quad \{ \text{going pointwise (2.60)} \} \\
& \left\{ \begin{array}{l} =_{[]} a \Rightarrow [\underline{[]}, \text{id}](i_1 a) \\ \neg(=_{[]} a) \Rightarrow [\underline{[]}, \text{id}](i_2 a) \end{array} \right. = a \\
\equiv & \quad \{ \text{+-cancellation (2.38)} \} \\
& \left\{ \begin{array}{l} =_{[]} a \Rightarrow \underline{[]} a \\ \neg(=_{[]} a) \Rightarrow \overline{\text{id}} a \end{array} \right. = a \\
\equiv & \quad \{ a = [] \text{ in one case and identity function (2.9) in the other} \} \\
& \left\{ \begin{array}{l} a = [] \Rightarrow a \\ \neg(a = []) \Rightarrow a \end{array} \right. = a \\
\equiv & \quad \{ \text{property } (p \rightarrow f, f) = f \text{ holds} \} \\
& a = a
\end{aligned}$$

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 (3.6)? 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) (3.13)
```

and generates diagram

$$\begin{array}{ccc}
1 & \xrightarrow{i_1} & 1 + A \times L \xleftarrow{i_2} A \times L \\
& \searrow \underline{Nil} & \downarrow in' \swarrow Cons \\
& & L
\end{array} \tag{3.14}$$

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 (3.6), 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 (3.13):

$$\text{data } T = \text{Nil} \mid \text{Cons } (A, T) \quad (3.15)$$

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 (3.15) 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

$$\begin{aligned} out_T : T &\longrightarrow 1 + A \times T \\ out_T Nil &= i_1 \text{NIL} \\ out_T(Cons(a,l)) &= i_2(a,l) \end{aligned} \quad (3.16)$$

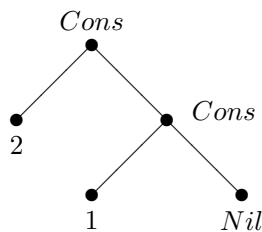
which can be straightforwardly encoded in HASKELL using the `Either` realization of  $+$  (recall sections 2.9 and 2.18):

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

In summary, isomorphism

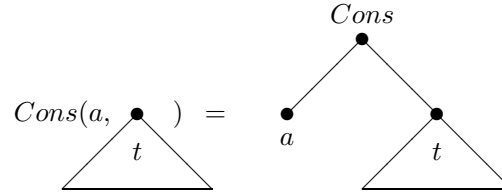
$$\begin{array}{ccc} & \xrightarrow{out_T} & \\ T & \cong & 1 + A \times T \\ & \xleftarrow{in_T} & \end{array} \quad (3.17)$$

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  $\mathbb{T}$ ,  $Nil$  and  $Cons$  was rather arbitrary in (3.15). Therefore, an alternative declaration such as, for instance,

$$\text{data } \mathbb{U} = \text{Stop} \mid \text{Join } (A, \mathbb{U}) \quad (3.18)$$

would have been perfectly acceptable, generating another solution for the equation under algebra  $[\text{Stop}, \text{Join}]$ . It is easy to check that (3.18) 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 7.

**Exercise 3.1.** *Obviously,*

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

*does not work as a HASKELL realization of the mediating arrow in diagram (3.9). What do you need to write instead?*

□

### 3.3 Observing an inductive datatype

Suppose that one is asked to express a particular *observation* of an inductive such as  $\mathbb{T}$  (3.15), that is, a function of signature  $B \xleftarrow{f} \mathbb{T}$  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  $\mathbb{N}_0$ ,

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

and that  $0 \in \mathbb{N}_0$  is a value denoting “the addition of nothing”. So constant arrow  $\mathbb{N}_0 \xleftarrow{\underline{0}} 1$  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  $\underline{0}$ ,

$$\mathbb{N}_0 \xleftarrow{[\underline{0}, add]} 1 + \mathbb{N}_0 \times \mathbb{N}_0 \quad (3.19)$$

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

$$T \xleftarrow{in_T} 1 + \mathbb{N}_0 \times T \quad (3.20)$$

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:

$$\begin{array}{ccc} T & \xleftarrow{in_T} & 1 + \mathbb{N}_0 \times T \\ f \downarrow & & \downarrow \text{---} \\ \mathbb{N}_0 & \xleftarrow{[\underline{0}, add]} & 1 + \mathbb{N}_0 \times \mathbb{N}_0 \end{array} \quad (3.21)$$

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 lefthand side, it will do the same job on the righthand 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

$$\begin{array}{ccc} T & \xleftarrow{[\underline{Nil}, Cons]} & 1 + \mathbb{N}_0 \times T \\ f \downarrow & & \downarrow id + id \times f \\ \mathbb{N}_0 & \xleftarrow{[\underline{0}, add]} & 1 + \mathbb{N}_0 \times \mathbb{N}_0 \end{array} \quad (3.22)$$

which pictures the following property of  $f$

$$f \cdot [\underline{Nil}, Cons] = [\underline{0}, add] \cdot (id + id \times f) \quad (3.23)$$

and is easy to convert to pointwise notation:

$$\begin{aligned}
& f \cdot [\underline{Nil}, Cons] = [\underline{0}, add] \cdot (id + id \times f) \\
\equiv & \quad \{ (2.40) \text{ on the lefthand side, (2.41) and identity } id \text{ on the righthand side} \} \\
& [f \cdot \underline{Nil}, f \cdot Cons] = [\underline{0}, add \cdot (id \times f)] \\
\equiv & \quad \{ \textit{either structural equality (2.58)} \} \\
& \begin{cases} f \cdot \underline{Nil} = \underline{0} \\ f \cdot Cons = add \cdot (id \times f) \end{cases} \\
\equiv & \quad \{ \textit{going pointwise} \} \\
& \begin{cases} (f \cdot \underline{Nil})x = \underline{0}x \\ (f \cdot Cons)(a, x) = (add \cdot (id \times f))(a, x) \end{cases} \\
\equiv & \quad \{ \textit{composition (2.6), constant (2.12), product (2.22) and definition of } add \} \\
& \begin{cases} f Nil = 0 \\ f(Cons(a, x)) = a + f x \end{cases}
\end{aligned}$$

Note that we could have used  $out_{\top}$  in diagram (3.21),

$$\begin{array}{ccc}
\top & \xrightarrow{out_{\top}} & 1 + \mathbb{N}_0 \times \top \\
f \downarrow & & \downarrow id + id \times f \\
\mathbb{N}_0 & \xleftarrow{[\underline{0}, add]} & 1 + \mathbb{N}_0 \times \mathbb{N}_0
\end{array} \tag{3.24}$$

obtaining another version of the *definition* of  $f$ ,

$$f = [\underline{0}, add] \cdot (id + id \times f) \cdot out_{\top} \tag{3.25}$$

which would lead to exactly the same pointwise recursive definition:

$$\begin{aligned}
& f = [\underline{0}, add] \cdot (id + id \times f) \cdot out_{\top} \\
\equiv & \quad \{ (2.41) \text{ and identity } id \text{ on the righthand side} \} \\
& f = [\underline{0}, add \cdot (id \times f)] \cdot out_{\top} \\
\equiv & \quad \{ \textit{going pointwise on } out_{\top} \text{ (3.16)} \} \\
& \begin{cases} f Nil = ([\underline{0}, add \cdot (id \times f)] \cdot out_{\top}) Nil \\ f(Cons(a, x)) = ([\underline{0}, add \cdot (id \times f)] \cdot out_{\top})(a, x) \end{cases} \\
\equiv & \quad \{ \textit{definition of } out_{\top} \text{ (3.16)} \}
\end{aligned}$$

$$\begin{aligned}
& \left\{ \begin{array}{l} f \text{ Nil} = ([\underline{0}, \text{add} \cdot (\text{id} \times f)] \cdot i_1) \text{ Nil} \\ f(\text{Cons}(a, x)) = ([\underline{0}, \text{add} \cdot (\text{id} \times f)] \cdot i_2)(a, x) \end{array} \right. \\
\equiv & \quad \{ \text{+-cancellation (2.38)} \} \\
& \left\{ \begin{array}{l} f \text{ Nil} = \underline{0} \text{ Nil} \\ f(\text{Cons}(a, x)) = (\text{add} \cdot (\text{id} \times f))(a, x) \end{array} \right. \\
\equiv & \quad \{ \text{simplification} \} \\
& \left\{ \begin{array}{l} f \text{ Nil} = 0 \\ f(\text{Cons}(a, x)) = a + f x \end{array} \right.
\end{aligned}$$

Pointwise  $f$  mirrors the structure of type  $\mathbb{T}$  in having as many definition clauses as constructors in  $\mathbb{T}$ . 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  $\mathbb{T}$ , that is, for

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

— recall (3.10) — taking place of  $\text{out}_{\mathbb{T}}$ , we get a “more algorithmic” version of  $f$ :

$$\begin{aligned}
& f = [\underline{0}, \text{add}] \cdot (\text{id} + \text{id} \times f) \cdot (! + \langle \text{hd}, \text{tl} \rangle) \cdot (=_{[]}?) \\
\equiv & \quad \{ \text{+-functor (2.42), identity and } \times \text{-absorption (2.25)} \} \\
& f = [\underline{0}, \text{add}] \cdot (! + \langle \text{hd}, f \cdot \text{tl} \rangle) \cdot (=_{[]}?) \\
\equiv & \quad \{ \text{+-absorption (2.41) and constant } \underline{0} \} \\
& f = [\underline{0}, \text{add} \cdot \langle \text{hd}, f \cdot \text{tl} \rangle] \cdot (=_{[]}?) \\
\equiv & \quad \{ \text{going pointwise on guard } =_{[]}? \text{ (2.60) and simplifying} \} \\
& f l = \begin{cases} l = [] & \Rightarrow \underline{0} l \\ \neg(l = []) & \Rightarrow (\text{add} \cdot \langle \text{hd}, f \cdot \text{tl} \rangle) l \end{cases} \\
\equiv & \quad \{ \text{simplification} \} \\
& f l = \begin{cases} l = [] & \Rightarrow 0 \\ \neg(l = []) & \Rightarrow \text{hd} l + f(\text{tl} l) \end{cases}
\end{aligned}$$

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

```
f l | l == [] = 0
    | otherwise = head l + f (tail l)
```

or

```
f l = if l == []
      then 0
      else head l + f (tail l)
```

both requiring the equality predicate “==” and destructors “head” and “tail”.

### 3.4 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  $\mathbb{T}$  (3.15), we want to be able to synthesize (generate) particular inhabitants of  $\mathbb{T}$ . In other words, we want to be able to specify functions with signature  $B \xrightarrow{f} \mathbb{T}$  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  $\text{Nil}$ , in case  $n = 0$ .

Let us try and draw a diagram similar to (3.24) 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  $\text{in}_{\mathbb{T}}$  takes the place of  $\text{out}_{\mathbb{T}}$ :

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

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,

$$f = \text{in}_{\mathbb{T}} \cdot (\text{id} + \text{id} \times f) \cdot g \tag{3.26}$$

which is equivalent to:

$$f = [\text{Nil}, \text{Cons} \cdot (\text{id} \times f)] \cdot g \tag{3.27}$$

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,

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

$\xrightarrow{g}$

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  $n > 0$ ,

$$pred\ n \stackrel{\text{def}}{=} n - 1 \quad (3.28)$$

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

$$g = (! + \langle id, pred \rangle) \cdot (=0?) \quad (3.29)$$

Filling this in (3.27) we get

$$\begin{aligned} f &= [ \underline{Nil}, Cons \cdot (id \times f) ] \cdot (! + \langle id, pred \rangle) \cdot (=0?) \\ &\equiv \{ \text{+ -absorption (2.41) followed by } \times \text{-absorption (2.25) etc. } \} \\ f &= [ \underline{Nil}, Cons \cdot \langle id, f \cdot pred \rangle ] \cdot (=0?) \\ &\equiv \{ \text{going pointwise on guard } =_0? \text{ (2.60) and simplifying } \} \\ f\ n &= \begin{cases} n = 0 & \Rightarrow Nil \\ \neg(n = 0) & \Rightarrow Cons(n, f(n-1)) \end{cases} \end{aligned}$$

which matches the function we had in mind:

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

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.

### 3.5 Introducing (list) catas, anas and hylos

Suppose that, back to section 3.3, we want to *multiply*, rather than add, the elements occurring in lists of type  $T$  (3.15). 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  $\mathbb{N}_0 \xleftarrow{[\underline{0}, add]} 1 + \mathbb{N}_0 \times \mathbb{N}_0$  of diagram (3.24) 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,

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

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 3.4). Intuition will once again help us in deciding that only arrow  $g$  in (3.26) needs to be updated.

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

$$\begin{array}{ccc} \mathbb{T} & \xrightarrow{out_{\mathbb{T}}} & 1 + \mathbb{N}_0 \times \mathbb{T} \\ f \downarrow & & \downarrow id + id \times f \\ B & \xleftarrow{g} & 1 + \mathbb{N}_0 \times B \end{array} \quad \begin{array}{ccc} \mathbb{T} & \xleftarrow{in_{\mathbb{T}}} & 1 + \mathbb{N}_0 \times \mathbb{T} \\ f \uparrow & & \uparrow id + id \times f \\ B & \xrightarrow{g} & 1 + \mathbb{N}_0 \times B \end{array} \quad (3.31)$$

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 define<sup>3</sup>.

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  $\mathbb{T}$ -*catamorphism* induced by  $g$ ” and  $\llbracket g \rrbracket$  as “the  $\mathbb{T}$ -*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  $\mathbb{T}$ . For the moment, it suffices to say that

- the  $\mathbb{T}$ -catamorphism induced by  $B \xleftarrow{g} 1 + \mathbb{N}_0 \times B$  is the unique function  $B \xleftarrow{\langle g \rangle} \mathbb{T}$

<sup>3</sup>The theory which supports the statements of this paragraph will not be dealt with until chapter 7.

which obeys to property (or is defined by)

$$\langle g \rangle = g \cdot (id + id \times \langle g \rangle) \cdot out_{\top} \quad (3.32)$$

which is the same as

$$\langle g \rangle \cdot in_{\top} = g \cdot (id + id \times \langle g \rangle) \quad (3.33)$$

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

$$\langle g \rangle = in_{\top} \cdot (id + id \times \langle g \rangle) \cdot g \quad (3.34)$$

From (3.31) it can be observed that  $\top$  can act as a mediator between any  $\top$ -anamorphism and any  $\top$ -catamorphism, that is to say,  $B \xleftarrow{\langle g \rangle} \top$  composes with  $\top \xleftarrow{\langle h \rangle} C$ , for some  $C \xrightarrow{h} 1 + \mathbb{N}_0 \times C$ . In other words, a  $\top$ -catamorphism call always observe (consume) the output of a  $\top$ -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:

$$\begin{array}{ccc} B & \xleftarrow{g} & 1 + \mathbb{N}_0 \times B \\ \langle g \rangle \uparrow & & \uparrow id + id \times \langle g \rangle \\ \top & \xleftarrow{in_{\top}} & 1 + \mathbb{N}_0 \times \top \\ \langle h \rangle \uparrow & & \uparrow id + id \times \langle h \rangle \\ C & \xrightarrow{h} & 1 + \mathbb{N}_0 \times C \end{array} \quad (3.35)$$

What can we say about the  $\langle g \rangle \cdot \langle h \rangle$  composition? It is a function from  $B$  to  $C$  which resorts to  $\top$  as an *intermediate* data-structure and can be subject to the following calculation (cf. outermost rectangle in (3.35)):

$$\begin{aligned} \langle g \rangle \cdot \langle h \rangle &= g \cdot (id + id \times \langle g \rangle) \cdot (id + id \times \langle h \rangle) \cdot h \\ \equiv & \quad \{ \text{+-functor (2.42)} \} \\ \langle g \rangle \cdot \langle h \rangle &= g \cdot ((id \cdot id) + (id \times \langle g \rangle) \cdot (id \times \langle h \rangle)) \cdot h \\ \equiv & \quad \{ \text{identity and } \times\text{-functor (2.28)} \} \\ \langle g \rangle \cdot \langle h \rangle &= g \cdot (id + id \times \langle g \rangle \cdot \langle h \rangle) \cdot h \end{aligned}$$

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

$$\begin{array}{ccc} B & \xleftarrow{g} & 1 + \mathbf{N}_0 \times B \\ \langle g \rangle \cdot \llbracket h \rrbracket \uparrow & & \uparrow id + id \times \langle g \rangle \cdot \llbracket h \rrbracket \\ C & \xrightarrow{h} & 1 + \mathbf{N}_0 \times C \end{array} \quad (3.36)$$

As an example, let us see what comes out of  $\langle g \rangle \cdot \llbracket h \rrbracket$  for  $h$  and  $g$  respectively given by (3.29) and (3.30):

$$\begin{aligned} \langle g \rangle \cdot \llbracket h \rrbracket &= g \cdot (id + id \times \langle g \rangle \cdot \llbracket h \rrbracket) \cdot h \\ \equiv & \quad \{ \langle g \rangle \cdot \llbracket h \rrbracket \text{ abbreviated to } f \text{ and instantiating } h \text{ and } g \} \\ f &= [\underline{1}, mul] \cdot (id + id \times f) \cdot (! + \langle id, pred \rangle) \cdot (=0?) \\ \equiv & \quad \{ +\text{-functor (2.42) and identity} \} \\ f &= [\underline{1}, mul] \cdot (! + (id \times f) \cdot \langle id, pred \rangle) \cdot (=0?) \\ \equiv & \quad \{ \times\text{-absorption (2.25) and identity} \} \\ f &= [\underline{1}, mul] \cdot (! + \langle id, f \cdot pred \rangle) \cdot (=0?) \\ \equiv & \quad \{ +\text{-absorption (2.41) and constant } \underline{1} \text{ (2.15)} \} \\ f &= [\underline{1}, mul \cdot \langle id, f \cdot pred \rangle] \cdot (=0?) \\ \equiv & \quad \{ \text{McCarthy conditional (2.59)} \} \\ f &= (=0?) \rightarrow \underline{1}, mul \cdot \langle id, f \cdot pred \rangle \end{aligned}$$

Going pointwise, we get — via (2.59) —

$$\begin{aligned} f 0 &= [\underline{1}, mul \cdot \langle id, f \cdot pred \rangle](i_1 0) \\ &= \quad \{ +\text{-cancellation (2.38)} \} \\ & \quad \underline{1} 0 \\ &= \quad \{ \text{constant function (2.12)} \} \\ & \quad 1 \end{aligned}$$

and

$$\begin{aligned} f(n+1) &= [\underline{1}, mul \cdot \langle id, f \cdot pred \rangle](i_2(n+1)) \\ &= \quad \{ +\text{-cancellation (2.38)} \} \end{aligned}$$



$$\begin{aligned} & mul \cdot \langle id, f \cdot pred \rangle (n + 1) \\ = & \{ \text{pointwise definitions of } split, \text{ identity, predecessor and } mul \} \\ & (n + 1) \times f n \end{aligned}$$

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 (3.35) for the particular  $g$  and  $h$  we have considered above and recall a popular “definition” of factorial:

$$n! = \underbrace{n \times (n - 1) \times \dots \times 1}_{n \text{ times}} \quad (3.37)$$

In fact,  $\llbracket h \rrbracket n$  produces T-list

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

as an intermediate data-structure which is consumed by  $\langle g \rangle$ , the effect of which is but the “replacement” of  $Cons$  by  $\times$  and  $Nil$  by 1, therefore accomplishing (3.37) 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  $\langle g \rangle \cdot \llbracket h \rrbracket$  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  $\langle g \rangle \cdot \llbracket h \rrbracket$  of a T-catamorphism with a T-anamorphism is called a *T-hylomorphism*<sup>4</sup> 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 3.2.** 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

$$\text{sq } n \stackrel{\text{def}}{=} n^2 \quad (3.38)$$

<sup>4</sup>This terminology is derived from the Greek word  $\nu\lambda\omicron\sigma$  (hylōs) meaning “matter”.

as a  $\mathbb{T}$ -hylomorphism and encode it in HASKELL.

□

---

**Exercise 3.3.** Write function  $x^n$  as a  $\mathbb{T}$ -hylomorphism and encode it in HASKELL.

□

---

**Exercise 3.4.** The following function in HASKELL computes the  $\mathbb{T}$ -sequence of all even numbers less or equal to  $n$ :

```
f n = if n <= 1
      then Nil
      else Cons(m, f(m-2))
      where m = if even n then n else n-1
```

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

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

□

---

### 3.6 Inductive types more generally

So far we have focussed our attention exclusively to a particular inductive type  $\mathbb{T}$  (3.20) — 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  $F T$

$$\begin{array}{ccc}
 & \text{out}_T & \\
 & \curvearrowright & \\
 T & \cong & F T \\
 & \curvearrowleft & \\
 & \text{in}_T & 
 \end{array} \tag{3.39}$$

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

$$F X \stackrel{\text{def}}{=} 1 + \mathbb{N}_0 \times X \tag{3.40}$$

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

$$X \cong 1 + X \times \mathbb{N}_0 \tag{3.41}$$

— 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*:

$$X \cong 1 + X \times \mathbb{N}_0 \times X \tag{3.42}$$

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

$$X \cong F X \tag{3.43}$$

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

## 3.7 Functors

The concept of a functor  $F$ , borrowed from category theory, is a most generic and useful device in programming<sup>5</sup>. As we have seen,  $F$  can be regarded as a datatype constructor

<sup>5</sup>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

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  $F B \xleftarrow{F f} F A$  exists for every such  $f$ :

$$\begin{array}{ccc} A & \cdots & F A \\ f \downarrow & & \downarrow F f \\ B & \cdots & F B \end{array}$$

$F f$  extends  $f$  to  $F$ -structures and will, by definition, obey to two very basic properties: it commutes with identity

$$F id_A = id_{(F A)} \quad (3.44)$$

and with composition

$$F(g \cdot h) = (F g) \cdot (F h) \quad (3.45)$$

Two simple examples of a functor follow:

- Identity functor: define  $F X = X$ , for every datatype  $X$ , and  $F f = f$ . Properties (3.44) and (3.45) 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:

$$\begin{array}{ccc} A & \cdots & C \\ f \downarrow & & \downarrow id_C \\ B & \cdots & C \end{array}$$

Properties (3.44) and (3.45) hold trivially again.

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 (3.44) and (3.45) have to hold for every parameter of an  $n$ -ary functor. For a binary functor  $B$ , for instance, equation (3.44) becomes

$$B(id_A, id_B) = id_{B(A,B)} \quad (3.46)$$

---

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
$B^A$	$f$	$f^A$	Exponential

Table 3.1: Datatype constructions and associated operators.

and equation (3.45) becomes

$$B(g \cdot h, i \cdot j) = B(g, i) \cdot B(h, j) \quad (3.47)$$

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 (2.22). Properties (2.29) and (2.28) instantiate (3.46) and (3.47), 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 (2.37) — and functorial properties (2.43) and (2.42). 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 (2.73) and (2.74). All this is summarized in table 3.1.

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

$$(F \cdot G)X \stackrel{\text{def}}{=} F(GX) \quad (3.48)$$

$$(F \cdot G)f \stackrel{\text{def}}{=} F(Gf) \quad (3.49)$$

## 3.8 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 (3.6). 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 (3.40) is assumed, the diagrams in (3.31) can be written as simply as

$$\begin{array}{ccc}
 T & \xrightarrow{out_\tau} & FT \\
 f \downarrow & & \downarrow Ff \\
 B & \xleftarrow{g} & FB
 \end{array}
 \qquad
 \begin{array}{ccc}
 T & \xleftarrow{in_\tau} & FT \\
 f \uparrow & & \uparrow Ff \\
 B & \xrightarrow{g} & FB
 \end{array}
 \qquad (3.50)$$

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,

$$\begin{aligned}
 FX &\cong C_0 + (C_1 \times X) + (C_2 \times X^2) + \cdots + (C_n \times X^n) \\
 &= \sum_{i=0}^n C_i \times X^i
 \end{aligned}
 \qquad (3.51)$$

and that *Newton’s binomial formula*

$$(A + B)^n \cong \sum_{p=0}^n {}^n C_p \times A^{n-p} \times B^p
 \qquad (3.52)$$

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

$$F X \stackrel{\text{def}}{=} A \times (1 + X)^2$$

(where  $A$  is a constant datatype) and check the following reasoning:

$$\begin{aligned}
 F X &= A \times (1 + X)^2 \\
 &\cong \{ \text{law (2.87)} \} \\
 &A \times ((1 + X) \times (1 + X)) \\
 &\cong \{ \text{law (2.50)} \} \\
 &A \times ((1 + X) \times 1 + (1 + X) \times X) \\
 &\cong \{ \text{laws (2.81), (2.31) and (2.50)} \} \\
 &A \times ((1 + X) + (1 \times X + X \times X)) \\
 &\cong \{ \text{laws (2.81) and (2.87)} \} \\
 &A \times ((1 + X) + (X + X^2)) \\
 &\cong \{ \text{law (2.46)} \} \\
 &A \times (1 + (X + X) + X^2) \\
 &\cong \{ \text{canonical form obtained via laws (2.50) and (2.88)} \} \\
 &\underbrace{A}_{C_0} + \underbrace{A \times 2 \times X}_{C_1} + \underbrace{A}_{C_2} \times X^2
 \end{aligned}$$

**Exercise 3.5.** 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.

□

## 3.9 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 (3.43) is polynomial. For instance, datatype  $T$  (3.20) is polynomial ( $n = 1$ ) and its associated polynomial

functor is canonically defined with coefficients  $C_0 = 1$  and  $C_1 = \mathbf{N}_0$ . 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

$$\mathbb{T} \xleftarrow[\text{in}_{\mathbb{T}}]{\cong} \sum_{i=0}^n C_i \times \mathbb{T}^i \quad (3.53)$$

Then we have

$$\text{in}_{\mathbb{T}} \stackrel{\text{def}}{=} [f_1, \dots, f_n]$$

where, for  $i = 1, n$ ,  $f_i$  is an arrow of type  $\mathbb{T} \longleftarrow C_i \times \mathbb{T}^i$ . Since  $n$  is finite, one may expand exponentials according to (2.87) and encode this in HASKELL as follows:

```
data T = C0 |
        C1 (C1, T) |
        C2 (C2, (T, T)) |
        ... |
        Cn (Cn, (T, ..., T))
```

Of course the choice of symbol  $C_i$  to realize each  $f_i$  is arbitrary<sup>6</sup>. Several instances of polynomial inductive types (in canonical form) will be mentioned in section 3.13. Section 3.17 will address the conversion between inductive datatypes induced by so-called *natural transformations*.

The concepts of catamorphism, anamorphism and hylomorphism introduced in section 3.5 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 7.

### 3.10 F-algebras and F-homomorphisms

Our interest in polynomial types is basically due to the fact that, for polynomial  $F$ , equation (3.43) always has a particularly interesting solution which corresponds to our notion of a recursive datatype.

---

<sup>6</sup>A more traditional (but less close to (3.53)) encoding will be

$$\text{data T} = \text{C0} \mid \text{C1 C1 T} \mid \text{C2 C2 T T} \mid \dots \mid \text{Cn Cn T} \dots \text{T} \quad (3.54)$$

delivering every constructor in curried form.



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  $\alpha$  of signature  $A \xleftarrow{\alpha} F A$ .  $A$  is called the *carrier* of F-algebra  $\alpha$  and contains the values which  $\alpha$  manipulates by computing new  $A$ -values out of existing ones, according to the F-pattern (the “type” of the algebra). As examples, consider  $[\underline{0}, \text{add}]$  (3.19) and  $\text{in}_\top$  (3.20), 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  $\alpha$  of type  $F X = 1 + X \times X$  will be of form  $[\alpha_1, \alpha_2]$ , where  $\alpha_1$  is a constant and  $\alpha_2$  is a binary operator. So monoids are algebras of this type <sup>7</sup>.

Secondly, we introduce the notion of an *F-homomorphism* which is but a function observing a particular F-algebra  $\alpha$  into another F-algebra  $\beta$ :

$$\begin{array}{ccc}
 A & \xleftarrow{\alpha} & F A \\
 f \downarrow & & \downarrow F f \\
 B & \xleftarrow{\beta} & F B
 \end{array}
 \quad f \cdot \alpha = \beta \cdot (F f)
 \quad (3.55)$$

Clearly,  $f$  can be regarded as a structural translation between  $A$  and  $B$ , that is,  $A$  and  $B$  have a similar structure <sup>8</sup>. Note that — thanks to (3.44) — identity functions are always (trivial) F-homomorphisms and that — thanks to (3.45) — these homomorphisms compose, that is, the composition of two F-homomorphisms is an F-homomorphism.

## 3.11 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 (3.43). Moreover, for polynomial F a particular iso F-algebra always exists, which is denoted by  $\mu F \xleftarrow{\text{in}} F \mu F$  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  $\mu F$  —  $\mu$  for “minimal” <sup>9</sup>. Second, it is the so-called *initial* F-algebra. What does this mean?

It means that, for every F-algebra  $\alpha$  there exists one and only one F-homomorphism between  $\text{in}$  and  $\alpha$ . This unique arrow mediating  $\text{in}$  and  $\alpha$  is therefore determined by  $\alpha$  itself, and is called the *F-catamorphism* generated by  $\alpha$ . This construct, which was introduced in 3.5, is in general denoted by  $(\alpha)_F$ :

<sup>7</sup>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.*

<sup>8</sup>Cf. *homomorphism* = *homo* (the same) + *morphos* (structure, shape).

<sup>9</sup> $\mu F$  means the least fixpoint solution of equation  $X \cong F X$ , as will be described in chapter 7.

$$\begin{array}{ccc}
 \mu F & \xleftarrow{in} & F \mu F \\
 f = \langle \alpha \rangle_F \downarrow & & \downarrow F \langle \alpha \rangle_F \\
 A & \xleftarrow{\alpha} & F A
 \end{array} \tag{3.56}$$

We will drop the  $F$  subscript in  $\langle \alpha \rangle_F$  wherever deducible from the context, and often call  $\alpha$  the “gene” of  $\langle \alpha \rangle_F$ .

As happens with *splits*, *eithers* and *transposes*, the uniqueness of the catamorphism construct is captured by a universal property established in the class of all  $F$ -homomorphisms:

$$k = \langle \alpha \rangle \Leftrightarrow k \cdot in = \alpha \cdot F k \tag{3.57}$$

According to the experience gathered from section 2.12 onwards, a few properties can be expected as consequences of (3.57). For instance, one may wonder about the “gene” of the identity catamorphism. Just let  $k = id$  in (3.57) and see what happens:

$$\begin{aligned}
 id &= \langle \alpha \rangle \Leftrightarrow id \cdot in = \alpha \cdot F id \\
 &= \{ \text{identity (2.10) and } F \text{ is a functor (3.44)} \} \\
 id &= \langle \alpha \rangle \Leftrightarrow in = \alpha \cdot id \\
 &= \{ \text{identity (2.10) once again} \} \\
 id &= \langle \alpha \rangle \Leftrightarrow in = \alpha \\
 &= \{ \alpha \text{ replaced by } in \text{ and simplifying} \} \\
 id &= \langle in \rangle
 \end{aligned}$$

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 :**

$$\begin{array}{ccc}
 \mu F & \xleftarrow{in} & F \mu F \\
 \langle in \rangle \downarrow & & \downarrow F \langle in \rangle \\
 \mu F & \xleftarrow{in} & F \mu F
 \end{array} \tag{3.58} \quad \langle in \rangle = id_{\mu F}$$

In a more intuitive way, one might have observed that  $\langle in \rangle$  is, by definition of  $in$ , the unique arrow mediating  $\mu F$  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 (3.57), for  $k = \langle \alpha \rangle$ , is

**Cata-cancellation :**

$$(\downarrow\alpha) \cdot in = \alpha \cdot F(\downarrow\alpha) \quad (3.59)$$

Because *in* is iso, this law can be rephrased as follows

$$(\downarrow\alpha) = \alpha \cdot F(\downarrow\alpha) \cdot out \quad (3.60)$$

where *out* denotes the inverse of *in*:

$$\begin{array}{ccc} & \xrightarrow{out} & \\ \mu F & \cong & F \mu F \\ & \xleftarrow{in} & \end{array}$$

Now, let  $f$  be F-homomorphism (3.55) between F-algebras  $\alpha$  and  $\beta$ . How does it relate to  $(\downarrow\alpha)$  and  $(\downarrow\beta)$ ? Note that  $f \cdot (\downarrow\alpha)$  is an arrow mediating  $\mu F$  and  $B$ . But  $B$  is the carrier of  $\beta$  and  $(\downarrow\beta)$  is the unique arrow mediating  $\mu F$  and  $B$ . So the two arrows are the same:

**Cata-fusion :**

$$\begin{array}{ccc} \mu F & \xleftarrow{in} & F \mu F \\ (\downarrow\alpha) \downarrow & & \downarrow F(\downarrow\alpha) \\ A & \xleftarrow{\alpha} & F A \\ f \downarrow & & \downarrow F f \\ B & \xleftarrow{\beta} & F B \end{array} \quad f \cdot (\downarrow\alpha) = (\downarrow\beta) \quad \text{if} \quad f \cdot \alpha = \beta \cdot F f \quad (3.61)$$

Of course, this law is also a consequence of the universal property, for  $k = f \cdot (\downarrow\alpha)$ :

$$\begin{aligned} f \cdot (\downarrow\alpha) = (\downarrow\beta) &\Leftrightarrow (f \cdot (\downarrow\alpha)) \cdot in = \beta \cdot F(f \cdot (\downarrow\alpha)) \\ &\Leftrightarrow \{ \text{composition is associative and } F \text{ is a functor (3.45)} \} \\ &\quad f \cdot ((\downarrow\alpha) \cdot in) = \beta \cdot (F f) \cdot (F(\downarrow\alpha)) \\ &\Leftrightarrow \{ \text{cata-cancellation (3.59)} \} \\ &\quad f \cdot \alpha \cdot F(\downarrow\alpha) = \beta \cdot F f \cdot F(\downarrow\alpha) \\ &\Leftrightarrow \{ \text{require } f \text{ to be a F-homomorphism (3.55)} \} \\ &\quad f \cdot \alpha \cdot F(\downarrow\alpha) = f \cdot \alpha \cdot F(\downarrow\alpha) \wedge f \cdot \alpha = \beta \cdot F f \\ &\Leftrightarrow \{ \text{simplify} \} \\ &\quad f \cdot \alpha = \beta \cdot F f \end{aligned}$$

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.

### 3.12 Parameterization and type functors

By analogy with what we have done about *splits* (product), *eithers* (coproduct) and *transposes* (exponential), we now look forward to identifying F-catamorphisms which exhibit functorial behaviour.

Suppose that one wishes to square all numbers which are members of lists of type  $\mathbb{T}$  (3.20). It can be checked that

$$(\llbracket \underline{Nil}, Cons \cdot (sq \times id) \rrbracket) \quad (3.62)$$

will do this for us, where  $\mathbb{N}_0 \xleftarrow{sq} \mathbb{N}_0$  is given by (3.38). This catamorphism, which converted to pointwise notation is nothing but function  $h$  which follows

$$\begin{cases} h Nil = Nil \\ h(Cons(a, l)) = Cons(sq a, h l) \end{cases}$$

maps type  $\mathbb{T}$  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 (3.62). However, the output type no longer is  $\mathbb{T}$ , but rather type  $\mathbb{T}' \cong 1 + B \times \mathbb{T}'$ .

Types  $\mathbb{T}$  and  $\mathbb{T}'$  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  $\mathbb{T}$  and  $B$  in  $\mathbb{T}'$ . This suggests that these two types can be regarded as instances of a more generic list datatype `List`

$$\text{List } X \xleftarrow[\text{in}=\llbracket \underline{Nil}, Cons \rrbracket]{\cong} 1 + X \times \text{List } X \quad (3.63)$$

in which the type of elements  $X$  is allowed to vary. Thus one has  $\mathbb{T} = \text{List } \mathbb{N}_0$  and  $\mathbb{T}' = \text{List } B$ .

By inspection, it can be checked that, for every  $B \xleftarrow{f} A$ ,

$$(\llbracket \underline{Nil}, Cons \cdot (f \times id) \rrbracket) \quad (3.64)$$

maps `List A` to `List B`. Moreover, for  $f = id$  one has:

$$\begin{aligned} & (\llbracket \underline{Nil}, Cons \cdot (id \times id) \rrbracket) \\ = & \quad \{ \text{by the } \times\text{-functor-id property (2.29) and identity} \} \\ & (\llbracket \underline{Nil}, Cons \rrbracket) \\ = & \quad \{ \text{cata-reflection (3.58)} \} \\ & id \end{aligned}$$

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 (3.44) for  $F$  replaced by  $\text{List}$ . 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.* (3.45). Altogether, this means that  $\text{List}$  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

$$T \cong FT$$

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  $T$ )

$$TX \cong B(X, TX)$$

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  $FY = B(\mathbb{N}_0, Y) = 1 + \mathbb{N}_0 \times Y$ , recall (3.40). Moreover, one has

$$Ff = B(id, f) \tag{3.65}$$

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 (3.55).

$\mathbb{T}X$  will be referred to as the *type functor* generated by  $B$ :

$$\underbrace{\mathbb{T}X}_{\text{type functor}} \cong \underbrace{B(X, \mathbb{T}X)}_{\text{base functor}}$$

We proceed to the description of its functorial behaviour —  $\mathbb{T}f$  — for a given  $B \xleftarrow{f} A$ . As far as typing rules are concerned, we shall have

$$\frac{B \xleftarrow{f} A}{\mathbb{T}B \xleftarrow{\mathbb{T}f} \mathbb{T}A}$$

So we should be able to express  $\mathbb{T}f$  as a  $B(A, -)$ -catamorphism  $(\llbracket g \rrbracket)$ :

$$\begin{array}{ccc} A & & \mathbb{T}A \xleftarrow{\text{in}_{\mathbb{T}A}} B(A, \mathbb{T}A) \\ f \downarrow & & \downarrow \mathbb{T}f = \llbracket g \rrbracket \quad \downarrow B(id, \mathbb{T}f) \\ B & & \mathbb{T}B \xleftarrow{g} B(A, \mathbb{T}B) \end{array}$$

As we know that  $\text{in}_{\mathbb{T}B}$  is the standard constructor of values of type  $\mathbb{T}B$ , we may put it into the diagram too:

$$\begin{array}{ccc} A & & \mathbb{T}A \xleftarrow{\text{in}_{\mathbb{T}A}} B(A, \mathbb{T}A) \\ f \downarrow & & \downarrow \mathbb{T}f = \llbracket g \rrbracket \quad \downarrow B(id, \mathbb{T}f) \\ B & & \mathbb{T}B \xleftarrow{g} B(A, \mathbb{T}B) \\ & & \swarrow \text{in}_{\mathbb{T}B} \quad \searrow \text{dotted} \\ & & B(B, \mathbb{T}B) \end{array}$$

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

$$\mathbb{T}f \stackrel{\text{def}}{=} (\llbracket \text{in}_{\mathbb{T}B} \cdot B(f, id) \rrbracket) \tag{3.66}$$

and a final diagram, where  $\text{in}_{\mathbb{T}A}$  is abbreviated by  $\text{in}_A$  (ibid.  $\text{in}_{\mathbb{T}B}$  by  $\text{in}_B$ ):

$$\begin{array}{ccc} A & & \mathbb{T}A \xleftarrow{\text{in}_A} B(A, \mathbb{T}A) \\ f \downarrow & & \downarrow \mathbb{T}f = \llbracket \text{in}_B \cdot B(f, id) \rrbracket \quad \downarrow B(id, \mathbb{T}f) \\ B & & \mathbb{T}B \xleftarrow{\text{in}_B} B(B, \mathbb{T}B) \xleftarrow{B(f, id)} B(A, \mathbb{T}B) \end{array}$$

Next, we proceed to derive the useful law of *cata-absorption*

$$\llbracket g \rrbracket \cdot \mathbb{T} f = \llbracket g \cdot \mathbb{B}(f, id) \rrbracket \quad (3.67)$$

as a consequence of the laws studied in section 3.11. Our target is to show that, for  $k = \llbracket g \rrbracket \cdot \mathbb{T} f$  in (3.57), one gets  $\alpha = g \cdot \mathbb{B}(f, id)$ :

$$\begin{aligned} & \llbracket g \rrbracket \cdot \mathbb{T} f = \llbracket \alpha \rrbracket \\ \Leftrightarrow & \quad \{ \text{type-functor definition (3.66)} \} \\ & \llbracket g \rrbracket \cdot \llbracket in_B \cdot \mathbb{B}(f, id) \rrbracket = \llbracket \alpha \rrbracket \\ \Leftarrow & \quad \{ \text{cata-fusion (3.61)} \} \\ & \llbracket g \rrbracket \cdot in_B \cdot \mathbb{B}(f, id) = \alpha \cdot \mathbb{B}(id, \llbracket g \rrbracket) \\ \Leftrightarrow & \quad \{ \text{cata-cancellation (3.59)} \} \\ & g \cdot \mathbb{B}(id, \llbracket g \rrbracket) \cdot \mathbb{B}(f, id) = \alpha \cdot \mathbb{B}(id, \llbracket g \rrbracket) \\ \Leftrightarrow & \quad \{ \mathbb{B} \text{ is a bi-functor (3.47)} \} \\ & g \cdot \mathbb{B}(id \cdot f, \llbracket g \rrbracket \cdot id) = \alpha \cdot \mathbb{B}(id, \llbracket g \rrbracket) \\ \Leftrightarrow & \quad \{ id \text{ is natural (2.11)} \} \\ & g \cdot \mathbb{B}(f \cdot id, id \cdot \llbracket g \rrbracket) = \alpha \cdot \mathbb{B}(id, \llbracket g \rrbracket) \\ \Leftrightarrow & \quad \{ (3.47) \text{ again, this time from left to right} \} \\ & g \cdot \mathbb{B}(f, id) \cdot \mathbb{B}(id, \llbracket g \rrbracket) = \alpha \cdot \mathbb{B}(id, \llbracket g \rrbracket) \\ \Leftarrow & \quad \{ \text{obvious} \} \\ & g \cdot \mathbb{B}(f, id) = \alpha \end{aligned}$$

The following diagram pictures this property of catamorphisms:

$$\begin{array}{ccccc} & & \mathbb{T} A & \xleftarrow{in_A} & \mathbb{B}(A, \mathbb{T} A) \\ & & \mathbb{T} f \downarrow & & \downarrow \mathbb{B}(id, \mathbb{T} f) \\ A & & \mathbb{T} C & \xleftarrow{in_C} & \mathbb{B}(C, \mathbb{T} C) \xleftarrow{\mathbb{B}(f, id)} \mathbb{B}(A, \mathbb{T} C) \\ \downarrow f & & \downarrow \llbracket g \rrbracket & & \downarrow \mathbb{B}(id, \llbracket g \rrbracket) \\ C & & D & \xleftarrow{g} & \mathbb{B}(C, D) \xleftarrow{\mathbb{B}(f, id)} \mathbb{B}(A, D) \end{array}$$

It remains to show that (3.66) indeed defines a functor. This can be verified by checking properties (3.44) and (3.45) for  $F = \mathbb{T}$ :

- Property **type-functor-id**, cf. (3.44):

$$\begin{aligned}
& \top id \\
= & \quad \{ \text{by definition (3.66)} \} \\
& (in_B \cdot B(id, id)) \\
= & \quad \{ B \text{ is a bi-functor (3.46)} \} \\
& (in_B \cdot id) \\
= & \quad \{ \text{identity and cata-reflection (3.58)} \} \\
& id
\end{aligned}$$

- Property **type-functor**, cf. (3.45) :

$$\begin{aligned}
& \top (f \cdot g) \\
= & \quad \{ \text{by definition (3.66)} \} \\
& (in_B \cdot B(f \cdot g, id)) \\
= & \quad \{ id \cdot id = id \text{ and } B \text{ is a bi-functor (3.47)} \} \\
& (in_B \cdot B(f, id) \cdot B(g, id)) \\
= & \quad \{ \text{cata-absorption (3.67)} \} \\
& (in_B \cdot B(f, id)) \cdot \top g \\
= & \quad \{ \text{again cata-absorption (3.67)} \} \\
& (in_B) \cdot \top f \cdot \top g \\
= & \quad \{ \text{cata-reflection (3.58) followed by identity} \} \\
& \top f \cdot \top g
\end{aligned}$$

### 3.13 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 (3.53). The table contains two extra columns which may be used as bookmarks for equations (3.65) and (3.66), respec-



### 3.13. A CATALOGUE OF STANDARD POLYNOMIAL INDUCTIVE TYPES<sup>91</sup>

tively<sup>10</sup>:

Description	$\mathbb{T} X$	$\mathbb{B}(X, Y)$	$\mathbb{B}(id, f)$	$\mathbb{B}(f, id)$
“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 $X$	$X + Y^2$	$id + f^2$	$f + id$

(3.68)

All type functors  $\mathbb{T}$  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  $\mathbb{T}$ , then this will be based on an  $n + 1$ -ary base functor  $\mathbb{B}$ , cf.

$$\mathbb{T}(X_1, \dots, X_n) \cong \mathbb{B}(X_1, \dots, X_n, \mathbb{T}(X_1, \dots, X_n))$$

So, every  $n + 1$ -ary polynomial functor  $\mathbb{B}(X_1, \dots, X_n, X_{n+1})$  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 (3.68) leads to the so-called “full tree” type bi-functor

Description	$\mathbb{T}(X_1, X_2)$	$\mathbb{B}(X_1, X_2, Y)$	$\mathbb{B}(id, id, f)$	$\mathbb{B}(f, g, id)$
“Full” Trees	FTree( $X_1, X_2$ )	$X_1 + X_2 \times Y^2$	$id + id \times f^2$	$f + g \times id$

(3.69)

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 (2.87), see for instance

```
data BTree a = Empty | Node(a, (BTree a, BTree a))
```

Moreover, one may chose to curry the type constructors as in, e.g.

```
data BTree a = Empty | Node a (BTree a) (BTree a)
```

**Exercise 3.6.** Write as a catamorphisms

- the function which counts the number of elements of a non-empty list (type NList in (3.68)).

<sup>10</sup>Since  $(id_A)^2 = id_{(A^2)}$  one writes  $id^2$  for  $id$  in this table.

- the function which computes the maximum element of a binary-tree of natural numbers.

□

**Exercise 3.7.** Characterize the function which is defined by  $([\ ] , h)$  for each of the following definitions of  $h$ :

$$h(x, (y_1, y_2)) = y_1 ++ [x] ++ y_2 \quad (3.70)$$

$$h = ++ \cdot (\text{singl} \times ++ ) \quad (3.71)$$

$$h = ++ \cdot ( ++ \times \text{singl} ) \cdot \text{swap} \quad (3.72)$$

assuming  $\text{singl } a = [a]$ . Identify in (3.68) which datatypes are involved as base functors.

□

**Exercise 3.8.** Write as a catamorphism the function which computes the frontier of a tree of type `LTree` (3.68), listed from left to right.

□

## 3.14 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  $Fg$  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 (3.68) 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 3.9.** *Declare all datatypes in (3.68) in HASKELL notation and turn them into HASKELL type functors, that is, define `fmap` in each case.*

□

---

**Exercise 3.10.** *Declare datatype (3.69) in HASKELL notation and turn it into an instance of class `BiFunctor`.*

□

---

## 3.15 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* (2.47) — which stemmed from putting together the two universal properties of product and coproduct, (2.55) and (2.57), 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 (2.55), such catamorphisms are bound to be *splits*, as are the corresponding *genes*:

$$\begin{array}{ccc} T & \xleftarrow{in} & FT \\ \downarrow \langle \langle h, k \rangle \rangle & & \downarrow F \langle \langle h, k \rangle \rangle \\ A \times B & \xleftarrow{\langle h, k \rangle} & F(A \times B) \end{array}$$

As we did for the exchange rule, we put (2.55) and the universal property of catamorphisms (3.57) against each other and calculate:

$$\begin{aligned} \langle f, g \rangle &= \langle \langle h, k \rangle \rangle \\ \equiv & \quad \{ \text{cata-universal (3.57)} \} \\ \langle f, g \rangle \cdot in &= \langle h, k \rangle \cdot F \langle f, g \rangle \\ \equiv & \quad \{ \times\text{-fusion (2.24) twice} \} \\ \langle f \cdot in, g \cdot in \rangle &= \langle h \cdot F \langle f, g \rangle, k \cdot F \langle f, g \rangle \rangle \\ \equiv & \quad \{ (2.56) \} \\ f \cdot in &= h \cdot F \langle f, g \rangle \quad \wedge \quad g \cdot in = k \cdot F \langle f, g \rangle \end{aligned}$$

The rule thus obtained,

$$\left\{ \begin{array}{l} f \cdot in = h \cdot F \langle f, g \rangle \\ g \cdot in = k \cdot F \langle f, g \rangle \end{array} \right. \equiv \langle f, g \rangle = \langle \langle h, k \rangle \rangle \quad (3.73)$$

is referred to as the *mutual recursion law* (or as “Fokkinga’s law”) and is useful in combining two mutually recursive functions  $f$  and  $g$

$$\begin{array}{ccc} T & \xleftarrow{in} & FT \\ f \downarrow & & \downarrow F \langle f, g \rangle \\ A & \xleftarrow{h} & F(A \times B) \end{array} \quad \begin{array}{ccc} T & \xleftarrow{in} & FT \\ g \downarrow & & \downarrow F \langle f, g \rangle \\ B & \xleftarrow{k} & F(A \times B) \end{array}$$

into a single catamorphism.

When applied from left to right, law (3.73) 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:

$$\begin{aligned} fib\ 0 &= 1 \\ fib\ 1 &= 1 \\ fib(n+2) &= fib(n+1) + fib\ n \end{aligned}$$

It can be shown that  $fib$  is a hylomorphism of type LTree (3.68),  $fib = \llbracket count, fibd \rrbracket$ , for  $count = [\underline{1}, add]$ ,  $add(x, y) = x + y$  and  $fibd\ n = if\ n < 2\ then\ i_1 Nil\ else\ i_2(n - 1, n - 2)$ . This hylo-factorization of  $fib$  tells its internal algorithmic structure: the *divide step*  $\llbracket fibd \rrbracket$  builds a tree whose number of leaves is a Fibonacci number; the *conquer step*  $\llbracket count \rrbracket$  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,

$$f\ n = fib(n + 1)$$

Clearly,  $f(n + 1) = fib(n + 2) = f\ n + fib\ n$  and  $f\ 0 = fib\ 1 = 1$ . Putting  $f$  and  $fib$  together,

$$\begin{aligned} f\ 0 &= 1 \\ f(n + 1) &= f\ n + fib\ n \\ fib\ 0 &= 1 \\ fib(n + 1) &= f\ n \end{aligned}$$

we obtain two mutually recursive functions over the natural numbers ( $\mathbb{N}_0$ ) which transform into pointfree equalities

$$\begin{aligned} f \cdot [\underline{0}, suc] &= [\underline{1}, add \cdot \langle f, fib \rangle] \\ fib \cdot [\underline{0}, suc] &= [\underline{1}, f] \end{aligned}$$

over

$$\mathbb{N}_0 \begin{array}{c} \xrightarrow{\quad} \\ \cong \\ \xleftarrow{in=[\underline{0}, suc]} \end{array} \underbrace{1 + \mathbb{N}_0}_{F\ \mathbb{N}_0} \quad (3.74)$$

Reverse +-absorption (2.41) will further enable us to rewrite the above into

$$\begin{aligned} f \cdot in &= [\underline{1}, add] \cdot F \langle f, fib \rangle \\ fib \cdot in &= [\underline{1}, \pi_1] \cdot F \langle f, fib \rangle \end{aligned}$$

thus bringing functor  $F\ f = id + f$  explicit and preparing for mutual recursion removal:

$$\begin{aligned} f \cdot in &= [\underline{1}, add] \cdot F \langle f, fib \rangle \\ fib \cdot in &= [\underline{1}, \pi_1] \cdot F \langle f, fib \rangle \end{aligned}$$

$$\begin{aligned}
&\equiv \{ (3.73) \} \\
&\langle f, fib \rangle = (\langle [\underline{1}, add], [\underline{1}, \pi_1] \rangle) \\
&\equiv \{ \text{exchange law (2.47)} \} \\
&\langle f, fib \rangle = (\langle [\underline{1}, \underline{1}], \langle add, \pi_1 \rangle \rangle) \\
&\equiv \{ \text{going pointwise and denoting } \langle f, fib \rangle \text{ by } fib' \} \\
&\left\{ \begin{array}{l} fib' 0 = (1, 1) \\ fib' (n + 1) = (x + y, x) \text{ where } (x, y) = fib' n \end{array} \right.
\end{aligned}$$

Since  $fib = \pi_2 \cdot fib'$  we easily recover  $fib$  from  $fib'$  and obtain the intended linear version of Fibonacci (encoded in Haskell):

```

fib n = y where (x,y) = 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” 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.

Our intuition above is confirmed by observing that all  $\mathbf{N}_0$  catamorphisms are of shape  $(\llbracket \underline{k}, g \rrbracket)$ , and that  $(\llbracket \underline{k}, g \rrbracket)n = g^n k$ , where  $g^n$  is the  $n$ -th iteration of  $g$ , that is,  $g^0 = id$  and  $g^{n+1} = g \cdot g^n$ . So  $g$  is the body of a “for-loop” which repeats itself  $n$ -times, starting with initial value  $k$ .

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 (3.74),

$$\begin{aligned} bnm\ 0 &= 1 \\ bnm(n+1) &= 2n+3 + bnm\ n \end{aligned}$$

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  $lin\ n \stackrel{\text{def}}{=} 2^n+3$  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 3.11 below shows how to extend (3.73) to three mutually recursive functions (3.75). (From this it is easy to extend it further to the  $n$ -ary case.) It is routine work to show that, by application of (3.75) to the above three functions, one obtains the linear version of  $ns$  which follows:

```

ns'' n = let (a,b,c) = aux n in a
      where
        aux 0 = (0,1,3)
        aux(n+1) = let (a,b,c) = aux n
                   in (a+b,b+c,2+c)

```

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  $\phi_1, \dots, \phi_n$  can be handled by a suitably extended version of (3.73). This only happens if all  $f_i$  have the same “shape”, that is, if they share the same base functor  $F$ .

**Exercise 3.11.** Show that law (3.73) generalizes to more than two mutually recursive functions, in this case three:

$$\begin{cases} f \cdot in = h \cdot F \langle f, \langle g, j \rangle \rangle \\ g \cdot in = k \cdot F \langle f, \langle g, j \rangle \rangle \\ j \cdot in = l \cdot F \langle f, \langle g, j \rangle \rangle \end{cases} \equiv \langle f, \langle g, j \rangle \rangle = \langle \langle h, \langle k, l \rangle \rangle \rangle \quad (3.75)$$

□

**Exercise 3.12.** 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:

$$e^x = \sum_{n=0}^{\infty} \frac{x^n}{n!} \quad (3.76)$$

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 `exp 1 1 = 2.0`, `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  $\text{exp}$  into mutually recursive functions so as to apply (3.75) and obtain the following linear version:

$$\begin{aligned} \text{exp } x \ n &= \text{let } (e, b, c) = \text{aux } x \ n \\ &\quad \text{in } e \text{ where} \\ &\quad \text{aux } x \ 0 = (1, 2, x) \\ &\quad \text{aux } x \ (n+1) = \text{let } (e, s, h) = \text{aux } x \ n \\ &\quad \quad \text{in } (e+h, s+1, (x/s)*h) \end{aligned}$$

□

**Exercise 3.13.** From the following basic properties of addition and multiplication,

$$a * 0 = 0 \tag{3.77}$$

$$a * 1 = a \tag{3.78}$$

$$a * (b + c) = a * b + a * c \tag{3.79}$$

show that  $a * n$  is the “for-loop”  $(a+)^n 0$ .

□

**Exercise 3.14.** Show that, for all  $n \in \mathbb{N}_0$ ,  $n = \text{suc}^n 0$ . **Hint:** use cata-reflexion (3.58).

□

As example of application of (3.73) for  $T$  other than  $\mathbb{N}_0$ , consider the following recursive predicate which checks whether a (non-empty) list is ordered,

$$\begin{aligned} \text{ord} : A^+ &\longrightarrow 2 \\ \text{ord } [a] &= \text{TRUE} \\ \text{ord } (\text{cons}(a, l)) &= a \geq (\text{listMax } l) \wedge (\text{ord } l) \end{aligned}$$

where  $\geq$  is assumed to be a total order on datatype  $A$  and

$$listMax = ([ [ id, max ] ]) \quad (3.80)$$

computes the greatest element of a given list of  $A$ s:

$$\begin{array}{ccc} A^+ & \xleftarrow{[singl, cons]} & A + A \times A^+ \\ listMax \downarrow & & \downarrow id + id \times listMax \\ A & \xleftarrow{[id, max]} & A + A \times A \end{array}$$

(In the diagram,  $singl\ a = [a]$ .)

Predicate  $ord$  is not a catamorphism because of the presence of  $listMax\ l$  in the recursive branch. However, the following diagram depicting  $ord$

$$\begin{array}{ccc} A^+ & \xleftarrow{[singl, cons]} & A + A \times A^+ \\ ord \downarrow & & \downarrow id + id \times (listMax, ord) \\ 2 & \xleftarrow{[TRUE, \alpha]} & A + A \times (A \times 2) \end{array}$$

(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

$$\begin{array}{ccc} A^+ & \xleftarrow{[singl, cons]} & A + A \times A^+ \\ listMax \downarrow & & \downarrow id + id \times (listMax, g) \\ A & \xleftarrow{[id, max \cdot (id \times \pi_1)]} & A + A \times (A \times B) \end{array}$$

where the extra presence of  $g$  is cancelled by projection  $\pi_1$ .

For  $B = 2$  and  $g = ord$  we are in position to apply Fokkinga's law and obtain:

$$\begin{aligned} \langle listMax, ord \rangle &= ([ [ id, max \cdot (id \times \pi_1) ], [ TRUE, \alpha ] ]) \\ &= \{ \text{exchange law (2.47)} \} \\ &= ([ [ id, TRUE ], [ max \cdot (id \times \pi_1), \alpha ] ]) \end{aligned}$$

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  $ord$ :

$$ord\ l = \text{let } (a, b) = aux\ l \\ \text{in } b$$

### 3.16. “BANANA-SPLIT”: A COROLLARY OF THE MUTUAL-RECURSION LAW 101

where

$$\begin{aligned}
 aux &: A^+ \longrightarrow A \times 2 \\
 aux [a] &= (a, \text{TRUE}) \\
 aux (\text{cons}(a, l)) &= \text{let } (m, b) = aux\ l \\
 &\quad \text{in } (\max(a, m), (a > m \wedge b))
 \end{aligned}$$

## 3.16 “Banana-split”: a corollary of the mutual-recursion law

Let  $h = i \cdot F \pi_1$  and  $k = j \cdot F \pi_2$  in (3.73). Then

$$\begin{aligned}
 f \cdot in &= (i \cdot F \pi_1) \cdot F \langle f, g \rangle \\
 \equiv &\quad \{ \text{composition is associative and } F \text{ is a functor } \} \\
 f \cdot in &= i \cdot F (\pi_1 \cdot \langle f, g \rangle) \\
 \equiv &\quad \{ \text{by } \times\text{-cancellation (2.20)} \} \\
 f \cdot in &= i \cdot F f \\
 \equiv &\quad \{ \text{by cata-cancellation} \} \\
 f &= \langle i \rangle
 \end{aligned}$$

Similarly, from  $k = j \cdot F \pi_2$  we get

$$g = \langle j \rangle$$

Then, from (3.73), 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

$$\langle \langle i \rangle, \langle j \rangle \rangle = \langle (i \times j) \cdot \langle F \pi_1, F \pi_2 \rangle \rangle \quad (3.81)$$

by (reverse)  $\times$ -absorption (2.25).

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,

$$average \stackrel{\text{def}}{=} (/) \cdot \langle sum, length \rangle \quad (3.82)$$

where  $sum$  and  $length$  are the expected  $\mathbb{N}^+$  catamorphisms:

$$\begin{aligned} sum &= ([id, +]) \\ length &= ([1, succ \cdot \pi_2]) \end{aligned}$$

As defined by (3.82), 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":

$$\begin{aligned} average\ l &= x/y \\ \text{where } (x, y) &= aux\ l \\ aux[a] &= (a, 1) \\ aux(cons(a, l)) &= (a + x, y + 1) \end{aligned} \tag{3.83}$$

*where*  $(x, y) = aux\ l$

**Exercise 3.15.** Calculate (3.83) from (3.82). Which of these two versions of the same function is easier to understand?

□

### 3.17 Inductive datatype isomorphism

not yet available

### 3.18 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 [MA86], was popularized by [Mal90] and reached textbook format in [BdM97]. A comfortable basis for exploiting *polymorphism* [Wad89], the “datatypes as functors” motto has proved beneficial at a higher level of abstraction, giving birth to *polytypism* [JJ96].

The literature on *anas*, *catas* and *hylos* is vast (see eg. [MH95], [JJ98], [GHA01]) and it is part of a broader discipline which has become known as the *mathematics of program construction* [Bac04]. This chapter provides an introduction to such a discipline. Only the calculus of catamorphisms is presented. The corresponding theory of anamorphisms and hylomorphisms demands further mathematical machinery and will not be dealt with before chapters 5 and 7. The results on mutual recursion presented in this chapter were pioneered by Maarten Fokkinga [Fok92].



# Chapter 4

## Why Monads Matter

In this chapter we present a powerful device in state-of-the-art 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, 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

Consider the  $\mathbb{R}$  to  $\mathbb{R}$  function

$$g\ x \stackrel{\text{def}}{=} 1/x$$

Clearly,  $g$  is undefined for  $x = 0$  because  $g\ 0 = 1/0$  is so big a real number that it cannot be properly evaluated. In fact, the HASKELL output for  $g\ 0 = 1/0$  is just “panic”:

```
Main> g 0
```

```
Program error: {primDivDouble 1.0 0.0}
```

```
Main>
```

Functions such as  $g$  above are called *partial functions* because they cannot be applied to all of their inputs (*i.e.*, they diverge for some of their inputs). Partial functions are very common in mathematics or programming — for other examples think of *e.g.* list-processing functions `head` and `tail`.

Panic is very dangerous in programming. In order to avoid this kind of behaviour one has two alternatives, either ensuring that every call to  $gx$  is *protected* — *i.e.*, the contexts which wrap up such calls ensure *pre-condition*  $x \neq 0$ , or one *raises* exceptions, *i.e.* explicit error values. In the former case, mathematical proofs need to be carried out in order to guarantee *safety* (that is, *pre-condition* compliance). The overall effect is to *restrict* 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 ExtReal = Ok Real | Error
```

and then redefine

```
g :: Real -> ExtReal
g 0 = Error
g n = Ok 1/n
```

In general, one might define parametric type

```
data Ext a = Ok a | Error
```

in order to extend an arbitrary data type  $a$  with its (polymorphic) exception (or error value). Clearly, one has

$$\text{Ext } A \cong \text{Maybe } A \cong 1 + A$$

So, in abstract terms, one may regard as *partial* every function of signature

$$1 + A \xleftarrow{g} B$$

for some  $A$  and  $B$ <sup>1</sup>.

## 4.2 Putting partial functions together

Do partial functions compose? Their types won't match in general:

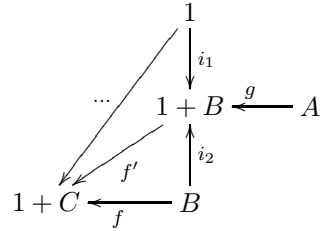
$$\begin{array}{c} 1 + B \xleftarrow{g} A \\ \vdots \\ 1 + C \xleftarrow{f} B \end{array}$$

---

<sup>1</sup>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.



Clearly, we have to extend  $f$  — which is itself a partial function — to some  $f'$  able to accept arguments from  $1 + B$ :



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  $f$ :

$$f \bullet g \stackrel{\text{def}}{=} [i_1, f] \cdot g \tag{4.1}$$

Expressed in terms of `Ext`, composite function  $f \bullet g$  works as follows:

$$(f \bullet g)a = f'(ga)$$

where

$$\begin{aligned} f' \text{Error} &= \text{Error} \\ f' (\text{Ok } b) &= f b \end{aligned}$$

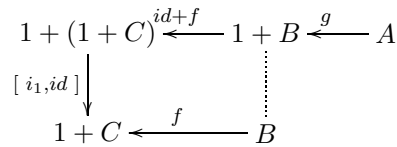
Altogether, we have the following Haskell expression for the meaning of  $f \bullet g$ :

```
\a -> f' (g a)
  where f' Error = Error
        f' (Ok b) = f b
```

Note that the adopted extension of  $f$  can be decomposed — by reverse  $+$ -absorption (2.41) — into

$$f' = [i_1, id] \cdot (id + f)$$

as displayed in diagram



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  $u$  such that

$$f \bullet u = f = u \bullet f \tag{4.2}$$

holds? Clearly, if it exists, it must bear type  $1 + A \xleftarrow{u} A$ . Let us *solve* (4.2) for  $u$ :

$$\begin{aligned} f \bullet u = f = u \bullet f & \\ \equiv \quad \{ \text{substitution} \} & \\ [i_1, f] \cdot u = f = [i_1, u] \cdot f & \\ \Leftarrow \quad \{ \text{let } u = i_2 \} & \\ [i_1, f] \cdot i_2 = f = [i_1, i_2] \cdot f \wedge u = i_2 & \\ \equiv \quad \{ \text{by } +\text{-cancellation (2.38) and } +\text{-reflection (2.39)} \} & \\ f = f = id \cdot f \wedge u = i_2 & \\ \Leftarrow \quad \{ \text{identity} \} & \\ u = i_2 & \end{aligned}$$

### 4.3 Lists

In contrast to partial functions, which can produce no output, let us now consider functions which 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  $C^* \xleftarrow{\text{concat}} (C^*)^*$ , where  $\text{concat} \stackrel{\text{def}}{=} (\llbracket \_ \rrbracket, \text{++})$ . In summary:

$$f \bullet g \stackrel{\text{def}}{=} \text{concat} \cdot f^* \cdot g \quad (4.3)$$

as captured in the following diagram:

$$\begin{array}{ccc} (C^*)^* & \xleftarrow{f^*} & B^* \xleftarrow{g} A \\ \text{concat} \downarrow & & \vdots \\ C^* & \xleftarrow{f} & B \end{array}$$

**Exercise 4.1.** Show that `singl` (recall exercise 3.7) is the unit  $u$  of  $\bullet$  in the context of (4.3).

□

---

**Exercise 4.2.** Encode in HASKELL a pointwise version of (4.3). **Hint:** first apply (list) *cata-absorption* (3.67).

□

---

## 4.4 Monads

Both function composition schemes (4.1) and (4.3) above share the same polytypic pattern: the output of the producer function is “F-times” more elaborate than the input of the consumer function, where  $F$  is some parametric datatype —  $F X = 1 + X$  in case of (4.1), and  $F X = X^*$  in case of (4.3). Then a composition scheme is devised for such functions, which is displayed in

$$\begin{array}{ccc} F(F C) & \xleftarrow{F f} & F B \xleftarrow{g} A \\ \mu \downarrow & & \vdots \\ F C & \xleftarrow{f} & B \end{array}$$

and is given by

$$f \bullet g \stackrel{\text{def}}{=} \mu \cdot F f \cdot g \quad (4.4)$$

where  $F A \xleftarrow{\mu} F^2 A$  is a suitable polymorphic function. Together with a unit function  $F A \xleftarrow{u} A$  and  $\mu$ , datatype  $F$  will form a so-called *monad* type, of which  $1 + \_$  and  $(\_)^*$  are the two examples seen above.

Arrow  $\mu \cdot F f$  is called the *extension* of  $f$ . Functions  $\mu$  and  $u$  are referred to as the monad's *multiplication* and *unit*, respectively. The monadic composition scheme (4.4) is called *Kleisli composition*.

A *monadic arrow*  $F B \xleftarrow{f} A$  conveys the idea of a function which produces an output of "type"  $B$  "wrapped by  $F$ ", where datatype  $F$  describes some kind of (computational) "effect". The monad's unit  $F B \xleftarrow{u} B$  is a primitive monadic arrow which produces (*i.e.* promotes, injects, wraps) data *together with* 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 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  $F$  to be regarded as a monad:

**Multiplication :**

$$\begin{array}{ccc} F^2 A & \xleftarrow{\mu} & F^3 A \\ \mu \downarrow & & \downarrow F \mu \\ F A & \xleftarrow{\mu} & F^2 A \end{array} \quad \mu \cdot \mu = \mu \cdot F \mu \quad (4.5)$$

**Unit :**

$$\begin{array}{ccc} F^2 A & \xleftarrow{u} & F A \\ \mu \downarrow & \swarrow id & \downarrow F u \\ F A & \xleftarrow{\mu} & F^2 A \end{array} \quad \mu \cdot u = \mu \cdot F u = id \quad (4.6)$$

Simple but beautiful symmetries apparent in these diagrams make it easy to memorize their laws and check them for particular cases. For instance, for the  $(1 + \_)$  monad, law (4.6) 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 (4.5) and (4.6), the different instances of  $\mu$  and  $u$  are differently typed, as these are polymorphic and exhibit natural properties:

**$\mu$ -natural :**

$$\begin{array}{ccc}
 A & F A \xleftarrow{\mu} F^2 A & \\
 f \downarrow & F f \downarrow \quad \downarrow F^2 f & \\
 B & F B \xleftarrow{\mu} F^2 B & 
 \end{array}
 \qquad
 F f \cdot \mu = \mu \cdot F^2 f
 \qquad
 (4.7)$$

**$u$ -natural :**

$$\begin{array}{ccc}
 A & F A \xleftarrow{u} A & \\
 f \downarrow & F f \downarrow \quad \downarrow f & \\
 B & F B \xleftarrow{u} B & 
 \end{array}
 \qquad
 F f \cdot u = u \cdot f
 \qquad
 (4.8)$$

The simplest of all monads is the *identity monad*  $F X \stackrel{\text{def}}{=} X$ , which is such that  $\mu = id$ ,  $u = id$  and  $f \bullet g = f \cdot g$ . So — in a sense — one may think of all the functional discipline studied so far as a particular case of a wider discipline in which an arbitrary monad is present.

#### 4.4.1 Properties involving (Kleisli) composition

The following properties arise from the definitions and monadic properties presented above:

$$f \bullet (g \bullet h) = (f \bullet g) \bullet h \qquad (4.9)$$

$$u \bullet f = f = f \bullet u \qquad (4.10)$$

$$(f \bullet g) \cdot h = f \bullet (g \cdot h) \qquad (4.11)$$

$$(f \cdot g) \bullet h = f \bullet (F g \cdot h) \qquad (4.12)$$

$$id \bullet id = \mu \qquad (4.13)$$

Properties (4.9) and (4.10) 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 (4.10) holds for the list monad, recall (4.2). A general proof can be produced similarly. We select property (4.9) as an illustration of the

rôle of the monadic properties:

$$\begin{aligned}
& f \bullet (g \bullet h) \\
= & \quad \{ \text{definition (4.4) twice} \} \\
& \mu \cdot F f \cdot (\mu \cdot F g \cdot h) \\
= & \quad \{ \mu \text{ is natural (4.7)} \} \\
& \mu \cdot \mu \cdot F(F f) \cdot F g \cdot h \\
= & \quad \{ \text{functor } F \} \\
& \mu \cdot \mu \cdot F(F f \cdot g) \cdot h \\
= & \quad \{ \text{definition (4.4)} \} \\
& \mu \cdot (F f \cdot g) \bullet h \\
= & \quad \{ \text{definition (4.4)} \} \\
& (f \bullet g) \bullet h
\end{aligned}$$

**Exercise 4.3.** Check the other laws above.

□

## 4.5 Monadic application (binding)

The monadic extension of functional application  $ap$  (2.67) is another operator  $ap'$  which is intended to be “tolerant” in face of any  $F$ 'ed argument  $x$ :

$$\begin{aligned}
(F B)^A \times F A & \xrightarrow{ap'} B \\
ap'(f, x) & = f' x = (\mu \cdot F f)x
\end{aligned} \tag{4.14}$$

If in curry/flipped format, monadic application is called *binding* and denoted by symbol “ $\gg=$ ”, looking very much like postfix functional application,

$$((F B)^A)^{F A} \xrightarrow{\gg=} F B \tag{4.15}$$

that is:

$$x \gg= f \stackrel{\text{def}}{=} (\mu \cdot F f)x \tag{4.16}$$

This operator will exhibit properties arising from its definition and the basic monadic properties, *e.g.*

$$\begin{aligned}
 x \gg= u & \\
 & \equiv \quad \{ \text{definition (4.16)} \} \\
 & \quad (\mu \cdot \mathbf{F} u)x \\
 & \equiv \quad \{ \text{law (4.6)} \} \\
 & \quad (id)x \\
 & \equiv \quad \{ \text{identity function} \} \\
 & \quad x
 \end{aligned}$$

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  $A \xrightarrow{f_1} \mathbf{F} B_1$ ,  $B_1 \xrightarrow{f_2} \mathbf{F} B_2$ ,  $\dots$ ,  $B_{n-1} \xrightarrow{f_n} \mathbf{F} B_n$ .

## 4.6 Sequencing and the `do`-notation

Given two monadic values  $x$  and  $y$ , it becomes possible to “sequence” them, thus obtaining another of such value, by defining the following operator:

$$x \gg y \stackrel{\text{def}}{=} x \gg= \underline{y} \tag{4.17}$$

For instance, within the finite-list monad, one has

$$[1, 2] \gg [3, 4] = (\text{concat} \cdot \underline{[3, 4]^*})[1, 2] = \text{concat}[[3, 4], [3, 4]] = [3, 4, 3, 4]$$

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:

$$\text{do } x_1; x_2; \dots; x_n \stackrel{\text{def}}{=} x_1 \gg \text{do } x_2; \dots; x_n$$

for  $n \geq 1$ . For  $n = 1$  one trivially has

$$\text{do } x_1 \stackrel{\text{def}}{=} x_1$$

## 4.7 Generators and comprehensions

The  $\text{do}$ -notation accepts a variant in which the arguments of the  $\gg$  operator are “generators” of the form

$$a \leftarrow x \tag{4.18}$$

where, for  $a$  of type  $A$ ,  $x$  is an inhabitant of monadic type  $F A$ . One may regard  $a \leftarrow x$  as meaning “let  $a$  be taken from  $x$ ”. Then the  $\text{do}$ -notation extends as follows:

$$\text{do } a \leftarrow x_1; x_2; \dots; x_n \stackrel{\text{def}}{=} x_1 \gg \lambda a. (\text{do } x_2; \dots; x_n) \tag{4.19}$$

Of course, we should now allow for the  $x_i$  to range over terms involving variable  $a$ . For instance (again in the list-monad), by writing

$$\text{do } a \leftarrow [1, 2, 3]; [a^2] \tag{4.20}$$

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,

$$[a^2 \mid a \leftarrow [1, 2, 3]] \tag{4.21}$$

as a “shorthand” of (4.20). This is an instance of the so-called *comprehension* notation, which can be defined in general as follows:

$$[e \mid a_1 \leftarrow x_1, \dots, a_n \leftarrow x_n] = \text{do } a_1 \leftarrow x_1; \dots; a_n \leftarrow x_n; u(e) \tag{4.22}$$

Alternatively, comprehensions can be defined as follows, where  $p, q$  stand for arbitrary generators:

$$[t] = u t \tag{4.23}$$

$$[f x \mid x \leftarrow l] = (F f) l \tag{4.24}$$

$$[t \mid p, q] = \mu [ [t \mid q] \mid p ] \tag{4.25}$$

Note, however, that comprehensions are not restricted to lists or sets — they can be defined for any monad  $F$ .



## 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

$$p \gg q = p \gg= \_ \rightarrow q$$

as expected. This class is instantiated for finite sequences (`[]`), `Maybe` and `IO`.

The  $\mu$  multiplication operator is function `join` in module `Monad.hs`:

```
join :: (Monad m) => m (m a) -> m a
join x = x >>= id
```

This is easily justified:

$$\begin{aligned}
 \text{join } x &= x \gg= id && (4.26) \\
 &= \{ \text{definition (4.16)} \} \\
 &\quad (\mu \cdot F id)x \\
 &= \{ \text{functors commute with identity (3.44)} \} \\
 &\quad (\mu \cdot id)x \\
 &= \{ \text{law (2.10)} \} \\
 &\quad \mu x
 \end{aligned}$$

In `Mpi.hs` we define (Kleisli) monadic composition in terms of the binding operator:

```
(.!) :: Monad a => (b -> a c) -> (d -> a b) -> d -> a c
(f .! g) a = (g a) >>= f
```

### 4.8.1 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, `getLine :: IO String` is a particular 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 (4.18). 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
  (>>=) = primbindIO
  return = primretIO
```

Rather interesting is the way IO is regarded as a functor:

$$\text{fmap } f \ x = x \gg= (\text{return} \ . \ f)$$

This goes the other way round, the monadic structure “helping” in defining the functor structure, everything consistent with the underlying theory:

$$\begin{aligned} x \gg= (u \cdot f) &= (\mu \cdot \text{IO}(u \cdot f))x \\ &= \quad \{ \text{functors commute with composition} \} \\ &\quad (\mu \cdot \text{IO } u \cdot \text{IO } f)x \\ &= \quad \{ \text{law (4.6) for } F = \text{IO} \} \\ &\quad (\text{IO } f)x \\ &= \quad \{ \text{definition of } \text{fmap} \} \\ &\quad (\text{fmap } f)x \end{aligned}$$

For enjoyable reading on monadic input/output in HASKELL see [Hud00], chapter 18.

**Exercise 4.4.** Use the `do`-notation and the comprehension notation to output the following truth-table, in HASKELL:

$p / q$	<i>False</i>	<i>True</i>
<i>False</i>	<i>False</i>	<i>False</i>
<i>True</i>	<i>False</i>	<i>True</i>

□

---

**Exercise 4.5.** 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?

□

---

## 4.9 The state monad

NB: this section is still very drafty

The so-called *state monad* is a monad whose inhabitants are state-transitions encoding a particular brand of state-based automaton known as *Mealy machine*. 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

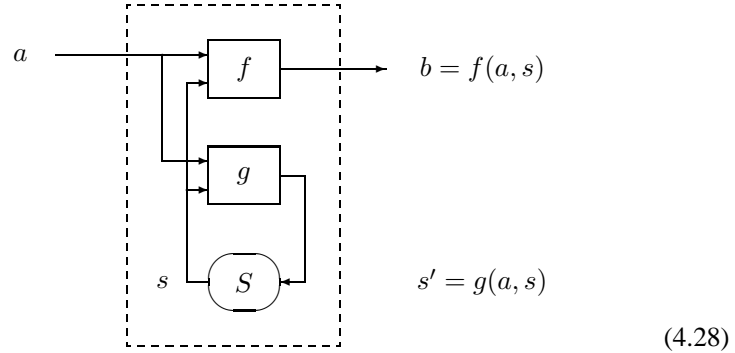
$$A \times S \xrightarrow{\delta} B \times S \quad (4.27)$$

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.

It is clear from (4.27) that  $\delta$  can be expressed as the *split* of two functions  $f$  and  $g$ ,  $\delta = \langle f, g \rangle$ , as depicted in the following diagram:



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 handled explicitly (as is the case of eg. the data kept in a large database). So, it is convenient to *abstract* from it. Such an abstraction leads to the *state monad* in the following way: recalling (2.75), we simply *curry*  $\delta$

$$A \xrightarrow{\bar{\delta}} \underbrace{(B \times S)^S}_{(\text{St } S) B} \quad (4.29)$$

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 which *appends* a new element to the back 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:

$$\begin{aligned} \text{enqueue} & : A \rightarrow (1 \times S)^S \\ \text{enqueue } a & \stackrel{\text{def}}{=} \langle !, (+[a]) \rangle \end{aligned} \quad (4.30)$$

Action *enqueue* performs *enq* on the state while acknowledging it by issuing an output of type 1.

**Unit and multiplication.** Let us show that

$$(\text{St } S) A \cong (A \times S)^S \quad (4.31)$$

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 2.14 is recommended.) Notation  $\widehat{f}$  will be used to abbreviate *uncurry*  $f$ . Thus the following variant of universal law (2.67),

$$\widehat{k} = f \Leftrightarrow f = ap \cdot (k \times id) \quad (4.32)$$

whose cancellation

$$\widehat{k} = ap \cdot (k \times id) \quad (4.33)$$

is written pointwise as follows:

$$\widehat{k}(c, a) = (k \ c)a \quad (4.34)$$

First of all, what is the functor behind (4.31)? Fixing the state space  $S$ , we obtain

$$FX \stackrel{\text{def}}{=} (X \times S)^S \quad (4.35)$$

on objects and

$$Ff \stackrel{\text{def}}{=} (f \times id)^S \quad (4.36)$$

on functions, where  $(-)^S$  is the exponential functor (2.71).

The unit of this monad is the transpose of the simplest of all Mealy machines — the identity:

$$\begin{aligned} u &: A \rightarrow (A \times S)^S \\ u &= \overline{id} \end{aligned} \quad (4.37)$$

Let us see what this means:

$$\begin{aligned} &u = \overline{id} \\ \equiv &\{ (2.67) \} \\ &ap \cdot (u \times id) = id \\ \equiv &\{ \text{introducing variables} \} \\ &ap(u \ a, s) = (a, s) \\ \equiv &\{ \text{definition of } a \} \\ &(u \ a)s = (a, s) \end{aligned}$$

From the type of  $\mu$ , 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 (2.71)) for  $x$  of type  $((A \times S)^S \times S) \xrightarrow{x} (A \times S)$ . This, on its turn, is easily recognized as an instance of the  $ap$  polymorphic function (2.67), which is such that  $ap = \widehat{id}$ , recall (2.69). Altogether, we define

$$\mu = ap^S \tag{4.38}$$

Let us check the meaning of  $\mu$  by applying it to an action expressed as in diagram (2.75):

$$\begin{aligned} & \mu\langle f, g \rangle = ap^S\langle f, g \rangle \\ \equiv & \quad \{ (2.71) \} \\ & \mu\langle f, g \rangle = ap \cdot \langle f, g \rangle \\ \equiv & \quad \{ \text{extensional equality (2.5)} \} \\ & \mu\langle f, g \rangle s = ap(f\ s, g\ s) \\ \equiv & \quad \{ \text{definition of } ap \} \\ & \mu\langle f, g \rangle s = (f\ s)(g\ s) \end{aligned}$$

We find out that  $\mu$  “unnests” the action inside  $f$  by applying it to the state delivered by  $g$ .

**Checking the monadic laws.** The calculation of (4.6) is made in two parts, checking  $\mu \cdot u = id$  first,

$$\begin{aligned} & \mu \cdot u \\ = & \quad \{ \text{definitions} \} \\ & ap^S \cdot \overline{id} \\ = & \quad \{ \text{exponentials absorption (2.72)} \} \\ & \overline{ap \cdot id} \\ = & \quad \{ \text{reflection (2.69)} \} \\ & id \end{aligned}$$

and then checking  $\mu \cdot (Fu) = id$ :

$$\begin{aligned}
 & \mu \cdot (Fu) \\
 = & \quad \{ (4.38, 4.36) \} \\
 & ap^S \cdot (\overline{id} \times id)^S \\
 = & \quad \{ \text{functor} \} \\
 & (ap \cdot (\overline{id} \times id))^S \\
 = & \quad \{ \text{cancellation (2.68)} \} \\
 & id^S \\
 = & \quad \{ \text{functor} \} \\
 & id
 \end{aligned}$$

The proof of (4.5) is also not difficult once supported by the laws of exponentials.

**Kleisli composition.** Let us calculate  $f \bullet g$  for this monad:

$$\begin{aligned}
 & f \bullet g \\
 = & \quad \{ (4.4) \} \\
 & \mu \cdot F f \cdot g \\
 = & \quad \{ (4.38); (4.36) \} \\
 & ap^S \cdot (f \times id)^S \cdot g \\
 = & \quad \{ (-)^S \text{ is a functor} \} \\
 & (ap \cdot (f \times id))^S \cdot g \\
 = & \quad \{ (4.32) \} \\
 & \widehat{f}^S \cdot g \\
 = & \quad \{ \text{cancellation} \} \\
 & \widehat{f}^S \cdot \overline{\widehat{g}} \\
 = & \quad \{ \text{absorption (2.72)} \} \\
 & \overline{\widehat{f} \cdot \widehat{g}}
 \end{aligned}$$

In summary, we have:

$$f \bullet g = \overline{\widehat{f} \cdot \widehat{g}} \quad (4.39)$$

Let us use this in calculating law

$$pop \bullet push = u \quad (4.40)$$

where  $push$  and  $pop$  are such that

$$\begin{aligned} push & : A \rightarrow (1 \times S)^S \\ \widehat{push} & \stackrel{\text{def}}{=} \langle !, \hat{\cdot} \rangle \end{aligned} \quad (4.41)$$

$$\begin{aligned} pop & : 1 \rightarrow (A \times S)^S \\ \widehat{pop} & \stackrel{\text{def}}{=} \langle head, tail \rangle \cdot \pi_2 \end{aligned} \quad (4.42)$$

for  $S$  the datatype of finite lists. We reason:

$$\begin{aligned} & pop \bullet push \\ = & \{ (4.39) \} \\ & \overline{\widehat{pop} \cdot \widehat{push}} \\ = & \{ (4.41, 4.42) \} \\ & \overline{\langle head, tail \rangle \cdot \pi_2 \cdot \langle !, \hat{\cdot} \rangle} \\ = & \{ (2.20, 2.24) \} \\ & \overline{\langle head, tail \rangle \cdot \hat{\cdot}} \\ = & \{ out \cdot in = id(\text{lists}) \} \\ & \overline{id} \\ = & \{ (4.37) \} \\ & u \end{aligned}$$

**Bind.** The effect of binding a state transition  $x$  to a state-monadic function  $h$  is calculated in a similar way:

$$\begin{aligned} & x \gg= h \\ = & \{ (4.16) \} \\ & (\mu \cdot Fh)x \\ = & \{ (4.38) \text{ and } (4.36) \} \\ & (ap^S \cdot (h \times id)^S)x \\ = & \{ (-)^S \text{ is a functor} \} \end{aligned}$$



$$\begin{aligned}
& (ap \cdot (h \times id))^S x \\
= & \{ \text{cancellation (4.33)} \} \\
& \widehat{h}^S x \\
= & \{ \text{exponential functor (2.71)} \} \\
& \widehat{h} \cdot x
\end{aligned}$$

Let us unfold  $\widehat{h} \cdot x$  by splitting  $x$  into its components two components  $f$  and  $g$ :

$$\begin{aligned}
& \langle f, g \rangle \gg= h = \widehat{h} \cdot \langle f, g \rangle \\
\equiv & \{ \text{go pointwise} \} \\
& (\langle f, g \rangle \gg= h)_s = \widehat{h}(\langle f, g \rangle_s) \\
\equiv & \{ (2.18) \} \\
& (\langle f, g \rangle \gg= h)_s = \widehat{h}(f\ s, g\ s) \\
\equiv & \{ (4.34) \} \\
& (\langle f, g \rangle \gg= h)_s = h(f\ s)(g\ s)
\end{aligned}$$

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:

$$get \stackrel{\text{def}}{=} \langle id, id \rangle \quad (4.43)$$

$$put \stackrel{\text{def}}{=} \overline{\langle !, \pi_1 \rangle} \quad (4.44)$$

Action  $g$  retrieves the data stored in the state while  $put$  (which can also be written

$$put\ s = modify(\underline{s}) \quad (4.45)$$

where

$$modify\ f \stackrel{\text{def}}{=} \langle !, f \rangle \quad (4.46)$$

updates the state via state-to-state function  $f$ ) stores a particular value in the state.

The following is an example, in Haskell, of the standard use of *get/put* in managing context data, in this case a counter. The function decorates each node of a *BTree* (recall this datatype from page 91) with its position in the tree:

```
decBTree Empty = return Empty
decBTree (Node (a,(t1,t2))) =
  do n <- get ;
    put(n+1) ;
    l <- decBTree t1 ;
    r <- decBTree t2 ;
    return (Node((a,n),(l,r)))
```

## 4.10 Bibliography notes

The use of monads in computer science started with Moggi [Mog89], who had the idea that monads should supply the extra semantic information needed to implement the lambda-calculus theory. Haskell [Jon03] is among the computer languages which make systematic use of monads for implementing effects and imperative constructs in an otherwise purely functional language.

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 [Wad89] made a great contribution in 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 [HM93], component-oriented programming [Bar01], strategic programming [LV03] and multimedia [Hud00].

**Part II**

**Moving Away From (Pure)  
Functions**



# Bibliography

- [ABH<sup>+</sup>92] 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](http://www.cs.nott.ac.uk/~rcb).
- [Bac78] 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.
- [Bac04] R.C. Backhouse. *Mathematics of Program Construction*. Univ. of Nottingham, 2004. Draft of book in preparation. 608 pages.
- [Bar01] L.S. Barbosa. *Components as Coalgebras*. University of Minho, December 2001. Ph. D. thesis.
- [BD77] R.M. Burstall and J. Darlington. A transformation system for developing recursive programs. *JACM*, 24(1):44–67, January 1977.
- [BdM97] R. Bird and O. de Moor. *Algebra of Programming*. Series in Computer Science. Prentice-Hall International, 1997. C.A.R. Hoare, series editor.
- [Bir98] R. Bird. *Introduction to Functional Programming*. Series in Computer Science. Prentice-Hall International, 2nd edition, 1998. C.A.R. Hoare, series editor.
- [Flo67] 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.
- [Fok92] M.M. Fokkinga. *Law and Order in Algorithmics*. PhD thesis, University of Twente, Dept INF, Enschede, The Netherlands, 1992.
- [GHA01] Jeremy Gibbons, Graham Hutton, and Thorsten Altenkirch. When is a function a fold or an unfold?, 2001. WGP, July 2001 (slides).

- [Gof84] J. Le Goff. *Calendário*, volume I, chapter 8, pages 260–292. I.N.-C.M., 1984. EINAUDI Encyclopedia (Portuguese translation).
- [HM93] Graham Hutton and Erik Meijer. Monadic parsing in Haskell. *Journal of Functional Programming*, 8(4), 1993.
- [Hud00] P. Hudak. *The Haskell School of Expression - Learning Functional Programming Through Multimedia*. Cambridge University Press, 1st edition, 2000. ISBN 0-521-64408-9.
- [JJ96] J. Jeuring and P. Jansson. Polytypic programming. In *Advanced Functional Programming*, number 1129 in LNCS, pages 68–114. Springer, 1996.
- [JJ98] P. Jansson and J. Jeuring. Polylib — a library of polytypic functions. In *Workshop on Generic Programming (WGP'98)*, Marstrand, Sweden, 1998.
- [Jon03] 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.
- [LV03] 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.
- [MA86] 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.
- [Mal90] G. Malcolm. Data structures and program transformation. *Science of Computer Programming*, 14:255–279, 1990.
- [McC63] J. McCarthy. Towards a mathematical science of computation. In C.M. Poplewell, editor, *Proc. of IFIP 62*, pages 21–28, Amsterdam-London, 1963. North-Holland Pub. Company.
- [MH95] 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.
- [Mog89] Eugenio 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.

- [NR69] 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.
- [Oli08a] J.N. Oliveira. Extended static checking by calculation using the pointfree transform, 2008. Tutorial paper (56 p.) accepted for publication by Springer-Verlag, LNCS series.
- [Oli08b] J.N. Oliveira. *Transforming Data by Calculation*. In *GTTSE'07*, volume 5235 of *LNCS*, pages 134–195. Springer, 2008.
- [OR06] J.N. Oliveira and C.J. Rodrigues. Pointfree factorization of operation refinement. In *FM'06*, volume 4085 of *LNCS*, pages 236–251. Springer-Verlag, 2006.
- [PH70] M.S. Paterson and C.E. Hewitt. Comparative schematology. In *Project MAC Conference on Concurrent Systems and Parallel Computation*, pages 119–127, August 1970.
- [Wad89] P.L. Wadler. Theorems for free! In *4th International Symposium on Functional Programming Languages and Computer Architecture*, pages 347–359, London, Sep. 1989. ACM.