

Chapter 4

Data representation and refinement

Isomorphism provides us with a notion of data-equivalence. How do we “order” data-structures according to their “size” or “expressive power”? Are there “better” or “worse” data-structures for a given purpose? In this chapter — which is currently in a very drafty version — we address the field of *formal* data processing and refinement by transformation.

4.1 Introduction

The definition of a function $B \xleftarrow{f} A$ can be regarded as a kind of “contract”: function f is committed to produce an A -value provided it is supplied with a B -value. Such a “functional contract” can be of two kinds: (a) f intentionally loses information, because B is found too detailed and one wants to capture only the A -aspect of B -values — so, A is an *abstraction* of B (f is non-injective); (b) f faithfully converts data from the B -format to the A -format — so, f is injective and, in the limit, the two formats are the same (f is the identity).

Case (a) above is perhaps more interesting than (b) and 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. In this way, it is close to *data-mining* which is nothing but a series of clever, highly scrutinized formal reductions of a complex input structure.

Case (a) also includes the process of *transforming* data from one format to another so that it fits in a particular context or matches a given requirement. One of the aspects of this is *data refinement*, which usually means a change of *data representation*: a problem is abstractly characterized in terms of clever, mathematically nice data models which either don’t fit in traditional data structuring techniques, or they do but at a severe cost of efficiency.

Refinement is the process of converting (with a variable concern for formality)

abstract software specifications into real applications intended to run on some available hardware. Formal refinement is one of the most relevant branches of formal methods. It splits into *data refinement* and *algorithmic refinement*, depending on where refinement decisions take place: changing the format of *real* data in the former case, or replacing a *virtual* data-structure¹ by a more efficient (though equally virtual) one in the latter.

Formal refinement is not yet widely used due to the difficulty of scaling its mathematical reasoning to real size problems. By contrast, informal refinement enforces so little discipline in software design that results are very often catastrophic, particularly in the case of very large applications which have a long life-span. As a rule, such applications lack in documentation, have bugs and raise serious maintenance problems. Using or selling them becomes a risky business.

In this chapter we concentrate on data representation and refinement because, in formal software design, data calculation precedes and paves the way to correct algorithm design by calculation.

4.2 Isomorphism

Isomorfism

$$A \begin{array}{c} \xrightarrow{f^{-1}} \\ \cong \\ \xleftarrow{f} \end{array} B \quad \text{that is} \quad \begin{cases} f \cdot f^{-1} = id_A \\ f^{-1} \cdot f = id_B \end{cases} \quad (4.1)$$

means that f^{-1} is both a *right-inverse* and a *left-inverse* of f . Further to many isomorphism laws studied in previous chapters, we are now interested in exploiting isomorphisms involving quasi-inductive datatypes:

- Powersets vs predicates: recall (3.28).
- Partial mappings vs total functions (for a finite A):

$$A \multimap B \begin{array}{c} \xrightarrow{tot} \\ \cong \\ \xleftarrow{untot} \end{array} (B + 1)^A \quad (4.2)$$

where

$$\begin{aligned} tot \sigma &\stackrel{\text{def}}{=} \lambda a. i f a \in \text{dom } \sigma \text{ then } i_1(\sigma a) \text{ else } i_2 \text{ NIL} \\ untot f &\stackrel{\text{def}}{=} \{a \mapsto b \mid a \in A \wedge f a = (i_1 b)\} \end{aligned}$$

- A consequence of the above,

$$\mathcal{P}A \begin{array}{c} \xrightarrow{set2fm} \\ \cong \\ \xleftarrow{dom} \end{array} A \multimap 1$$

¹Recall section 2.5.

cf.

$$\begin{aligned}
 & A \multimap 1 \\
 \cong & \quad \{ \text{converting to total maps (4.2)} \} \\
 & (1 + 1)^A \\
 \cong & \quad \{ \text{basic} \} \\
 & 2^A \\
 \cong & \quad \{ \text{by (3.28)} \} \\
 & \mathcal{P}A
 \end{aligned}$$

- Another example, extending *either* (1.35) to partial maps:

$$(B + C) \multimap A \begin{array}{c} \xrightarrow{\text{unpeither}} \\ \cong \\ \xleftarrow{\text{peither}} \end{array} (B \multimap A) \times (C \multimap A) \quad (4.3)$$

Why an isomorphism? Check:

$$\begin{aligned}
 & (B + C) \multimap A \\
 \cong & \quad \{ \text{converting to total maps} \} \\
 & (A + 1)^{(B+C)} \\
 \cong & \quad \{ \text{exponentials (1.77)} \} \\
 & (A + 1)^B \times (A + 1)^C \\
 \cong & \quad \{ \text{back to partial maps} \} \\
 & (B \multimap A) \times (C \multimap A)
 \end{aligned}$$

Abstraction and representation functions:

$$\text{peither} \stackrel{\text{def}}{=} \cup \cdot ((i_1 \multimap A) \times (i_2 \multimap A))$$

and, expressed in CAMILA,

```

; unpeither : ( ( A + B ) -> C ) --> ( A -> C ) x ( B -
> C )
; The inverse of (unary) unpeither
unpeither(f) = let (a=[ p2(x) -> f[x] | x <- dom(f): p1(x) == 1 ],
                    b=[ p2(x) -> f[x] | x <- dom(f): p1(x) == 2 ])
in < a , b >;

```

4.3 Representation

Representation of A in B means that there exists $A \xleftarrow{f} B$ which has a *right-inverse* r but no *left-inverse*. So, this means that B contains “more information” than A and, as in (3.5), we convey this fact by writing:

$$A \begin{array}{c} \xrightarrow{r} \\ \leq \\ \xleftarrow{f} \end{array} B \quad \text{that is} \quad \{ f \cdot r = id_A \} \quad (4.4)$$

Inequation (4.4) has several informal interpretations:

- A is “smaller” than B
- B is able to “represent” A
- B is “abstracted” by A
- A is “implemented” by B
- B is a refinement (“refines”) A

In set theory one says that f is surjective². Note that the same surjective f can have many right-inverses. Consider, for instance, the following inequation expressing the implementation of Booleans by integers,

$$\text{Bool} \begin{array}{c} \xrightarrow{r} \\ \leq \\ \xleftarrow{f} \end{array} \mathbb{Z}$$

where

$$f \stackrel{\text{def}}{=} =_0 \rightarrow \underline{\text{FALSE}}, \underline{\text{TRUE}}$$

(informal meaning: 0 represents FALSE, any other value represents TRUE). It is easy to see that, for every $k \in \mathbb{Z} - \{0\}$,

$$r \stackrel{\text{def}}{=} id \rightarrow \underline{k}, \underline{0}$$

is a right-inverse of f .

Right-inverses are always *left-invertible*, i.e. *injective* in set theory terminology³. This is a very important property of data representation — no two different abstract values $a, a' \in A$ get mixed up along the representation process:

$$a \neq a' \Rightarrow (r a) \neq (r a')$$

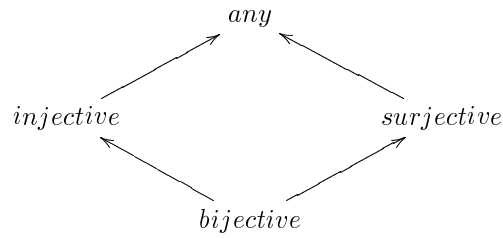
²In category theory one would say that f is *split-epic*, a concept stronger than *epic* (cf. section 1.6).

³*split-monic* in category theory terminology.

An isomorphism f is an arrow which has both a right-inverse r and a left-inverse s — a bijection in set theory terminology. It is easy to show that $r = s = f^{-1}$, as in (4.1)⁴:

$$\begin{aligned}
 & \begin{cases} f \cdot r = id_A \\ s \cdot f = id_B \end{cases} \\
 \Leftrightarrow & \quad \{ \text{composition and identity} \} \\
 & s \cdot (f \cdot r) = s \\
 \Leftrightarrow & \quad \{ \text{associativity} \} \\
 & (s \cdot f) \cdot r = s \\
 \Leftrightarrow & \quad \{ s \text{ is right-inverse of } f \text{ and identity} \} \\
 & r = s
 \end{aligned}$$

In summary, we have the following function terminological taxonomy:



In formal data processing, wherever f has a right-inverse r we say that f is an *abstraction function* (also called a *retrieve map*) and r a (particular choice of a) *representation function*.

4.3.1 Examples

Let us recall

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

and try to extend this to partial functions. Rather than an isomorphism, we get an inequation:

$$\begin{array}{ccc}
 & \text{pcurry} & \\
 (C \times A) \multimap B & \xrightarrow{\quad} & C \multimap (A \multimap B) \\
 & \xleftarrow{\text{unpcurry}} & \\
 & \leq &
 \end{array}$$

⁴Notation f^{-1} expresses the uniqueness of the inverse function of an isomorphism f .

where

$$\begin{aligned} pcurry & : (A \times B) \multimap C \multimap A \multimap (B \multimap C) \\ pcurry f & \stackrel{\text{def}}{=} \left(\begin{array}{c} a \\ \left(\begin{array}{c} \pi_2 p \\ f p \end{array} \right)_{p \in \text{dom } f \wedge (\pi_1 p) = a} \end{array} \right)_{a \in 2^{\pi_1(\text{dom } f)}} \end{aligned} \quad (4.5)$$

and

$$unpcurry \stackrel{\text{def}}{=} \lambda \sigma. \left(\begin{array}{c} (a, b) \\ (\sigma a) b \end{array} \right)_{a \in \text{dom } \sigma \wedge b \in \text{dom } (\sigma a)}$$

Clearly,

$$pcurry \cdot unpcurry \neq id$$

Doing the same with respect to *split* (1.18):

$$A \multimap (B \times C) \begin{array}{c} \xrightarrow{\text{unjoin}} \\ \leq \\ \xleftarrow{\boxtimes} \end{array} (A \multimap B) \times (A \multimap C) \quad (4.6)$$

where the finite mapping “join” operator (3.40) establishes how \multimap “distributes” over product \times . Its right-inverse is

$$unjoin \stackrel{\text{def}}{=} \langle A \multimap \pi_1, A \multimap \pi_2 \rangle \quad (4.7)$$

i.e. $unjoin \leftarrow split(* \multimap p1, * \multimap p2)$ in CAMILA notation.

The \leq -subcalculus encompasses other \multimap -laws which cannot be adapted from similar results concerning exponentials. For instance,

$$(B + C)^A \not\cong B^A \times B^C$$

and yet we have the following map-decomposition law:

$$A \multimap (B + C) \begin{array}{c} \xrightarrow{\text{uncojoin}} \\ \leq \\ \xleftarrow{\boxplus} \end{array} (A \multimap B) \times (A \multimap C) \quad (4.8)$$

where

$$\boxplus \stackrel{\text{def}}{=} \dagger \cdot ((A \multimap i_1) \times (A \multimap i_2))$$

A very important law is

$$A \multimap D \times (B \multimap C) \begin{array}{c} \xrightarrow{\text{unnjoin}} \\ \leq \\ \xleftarrow{\boxtimes_n} \end{array} (A \multimap D) \times ((A \times B) \multimap C) \quad (4.9)$$

which is established by an elaboration of \bowtie

$$\bowtie_n \stackrel{\text{def}}{=} \bowtie \cdot \langle \pi_1, \dagger \cdot ((A \rightarrow \perp) \times pcurry) \rangle \quad (4.10)$$

such that

$$\begin{aligned} unjoin &\stackrel{\text{def}}{=} (id \times unpcurry) \cdot unjoin \\ &= \{ \text{definition of } unjoin \text{ (4.7)} \} \\ &\quad \langle A \rightarrow \pi_1, unpcurry \cdot (A \rightarrow \pi_2) \rangle \end{aligned}$$

is one of its right-inverses. Other standard representations (refinements) follow: for finite A

$$\mathcal{P}A \begin{array}{c} \xrightarrow{\leq} \\ \xleftarrow{dom} \end{array} A \rightarrow B \quad (4.11)$$

(representation consists of classifying data, $b_0 \in B$).

$$\mathcal{P}A \begin{array}{c} \xrightarrow{\leq} \\ \xleftarrow{ran} \end{array} \mathbb{N} \rightarrow A \quad (4.12)$$

(representation consists of indexing data).

$$A^* \begin{array}{c} \xrightarrow{\leq} \\ \xleftarrow{\quad} \end{array} \mathbb{N} \rightarrow A \quad (4.13)$$

(refinement consists of rendering sequence order explicit).

$$A \rightarrow B \begin{array}{c} \xrightarrow{\leq} \\ \xleftarrow{\quad} \end{array} B \rightarrow \mathcal{P}A \quad (4.14)$$

(refinement consists of representing a function by its “kernel”).

Exercise 4.1 Guess a formal definition of *uncojoin*.

□

Exercise 4.2 Show that

$$A \rightarrow (B \rightarrow C) \leq (\mathcal{P}A) \times ((A \times B) \rightarrow C) \quad (4.15)$$

holds.

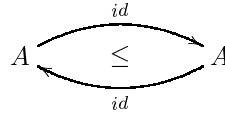
□

Exercise 4.3 Infer abstraction functions θ_1 and θ_2 whose righ-inverses are $A \rightarrow i_1$ and $A \rightarrow i_2$, respectively. Infer the signatures of θ_1 and θ_2 . Explain why θ_1 and θ_2 are not the righ-inverses of $A \rightarrow i_1$ and $A \rightarrow i_2$, respectively.

□

4.3.2 Basic properties of \leq :

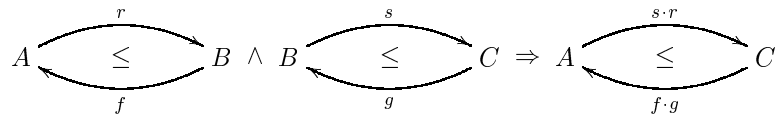
Reflexivity



cf.

$$id \cdot id = id$$

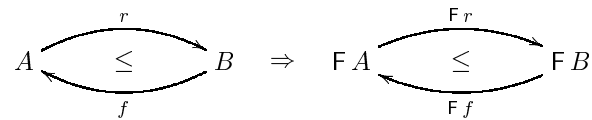
Transitivity



cf.

$$\begin{aligned}
 & (f \cdot g) \cdot (s \cdot r) \\
 = & \quad \{ \text{associativity} \} \\
 & f \cdot (g \cdot s) \cdot r \\
 = & \quad \{ s \text{ is right-inverse of } g \} \\
 & f \cdot id \cdot r \\
 = & \quad \{ \text{identity} \} \\
 & f \cdot r \\
 = & \quad \{ r \text{ is right-inverse of } f \} \\
 & id
 \end{aligned}$$

Monotonicity



where F is an arbitrary functor:

$$\begin{aligned}
 & id \\
 = & \quad \{ \text{functors commute with } id \}
 \end{aligned}$$

$$\begin{aligned}
 & \text{F id} \\
 = & \{ r \text{ is right-inverse of } f \} \\
 & \text{F}(f \cdot r) \\
 = & \{ \text{functors commute with composition} \} \\
 & (\text{F} f) \cdot (\text{F} r)
 \end{aligned}$$

therefore $\text{F} r$ is right-inverse of $\text{F} f$. Of course, this result extends to bifunctors, *e.g.*

$$A \begin{array}{c} \xrightarrow{r} \\ \leq \\ \xleftarrow{f} \end{array} B \wedge C \begin{array}{c} \xrightarrow{s} \\ \leq \\ \xleftarrow{g} \end{array} D \Rightarrow B(A, C) \begin{array}{c} \xrightarrow{B(r,s)} \\ \leq \\ \xleftarrow{B(f,g)} \end{array} B(B, D)$$

or n -ary functors in general.

4.3.3 Induction (Structural monotonicity)

Given two inductive datatypes μF and μG defined by functors F and G , respectively,

$$\begin{array}{ccc}
 \mu F & \xrightarrow{\cong} & F \mu F \\
 & \text{in}_{\mu F} & \\
 \mu G & \xrightarrow{\cong} & G \mu G \\
 & \text{in}_{\mu G} &
 \end{array}$$

and such that

$$F X \begin{array}{c} \xrightarrow{s} \\ \leq \\ \xleftarrow{g} \end{array} G X \tag{4.16}$$

Then

$$\mu F \begin{array}{c} \xrightarrow{r=(\text{in}_{\mu G} \cdot s)_F} \\ \leq \\ \xleftarrow{f=(\text{in}_{\mu F} \cdot g)_G} \end{array} \mu G \tag{4.17}$$

holds, *cf.* diagram

$$\begin{array}{ccccc}
 & & \text{in}_{\mu G} & & \\
 & & \curvearrowright & & \\
 \mu G & & F \mu G & \xleftarrow{g} & G \mu G \\
 \downarrow (\text{in}_{\mu F} \cdot g)_G & & \downarrow F(\text{in}_{\mu F} \cdot g)_G & & \downarrow G(\text{in}_{\mu F} \cdot g)_G \\
 \mu F & \xleftarrow{\text{in}_{\mu F}} & F \mu F & \xleftarrow{g} & G \mu F
 \end{array}$$

This property is very important in data-refinement calculations. Let us see some examples of application:

1. Empty sequences represented by non-empty ones — $\mu F = A^*$ and $\mu G = A^+$, for $1 \leq A$:

$$\begin{array}{ccc} A^* & \xrightarrow[\cong]{[\underline{\square}, cons]} & 1 + A \times A^* \\ A^+ & \xrightarrow[\cong]{[sing, cons]} & A + A \times A^+ \end{array}$$

where $sing\ a = [a]$.

Of course

$$\begin{array}{ccc} 1 + A \times X & \xrightarrow[\leq]{\underline{a_0} + id} & A + A \times X \\ & \xleftarrow{! + id} & \end{array}$$

holds for some $a_0 \in A$ since

$$\begin{array}{ccc} 1 & \xrightarrow[\leq]{a_0} & A \\ & \xleftarrow{!} & \end{array}$$

by hypothesis. According to (4.17) we infer not only that non-empty lists implement empty lists

$$\begin{array}{ccc} A^* & \xrightarrow[\leq]{embed} & A^+ \\ & \xleftarrow{blast} & \end{array}$$

(obvious!) but how they do it:

$$\begin{aligned} blast &= ([\underline{\square}, cons] \cdot (! + id)) \\ embed &= ([sing, cons] \cdot (\underline{a_0} + id)) \end{aligned}$$

(less obvious). Let us calculate further:

$$\begin{aligned} & blast \cdot [sing, cons] \\ = & \quad \{ \text{cata-cancellation} \} \\ & [\underline{\square}, cons] \cdot (! + id) \cdot (id + id \times blast) \\ = & \quad \{ +\text{-bifunctor and trivias} \} \\ & [\underline{\square}, cons] \cdot (! + id \times blast) \\ = & \quad \{ +\text{-absorption and trivias} \} \\ & [\underline{\square}, cons \cdot (id \times blast)] \end{aligned}$$

that is

$$\begin{aligned} \text{blast}[a] &= [] \\ \text{blast}(\text{cons}(a, l)) &= \text{cons}(a, \text{blast } l) \end{aligned}$$

(blast="all but last")

Similarly:

$$\begin{aligned} \text{embed}[] &= [a_0] \\ \text{embed}(\text{cons}(a, l)) &= \text{cons}(a, \text{embed } l) \end{aligned}$$

2. Representing *GenDia* (genealogical diagrams) by *DecTree* (decision trees):

Genealogical diagrams:

$$\begin{aligned} \text{GenDia} :: \quad & I : \text{Ind} && \text{/*data about an individual */} \\ & F : [\text{GenDia}] && \text{/*genealogy of his/her father (if known) */} \\ & M : [\text{GenDia}] && \text{/*genealogy of his/her mother (if known) */} \end{aligned}$$

i.e.

$$\text{GenDia} \cong \text{Ind} \times (\text{GenDia} + 1) \times (\text{GenDia} + 1) \quad (4.18)$$

Decision trees:

$$\begin{aligned} \text{DecTree} :: \quad & Q : \text{What} && \text{/*Question or Decision */} \\ & R : \text{Answer} \rightarrow \text{DecTree} && \text{/*Subtrees */} \\ \text{What} &= \dots \\ \text{Answer} &= \dots \end{aligned}$$

i.e.

$$\text{DecTree} \cong \text{What} \times (\text{Answer} \rightarrow \text{DecTree}) \quad (4.19)$$

Clearly,

$$\text{GenDia} \cong F(\text{Ind}, \text{GenDia})$$

where

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

and

$$\text{DecTree} \cong G(\text{What}, \text{Answer}, \text{DecTree})$$

where

$$G(A, B, X) \stackrel{\text{def}}{=} A \times (B \rightarrow X)$$

Reasoning: instantiate $B = 2$. Then

$$\begin{aligned} & G(A, 2, X) \\ = & \quad \{ \text{instantiation} \} \\ & A \times (2 \rightarrow X) \\ \cong & \quad \{ \text{totalizing maps} \} \\ & A \times (X + 1)^2 \\ \cong & \quad \{ \text{expanding square} \} \\ & A \times ((X + 1) \times (X + 1)) \\ \cong & \quad \{ \text{flattening (removing parentheses)} \} \\ & A \times (X + 1) \times (X + 1) \\ = & \quad \{ \text{instantiation} \} \\ & F(A, X) \end{aligned}$$

Since

$$F(A, X) \overset{\cong}{\rightleftarrows} G(A, 2, X)$$

we can deduce:

$$GenDia A \overset{\cong}{\rightleftarrows} DecTree(A, 2)$$

Challenge: calculate abstraction/representation maps.

3. Last but not least: finite sets represented by finite lists:

$$\mathcal{P}A \overset{\{in\}}{\underset{elems=\{ins\}}{\leq}} A^* \tag{4.20}$$

where

$$in = [[], cons]$$

and

$$ins = [\emptyset, puts]$$

— recall (3.2).

4.3.4 Constructive proof

We want to find a right-inverse for $(in_{\mu F} \cdot g)$ in (4.17) expressible as a catamorphism (α) , cf. diagram

$$\begin{array}{ccc}
 \mu F & \xleftarrow{in_{\mu F}} & F \mu F \\
 (\alpha)_F \downarrow & & \downarrow F(\alpha)_F \\
 \mu G & \xleftarrow{\alpha} & F \mu G \\
 (in_{\mu F} \cdot g)_G \downarrow & & \downarrow F(in_{\mu F} \cdot g)_G \\
 \mu F & \xleftarrow{in_{\mu F} \cdot g} & F \mu F
 \end{array}$$

So α is the unknown. Calculation of α :

$$\begin{aligned}
 & (in_{\mu F} \cdot g)_G \cdot (\alpha)_F = id \\
 \Leftrightarrow & \quad \{ \text{by F-cata-reflexion} \} \\
 & (in_{\mu F} \cdot g)_G \cdot (\alpha)_F = (in_{\mu F})_F \\
 \Leftarrow & \quad \{ \text{by F-cata-fusion} \} \\
 & (in_{\mu F} \cdot g)_G \cdot \alpha = in_{\mu F} \cdot F (in_{\mu F} \cdot g)_G \\
 \Leftrightarrow & \quad \{ \text{decompose } \alpha = in_{\mu G} \cdot \alpha' \} \\
 & (in_{\mu F} \cdot g)_G \cdot in_{\mu G} \cdot \alpha' = in_{\mu F} \cdot F (in_{\mu F} \cdot g)_G \\
 \Leftrightarrow & \quad \{ \text{by G-cata-cancellation} \} \\
 & in_{\mu F} \cdot g \cdot G (in_{\mu F} \cdot g)_G \cdot \alpha' = in_{\mu F} \cdot F (in_{\mu F} \cdot g)_G \\
 \Leftrightarrow & \quad \{ g \text{ is natural (4.21)} \} \\
 & in_{\mu F} \cdot F (in_{\mu F} \cdot g)_G \cdot g \cdot \alpha' = in_{\mu F} \cdot F (in_{\mu F} \cdot g)_G \\
 \Leftarrow & \quad \{ g \cdot s = id \} \\
 & \alpha' = s
 \end{aligned}$$

So

$$\alpha = in_{\mu G} \cdot s$$

NB: both g and s are natural, e.g.

$$(F f) \cdot g = g \cdot (G f) \quad \begin{array}{ccc} A & & G A \xrightarrow{g_A} F A \\ f \downarrow & & \downarrow F f \\ B & & G B \xrightarrow{g_B} F B \end{array} \quad (4.21)$$

holds.

4.4 Concrete invariants

Given some abstraction map f , the choice of a particular representation r induces a low-level invariant ϕ

$$A \begin{array}{c} \xrightarrow{r} \\ \cong \\ \xleftarrow{f} \end{array} B \xrightarrow{\phi} \mathbb{2}$$

called the *concrete invariant* induced by r . It emerges as characteristic predicate of the range of r , that is to say, ϕb holds wherever there exists at least one $a \in A$ such that $b = r a$. Because r is always injective, one has

$$A \begin{array}{c} \xrightarrow{r} \\ \cong \\ \xleftarrow{f} \end{array} B_\phi$$

where B_ϕ denotes the subset of B which satisfies concrete-invariant ϕ . So the replacement of low-level structure B by abstract structure A is safe provided ϕ is known to hold in the particular situation in hand — a kind of reverse specification via “concrete invariant discharge”. That is to say, we have

$$\begin{aligned} f \cdot r &= id \\ \phi b &\stackrel{\text{def}}{=} (r \cdot f)b = b \end{aligned}$$

Some examples follow:

$$A \rightarrow (B \times C) \begin{array}{c} \xrightarrow{\text{unjoin}} \\ \cong \\ \xleftarrow{\text{⋈}} \end{array} ((A \rightarrow B) \times (A \rightarrow C))_{eqd}$$

where

$$eqd \stackrel{\text{def}}{=} (=) \cdot (dom \times dom)$$

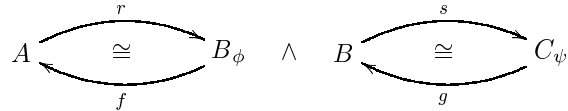
while

$$A \rightarrow (B + C) \begin{array}{c} \xrightarrow{\text{uncojoin}} \\ \cong \\ \xleftarrow{\text{⋈}^\dagger} \end{array} ((A \rightarrow B) \times (A \rightarrow C))_{djd}$$

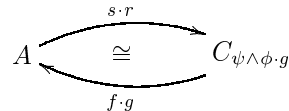
where

$$djd \stackrel{\text{def}}{=} =_\emptyset \cdot \cap \cdot (dom \times dom)$$

4.4.1 Concrete invariant handling rule



implies



For example, from (3.29) and (4.20) we infer:

$$A \multimap B \cong (A \times B)_{f \text{ dp. elems}}^*$$

4.5 Representing inductive datatypes in ordinary programming languages

4.5.1 The relational database model

A data model is said to be compliant with the *relational database model* if it is a finitary product of expressions of the following kind:

- **relations** — $\mathcal{P}(A_1 \times \dots \times A_n)$ where every A_i is regarded as an “atomic type” in the target relational environment;
- **mappings** — $(A_1 \times \dots \times A_n) \multimap (B_1 \times \dots \times B_m)$ where every A_i and B_j is regarded as an “atomic type” in the target relational environment.

Then:

- Laws such as (4.6,4.15,4.9,4.8) and (4.13) play an important rôle in derived construct reification. They can be used to “decompose” complex/nested mappings or sequences into tuples of simpler maps.
- Since finite maps are easily refined into binary relations — recall (3.29) — our calculus fragment above is useful in refining elaborate, finite-mapping/sequence based data models into relational database schemata.
- It has been proved elsewhere [3] that **3NF** is guaranteed by mere application of the above laws, among others.

Exercise 4.4 In relational data refinement (*i.e.* the representation of data in the relational database model) one often ignores normalization wherever data storage is not a problem, *i.e.*, one is dealing with small data sets. The following law describes one such refinement step:

$$(A \multimap B) \times (A \multimap C) \leq A \multimap ((B + 1) \times (C + 1)) \tag{4.22}$$

It merges two tables into a single one by, for instance, representing

NÚMERO	NOME
1010	Manuel
11230	Maria

and

NÚMERO	CURSO
11230	LESI
15234	LMCC

by

NÚMERO	NOME	CURSO
1010	Manuel	NIL
11230	Maria	LESI
15234	NIL	LMCC

1. Resort to θ_1 of exercise 4.3 to define an abstraction map f for (4.22) in point-free notation.
2. Define the obvious right-inverse of f .

□

4.5.2 Pointer introduction

For target languages like C/C++ we need a different approach. Recursive data models are coped with by resorting to pointers as shown next. First of all, we have the following basic fact: every non-empty datatype A can be represented by a “pointer”:

$$A \begin{array}{c} \xrightarrow{r=i_2} \\ \leq \\ \xleftarrow{f=[a_0, id]} \end{array} 1 + A \quad (a_0 \in A)$$

i.e.

$$A \cong (1 + A)_{\phi \ x \stackrel{\text{def}}{=} \exists a \in A : x = (i_2 \ a)}$$

(f abstracts both NIL and a_0 into a_0 . It is therefore non-injective.) Then we have the following implementation rules:

- Inductive datatypes of form

$$T \begin{array}{c} \xrightarrow{out} \\ \cong \\ \xleftarrow{in} \end{array} 1 + G T$$

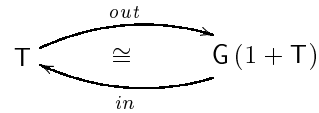
can be directly implemented in such languages, *e.g.*

$$A^* \cong 1 + \underbrace{A \times A^*}_{G A^*}$$

implemented by


```
typedef struct G {
    A first;
    struct G *next;
} *T;
```

- Inductive datatypes of form



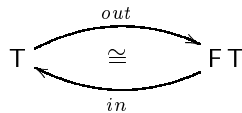
are also directly implementable, e.g.

$$\text{GenDia} \cong \text{Ind} \times (1 + \text{GenDia}) \times (1 + \text{GenDia}) \quad (4.23)$$

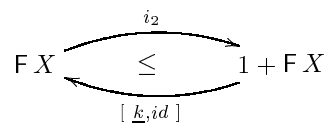
implemented by

```
typedef struct GenDia {
    Ind individual;          /* data about an individual */
    struct GenDia *father; /* genealogy of his/her father
                           (if known) */
    struct GenDia *mother; /* genealogy of his/her mother
                           (if known) */
};
```

- All other cases:



implemented by forcing a pointer structure:



and therefore (4.17)

$$\begin{array}{ccc}
 & \xrightarrow{(in \cdot i_2)_G} & \\
 T & \leq & U \\
 & \xleftarrow{(in \cdot [\underline{k}, id])_F} &
 \end{array} \quad (4.24)$$

where

$$U \cong 1 + FU$$

Example:

$$DecTree \cong What \times (Answer \rightarrow DecTree)$$

Via (3.29) this is implemented by

$$DecTree_1 \cong What \times (Answer \times DecTree_1)^*$$

On its turn, via (4.24), $DecTree_1$ is implemented by

$$DecTree_2 \cong 1 + What \times (Answer \times DecTree_2)^*$$

that is

$$\begin{cases} DecTree_2 \cong 1 + What \times subTrees \\ subTrees \cong 1 + (Answer \times DecTree_2) \times subTrees \end{cases}$$

This has a direct encoding in C as follows (for *Answer* and *What* = strings):

```
typedef struct DecTree2{
    char *What; /* Question or Decision */
    struct subTrees *R;
};

typedef struct subTrees {
    struct node *first;
    struct subTrees *next;
};

typedef struct node {
    char *Answer;
    struct DecTree2 *SubTree;
};
```

4.5.3 Recursion “removal”

What about SQL-like languages? How do we convert recursive structures in non-recursive tabular ones?

Let

$$\begin{array}{ccc} X_0 & \xrightarrow{\cong} & F X_0 \\ & \xleftarrow{in} & \end{array}$$

be a solution to

$$X \cong F X \tag{4.25}$$

Implementation step:

$$\begin{array}{ccc} X_0 & \xrightarrow{\cong} & (K \times (K \rightarrow F K))_\phi \\ & \xleftarrow{f} & \end{array} \tag{4.26}$$

for K a domain of “pointers” such that $K \cong \mathbb{N}$.

Abstraction Function f

We need to define a surjection

$$(K \times (K \rightarrow F K))_{\phi} \xrightarrow{f} X_0$$

that is, to find appropriate f and ϕ .

Temporarily assume that $\sigma \in K \rightarrow F K$ is a total function, and draw the diagram of coalgebra (K, σ) :

$$\begin{array}{c} K \\ \downarrow \sigma \\ F K \end{array}$$

Assuming a given piece of “linear storage” σ (“database”), let f_{σ} denote the function which, for each “pointer” $k \in K$, retrieves the value of X_0 corresponding to a σ “scan” starting from k :

$$\begin{array}{ccc} X_0 & \xleftarrow{f_{\sigma}} & K \\ & & \downarrow \sigma \\ & & F K \end{array}$$

The diagram can be immediately completed,

$$\begin{array}{ccc} X_0 & \xleftarrow{f_{\sigma}} & K \\ \uparrow in & & \downarrow \sigma \\ F X_0 & \xleftarrow{F f_{\sigma}} & F K \end{array}$$

expressing f_{σ} as F-hylomorphism $\llbracket \sigma, in \rrbracket_F$, that is,

$$f_{\sigma} = in \cdot (F f_{\sigma}) \cdot \sigma$$

or

$$f \quad : \quad K \times (K \rightarrow F K) \longrightarrow X_0 \tag{4.27}$$

$$f(k, \sigma) \stackrel{\text{def}}{=} in((F f_{\sigma})(\sigma k))$$

writing $f(k, \sigma)$ instead of “curried” $f_{\sigma} k$.

Concrete Invariant ϕ

Face pointer undefinedness in case σ is partial — σk in (4.27) is undefined wherever $k \notin \text{dom } \sigma$ and other pointers k' in the range of σ , reachable from k , may be in the same situation.

- Pointer reachability — transitive closure of

$$k_1 <_F k \stackrel{\text{def}}{=} k \in \text{dom } \sigma \wedge k_1 \in_F \sigma k \quad (4.28)$$

(recall that \in_F is given in section 3.3), that is

$$\phi(k, \sigma) \stackrel{\text{def}}{=} \begin{array}{l} \text{let } P = \{k\} \cup \{k' \in K \mid k' <_F^+ k\} \\ \text{in } P \subseteq \text{dom } \sigma \end{array}$$

- Well-foundedness — Restrict to least fixpoints, *i.e.* guarantee that $f(k, \sigma)$ does not yield infinite results, forcing P to be well-founded wrt. $<_F$:

$$\phi(k, \sigma) \stackrel{\text{def}}{=} \begin{array}{l} \text{let } P = \{k\} \cup \{k' \in K \mid k' <_F^+ k\} \\ \text{in } P \subseteq \text{dom } \sigma \wedge \\ \forall \emptyset \subset C \subseteq P : \exists m \in C : \forall k' <_F m : k' \notin C \end{array} \quad (4.29)$$

The proof of surjectiveness of f (4.27) for $(X_0, F X_0 \xrightarrow{\text{in}} X_0)$ a finite or denumerable fixpoint of F , is a lengthy one!

4.5.4 Trivial Example

(C, id_C) is of course a fixpoint of constant functor $F X = C$. Then

$$\begin{aligned} f(k, \sigma) &\stackrel{\text{def}}{=} \text{in}((F f_\sigma)(\sigma k)) \\ &= id_C((C f_\sigma)(\sigma k)) \\ &= id_C(id_C(\sigma k)) \\ &= \sigma k \end{aligned}$$

Since $<_F^+$ is empty, we have

$$\phi(k, \sigma) \stackrel{\text{def}}{=} k \in \text{dom } \sigma \quad (4.30)$$

as expected — *cf.* dynamic programming typical of *e.g.* C or PASCAL— instead of handling C data directly (statically stored), the program handles *dynamic* references to them.

4.5.5 Implementing *DecTree*

Recall

$$DecTree \cong What \times (Answer \rightarrow DecTree)$$

4.5.6 Data-level Reification

First transformation is recursion removal according to (4.26)

$$DecTree \leq_1 DecTree_1$$

leading to

$$DecTree_1 = K \times (K \rightarrow What \times (Answer \rightarrow K)) \quad (4.31)$$

What have we achieved? Two interpretations (at least) are admissible:

1. For K a domain of *pointers*,

$$K \rightarrow What \times (Answer \rightarrow K) \quad (4.32)$$

— that is:

$$K \rightarrow What \times (K + 1)^{Answer}$$

— models a heap-segment of dynamic storage. In PASCAL (the heap is “hidden” in the run-time system):

```

type DecTree1 = ^DecTree;
   DecTree = record
               Q: What;
               R: array [Answer] of ^DecTree
             end;

```

2. For K a domain of *object names* or “unique identifiers”, (4.32) models the object database

$$name \rightarrow object$$

implicit in an object-oriented programming environment, cf. the following CAMILA:

```

DecTree_1 :: ObjName: K
           Archive: ObjBase ;
ObjBase   = K -> Attributes
Attributes :: Q: What
           SubObjs: Answer -> K ;

```

Proceeding via law (4.9):

$$\begin{aligned}
 DecTree_1 &= K \times (K \rightarrow What \times (Answer \rightarrow K)) \\
 &\leq_2 K \times ((K \rightarrow What) \times ((K \times Answer) \rightarrow K)) \\
 &= DecTree_2
 \end{aligned} \quad (4.33)$$

For K a state descriptor, $DecTree_2$ accepts the following interpretation:

- $(K \times Answer) \rightarrow K$ = *state transition diagram* of a (deterministic) *finite state automaton* ($Answer$ = input stimuli);

- first factor $k \in K$ in (4.33) = current state of the automaton;
- $K \rightarrow What$ = semantic table assigning a meaning to each state.

Carrying on our reasoning,

$$\begin{aligned}
 DecTree_2 &= K \times ((K \rightarrow What) \times ((K \times Answer) \rightarrow K)) \\
 &\leq_3 K \times (\mathcal{P}(K \times What) \times \mathcal{P}((K \times Answer) \times K)) \\
 &= DecTree_3 \\
 &\cong_4 K \times (\mathcal{P}(K \times What) \times \mathcal{P}(K \times Answer \times K)) \\
 &= DecTree_4
 \end{aligned}$$

we obtain a final model,

$$DecTree_4 = K \times (\mathcal{P}(K \times What) \times \mathcal{P}(K \times Answer \times K)) \quad (4.34)$$

which is a relational database schema implementing $DecTree$ in terms of two database files where K plays the rôle of a domain of *keys*.

In summary,

$$DecTree \leq_1 DecTree_1 \leq_2 DecTree_2 \leq_3 DecTree_3 \cong_4 DecTree_4$$

4.6 Exercises

Exercise 4.5 Let f and g be two abstraction (i.e. surjective) functions. Identify which of the following combination of f and g are also abstraction functions:

$$[f, g] \quad (4.35)$$

$$f \times g \quad (4.36)$$

$$\langle f, g \rangle \quad (4.37)$$

$$f + g \quad (4.38)$$

□

Exercise 4.6

1. Apply law (4.17) to $T = \mathbb{N} \text{ — i.e. } in_F = [\underline{0}, succ] \text{ —, } U = A^*$ and $f = id + \pi_2$. Calculate a pointwise definition of abstraction function $l = \langle in_F \cdot f \rangle$. Which well-known function is l , after all?
2. Knowing that $U \cong T$ holds wherever f é iso (bijective) show that

$$\mathbb{N} \cong 1^*$$

and

$$1 \cong 0^*$$

hold.

□

Exercise 4.7 Show that *untot* can be defined by

$$\text{untot } \sigma \stackrel{\text{def}}{=} ((\overline{\text{pap}} \sigma) + !_A) \cdot (\lambda a. a \in \text{dom } \sigma)$$

for

$$\text{pap}(\sigma, a) \stackrel{\text{def}}{=} \{ a \in \text{dom } \sigma \Rightarrow \sigma a$$

□

Exercise 4.8 The *feature bundle* is a parametric datatype

$$\text{FB}(A, B) = A \multimap (B + \text{FB}(A, B))$$

which is widely useful in specifications, e.g.

$$\text{FileSystem} = \text{FB}(\text{FileName}, \text{File})$$

1. Define algebra

$$\mathcal{P}\text{FileName} \xleftarrow{g} A \multimap (B + \mathcal{P}\text{FileName})$$

so that $\{g\}$ specifies the function which collects all filenames in a hierarchical *FileSystem*.

2. FB can be used to specify “universal” data-representation media, e.g. the abstract interface data description language IDL:

$$\text{IDL} = \text{FB}(\text{STR}, \text{STR})$$

Calculate a C/C++ pointer-based implementation of IDL.

3. Knowing that representation

$$\begin{array}{ccc} B^* & \begin{array}{c} \xrightarrow{r} \\ \leq \\ \xleftarrow{f} \end{array} & \text{FB}(\mathbb{N}, B) \end{array}$$

holds, express *r* as a list-catamorphism.

4. Resort to (4.17) in order to calculate *f* and *r* in the obvious facts

$$\begin{array}{ccc} A \multimap B & \begin{array}{c} \xrightarrow{r} \\ \leq \\ \xleftarrow{f} \end{array} & \text{FB}(A, B) \\ \\ \mathcal{P}A & \begin{array}{c} \xrightarrow{r} \\ \leq \\ \xleftarrow{f} \end{array} & \text{FB}(A, 1) \end{array}$$

5. Can you represent $A \times B$ or $A + B$ in FB?

□

Exercise 4.9 Let

$$A + A \times X \xleftarrow{g=id+\pi_2} A$$

in (4.16). What abstraction function *f* do you obtain by (4.17) in this case? Justify.

□

Bibliography

- [1] C. J. Rodrigues. Sobre o desenvolvimento formal de bases de dados. Master's thesis, University of Minho, 1993. (In Portuguese).