

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

| | | |
|----------|--|-----------|
| 1 | Introduction | 3 |
| 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 |
| | | |
| I | Calculating with Functions | 9 |
| | | |
| 2 | An Introduction to Pointfree Programming | 11 |
| 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 |
| 3 | Recursion in the Pointfree Style | 55 |
| 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 |
| 4 | Why Monads Matter | 105 |
| 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 |

| | | |
|------------|---|------------|
| II | Moving Away From (Pure) Functions | 125 |
| 5 | Quasi-inductive datatypes | 127 |
| 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 |
| III | Calculating with Relations | 145 |
| 6 | Introduction to Relational Calculation | 147 |
| 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 |
| 7 | An Introduction to Relational Hyломorphisms | 197 |
| 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 | μ -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 |
| 8 | Theorems for Free | 225 |
| 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 |
|----------------------------------|-----|

IV Data Refinement by Calculation 233

9 On Data Representation 235

| | |
|--|-----|
| 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 |

10 On Algorithmic Refinement 265

| | |
|--|-----|
| 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/for</code> loops | 277 |
| 10.13 | Bibliography notes | 277 |
| A | Haskell Support Library in Haskell | 279 |
| A.0.1 | Set.hs | 279 |
| B | Solutions to Selected Exercises | 281 |

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 "**formulae and equations** (...) which share the **elegance** of those which underlie **physics and chemistry** or any other branch of basic science".*

The formulae 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⁺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*).

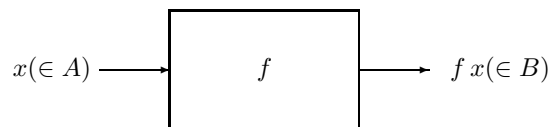
Many technical issues are deliberately ignored and deferred to later chapters. Most programming examples will be provided in the HASKELL functional programming language. Appendix A includes the listings of some HASKELL modules which complement the HUGS *Standard Prelude* (which is based very closely 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,

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

expresses the computation rule of the *successor* function — the function 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 succ and which, by the way, can be shared by other functions, *e.g.* $\text{sq } n \stackrel{\text{def}}{=} n^2$.

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

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

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 $12.cu.\tilde{p}.18.ce.\tilde{p}.27.co.\tilde{p}.17$ for what is nowadays written as $12x^3 + 18x^2 + 27x + 17$. We may find such *syncopated* notation odd, but 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.

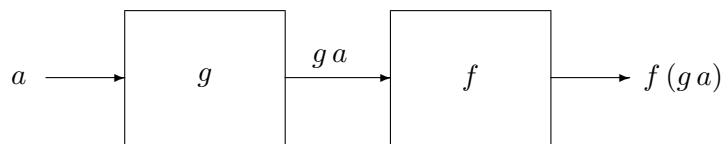
Concerning application, the pointfree style replaces it by the more generic concept of functional *composition* suggested by function-arrow chaining: wherever two functions are such that the target type of one of them, say $B \xleftarrow{g} A$ is the same as the source type of the other, say $C \xleftarrow{f}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

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

or by block-diagram



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

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

Composition is certainly the most basic of all functional combinators. It is the first kind of “glue” which comes to mind when programmers need to combine, or chain functions (or processes) to obtain more elaborate functions (or processes) ¹. 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 $f_1 \cdot \dots \cdot f_n$.

¹It even has a place in script languages such as UNIX’s, where $\mathbf{f} \mid \mathbf{g}$ 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):

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

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 `f . 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

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

(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° 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,

$$\text{Weekday} = \{ \textit{Sunday}, \textit{Monday}, \textit{Tuesday}, \textit{Wednesday}, \textit{Thursday}, \textit{Friday}, \textit{Saturday} \}$$

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

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

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

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

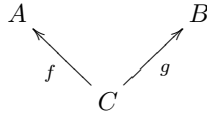
Finally, note that all classes of functions referred to so far — constants, identities, epics, monics and isos — are closed under composition, that is, the composition of two constants is a constant, the composition of two epics is epic, *etc.*

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



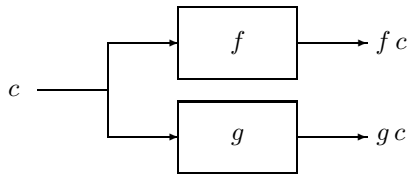
Because f and g share the same domain, their outputs can be paired, that is, we may write ordered pair $(f\ c, g\ c)$ for each $c \in C$. Such pairs belong to the Cartesian product of A and B , that is, to the set

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

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:

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

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* π_1 or π_2 ,

$$A \xleftarrow{\pi_1} A \times B \xrightarrow{\pi_2} B \tag{2.19}$$

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:

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

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:

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

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 π_1 and π_2 of $C \times D$:

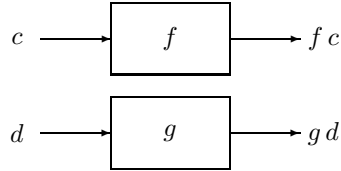
$$\begin{array}{ccccc}
 A & \xleftarrow{\pi_1} & A \times B & \xrightarrow{\pi_2} & B \\
 \uparrow f & & & & \uparrow g \\
 C & \xleftarrow{\pi_1} & C \times D & \xrightarrow{\pi_2} & D
 \end{array}$$

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

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

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

$$\frac{\begin{array}{c} A \xleftarrow{f} C \\ B \xleftarrow{g} D \end{array}}{A \times B \xleftarrow{f \times g} C \times D} \tag{2.23}$$

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*:

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

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:

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

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 π_1 and π_2 only. These exhibit interesting properties, for instance $\langle \pi_1, \pi_2 \rangle = id$. This property is known as \times -*reflexion* and is depicted as follows:

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

2.8. GLUING FUNCTIONS WHICH DO NOT COMPOSE — PRODUCTS 25

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

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

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

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

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 π_1 and `snd` to denote π_2 . In the C programming language this datatype is called the “struct datatype”,

```

struct {
    A first;
    B second;
};

```

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

```

record
  first: A;
  second: B
end;

```

Isomorphism (2.31) 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.

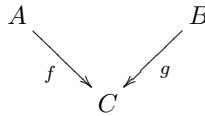
□

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 $f : C \leftarrow A$ and $g : C \leftarrow B$ which however share the same codomain is depicted in



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

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

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

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*:

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

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 :

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

+reflexion :

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

+fusion :

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

+absorption :

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

+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:

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

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 *either*s ($A + B$). 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 *either*s.

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

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

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

$$\begin{array}{ccc}
 ? & \xrightarrow{h+(i,j)} & ? \\
 \swarrow x+y & & \downarrow id+id \times f \\
 & & ?
 \end{array}$$

□

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 ², 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 $B \xleftarrow{g} C$, there exists an arrow $A \times B \xleftarrow{\langle f, g \rangle} C$ such that

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

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

²In 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

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

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

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

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

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

if ($p\ x$) *then* ($g\ x$) *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 $B \xleftarrow{h} A$. Bool is the primitive datatype containing truth values FALSE and 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:

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

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 `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*

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

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 $\text{Bool} \xleftarrow{=c} A$ 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 — EXPONENTIALS39

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 :

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

Exponentials reflexion :

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

Exponentials fusion :

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

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:

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

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

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

2.14. GLUING FUNCTIONS WHICH DO NOT COMPOSE — EXPONENTIALS41

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:

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

Now we are ready to understand the laws which follow:

Exponentials absorption :

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

Exponentials-functor :

$$(g \cdot h)^A = g^A \cdot h^A \quad (2.73)$$

Exponentials-functor-id :

$$id^A = id \quad (2.74)$$

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

$$\text{curry } f \stackrel{\text{def}}{=} \overline{f} \quad (2.75)$$

which maps a function f to its transpose \overline{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 $\widehat{}$ -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 \bar{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 $\hat{}$ 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 \mathbb{R} , `Bool` and \mathbb{N} and, in particular, the fact that we can associate to each natural number n its *initial segment* $n = \{1, 2, \dots, n\}$. We extend this to \mathbb{N}_0 by stating $0 = \{\}$ and, for $n > 0$, $n + 1 = \{n + 1\} \cup n$.

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

Datatype 0

Among such enumerated types, 0 is the smallest because it is empty. This is the `Void` datatype in HASKELL, which has no constructor at all. Datatype 0 (which we tend to write simply as 0) may not seem very “useful” in practice but it is of theoretical interest. For instance, it is easy to check that the following “obvious” properties hold:

$$A + 0 \cong A \tag{2.79}$$

$$A \times 0 \cong 0 \tag{2.80}$$

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:

$$A \times 1 \cong A \tag{2.81}$$

$$A^0 \cong 1 \tag{2.82}$$

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 π_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* $[b_0, g]$ 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: $f = [b_1, b_2]$. Now suppose that $B = \{b_1, b_2\}$ (for $b_1 \neq b_2$). Then $1 + 1 \cong B$ will hold, for whatever choice of inhabitants b_1 and b_2 . So we are in a situation similar to 1: we will use symbol `2` to represent the abstract class of all such B s containing exactly two elements. Therefore, we can write:

$$1 + 1 \cong 2$$

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

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

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

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 *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 :

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

Initial datatype fusion :

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

Terminal datatype reflexion :

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

Terminal datatype fusion :

$$\begin{array}{c}
 1 \\
 \uparrow \!_A \\
 A \xleftarrow{f} B
 \end{array}
 \begin{array}{l}
 \nearrow \!_B \\
 \end{array}
 \qquad \!_A \cdot f = \!_B \qquad (2.92)$$

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:

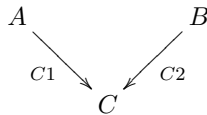
```
Main> :i CostumerId
-- type constructor
data CostumerId

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

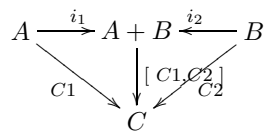
In general, let A and B be two known datatypes. Via declaration

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

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 :

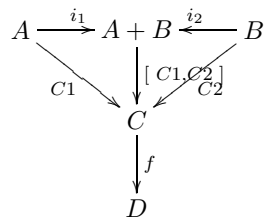


This diagram leads to an obvious instance of coproduct diagram (2.36),



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 $[C1, C2]$:

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

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 π_1, π_2 “is” such a datatype, one may be interested in a suitably named replica of $A \times A$ in which points are built explicitly by some constructor (say $Point$) and observed by dedicated selectors (say x and y):

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

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

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

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 $distl$ (read: ‘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* trend³. 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

³Even 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.