

Calculate databases with ‘simplicity’

J.N. Oliveira

DI/U.Minho   Braga, Portugal

IFIP 2.1 #59 Meeting,
Nottingham, UK

September 2004

Abstract

- Effort to replace “à la Codd” database schema design (normalization etc) by **calculation** based on **simple** (dually, **injective**) binary relations.
- Simple relations relevant because database entities can be modelled as finite such relations.
- **(Pointfree)** calculus simpler to use than the standard theory.
- Generic result which enables the refinement of **recursive** data models
- Prospect of automatic SQL generation (using **Strafunski / Haskell**) based on results so far.

Motivation

- **SQL** — **data-processing** standard “de facto”
- **XML** — **abstract syntax** “made popular”
- Can XML be trusted as a data-storage technology?
- Ad hoc XML ↔ SQL conversion
- Need for reliable XML ↔ SQL data exchange technology

Example

(Haskell instead of XML, if you don't mind):

```
type StringExp = Exp String String

data Exp v o = Var v |
             Term o [Exp v o]
```

How do you SQL-archive StringExp data?

Example — SQL

```
CREATE TABLE SYMBOLS (
  Symbol CHAR (20) NOT NULL,
  NodeId NUMERIC (10) NOT NULL,
  IfVar BOOLEAN NOT NULL
  CONSTRAINT SYMBOLS_pk
    PRIMARY KEY(NodeId,IfVar)
);

CREATE TABLE EXPRESSIONS (
  FatherId NUMERIC (10) NOT NULL,
  ArgNr NUMERIC (10) NOT NULL,
  ChildId NUMERIC (10) NOT NULL
  CONSTRAINT EXPRESSIONS_pk
    PRIMARY KEY (FatherId,ArgNr)
);
```

Can you **rely** on this implementation?

Overall idea

- Calculate implementations from specifications

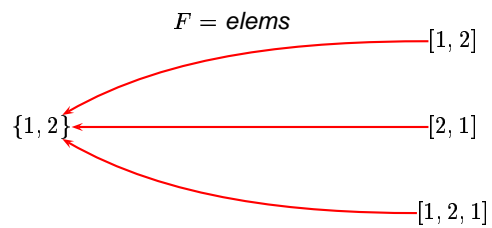
$$\begin{aligned} \text{Spec} &= X \\ &\leq X' \\ &\leq X'' \\ &\leq \dots \\ &\leq \text{Imp} \end{aligned}$$

by adding details in a controlled manner.

- Define a suitable ordering \leq on datatypes and develop corresponding data refinement theory

Example of data refinement

Finite sets represented by finite lists:



Refinement inequation

$$\mathcal{P}_f A \overset{\leq}{\underset{\text{elems}}{\curvearrowright}} A^*$$

meaning that

- sets are “implemented” by lists
- A^* is able to “represent” $\mathcal{P}_f A$
- A^* is “abstracted” by $\mathcal{P}_f A$
- A^* is a refinement (“refines”) $\mathcal{P}_f A$

Refinement inequations

A is implemented by B , as witnessed by pair f, r , iff

$$A \overset{r}{\underset{f}{\curvearrowright}} B$$

such that

- **representation** r is injective
- **abstraction** f is surjective
- that is,

$$f \cdot r = id$$

Not general enough (I)

In the following inequation

$$A \overset{i_1}{\underset{i_1^\circ}{\curvearrowright}} A + 1$$

expressing the fact that every element of datatype A can be represented by a “pointer”,

- $r = i_1$ is injective, but
- its converse i_1° is **partial** (=not entirely defined)

Not general enough (II)

Representations r need not be functions. Back to

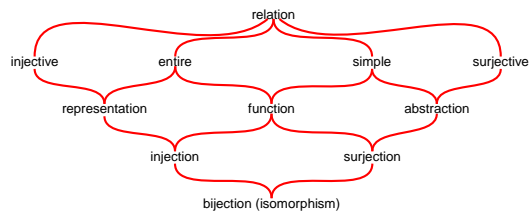
$$\mathcal{P}_f A \begin{array}{c} \xrightarrow{R} \\ \xleftarrow{\leq} \\ \text{elems} \end{array} A^*$$

relation $R = \text{elems}^\circ$ will be perfectly acceptable as a representation since

$$\text{elems} \cdot \text{elems}^\circ = id$$

because elems is a surjection.

Binary relation taxonomy



Terminology: simple / entire relation instead of partial / total function (or relation)

Taxonomy basis

	Reflexive	Coreflexive
$\ker R$	entire R	injective R
$\text{img } R$	surjective R	simple R

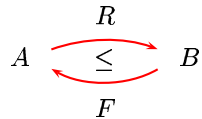
where

- **Reflexive** relation: $id \subseteq R$
- **Coreflexive** relation (or **partial identity**): $R \subseteq id$
- **Kernel** and **image**:

$$\ker R = R^\circ \cdot R$$

$$\text{img } R = R \cdot R^\circ \quad (= \ker (R^\circ))$$

Principle of data abstraction



where

- $A \xleftarrow{F} B$ is a **surjective + simple abstraction** relation
- R is **entire** and $R \subseteq F^\circ$ — it is said to be a **representation** for F .

The fact that R is **injective** follows from $R \subseteq F^\circ$.

Summary

$\ker R = id$	entire $R \wedge$ injective R	representation R
$img F = id$	surjective $F \wedge$ simple F	abstraction F

It follows that R is a **right-inverse** of F , that is

$$F \cdot R = id$$

This is proved by circular inclusion

$$F \cdot R \subseteq id \subseteq F \cdot R$$

in the next slide.

Right invertibility

$$\begin{aligned}
 & F \cdot R \subseteq id \wedge id \subseteq F \cdot R \\
 \equiv & \quad \{ \text{img } F = id \text{ and } \text{ker } R = id \} \\
 & F \cdot R \subseteq F \cdot F^\circ \wedge R^\circ \cdot R \subseteq F \cdot R \\
 \equiv & \quad \{ \text{converse of right conjunct} \} \\
 & F \cdot R \subseteq F \cdot F^\circ \wedge R^\circ \cdot R \subseteq R^\circ \cdot F^\circ \\
 \Leftarrow & \quad \{ (F \cdot) \text{ and } (R^\circ \cdot) \text{ are monotonic} \} \\
 & R \subseteq F^\circ \wedge R \subseteq F^\circ \\
 \equiv & \quad \{ R \subseteq F^\circ \text{ is assumed} \} \\
 & \text{TRUE}
 \end{aligned}$$

Functional abstractions

Quotation from [Mor90], chapter 17, pp. 173–174:

[It is] common for the coupling invariant to be functional from concrete to abstract. [...]

*The general form for such coupling invariants, called **functional abstractions** is*

$$a = af \cdot c \wedge dti \cdot c, \quad (17.1)$$

*where af is a function, called the **abstraction function** and dti is a predicate, in which a does not appear, called the **data-type invariant**. [...]*

That is,

$$F = af \cdot [dti]$$

Functional abstractions

- **Galois abstractions** — let $R, F := f^\flat, f$ be Galois connected functions where the connection is **perfect** on the “abstract side”,

$$f \cdot f^\flat = id$$

Example: **hash-table** representation of a data collection [OR04]

- **Isomorphisms** —

$$\begin{array}{ccc}
 & r & \\
 A & \xrightarrow{\cong} & B \quad \text{such that } r = f^\circ \\
 & \xleftarrow{f} &
 \end{array}$$

\leq is a preorder

Reflexivity

$$A \begin{array}{c} \xrightarrow{id} \\ \leq \\ \xleftarrow{id} \end{array} A$$

Transitivity

$$A \begin{array}{c} \xrightarrow{R} \\ \leq \\ \xleftarrow{F} \end{array} B \wedge B \begin{array}{c} \xrightarrow{S} \\ \leq \\ \xleftarrow{G} \end{array} C \Rightarrow A \begin{array}{c} \xrightarrow{S \cdot R} \\ \leq \\ \xleftarrow{F \cdot G} \end{array} C$$

Proof of transitivity

a) Composition preserves simplicity and surjectiveness:

$$\begin{aligned} & \text{img}(F \cdot G) = id \\ \equiv & \quad \{ \text{expand } \text{img} ; \text{converses} \} \\ & F \cdot (\text{img } G) \cdot F^\circ = id \\ \equiv & \quad \{ G \text{ is simple and surjective} \} \\ & \text{img } F = id \\ \equiv & \quad \{ F \text{ is simple and surjective} \} \\ & id = id \end{aligned}$$

b) $S \cdot R \subseteq (F \cdot G)^\circ$ by monotonicity.

Structural data refinement

For F a relator,

$$A \begin{array}{c} \xrightarrow{R} \\ \leq \\ \xleftarrow{F} \end{array} B \Rightarrow F A \begin{array}{c} \xrightarrow{F R} \\ \leq \\ \xleftarrow{F F} \end{array} F B$$

- Easy proof next slides.
- Also valid for n -ary relators such as \times , $+$ etc.

Structural data refinement

a) **Abstraction:**

$$\begin{aligned}
 & \text{img}(F F) \\
 = & \quad \{ \text{image definition ; relators commute with converse} \} \\
 & (F(F^\circ)) \cdot (F F) \\
 = & \quad \{ \text{relators commute with composition} \} \\
 & F(F^\circ \cdot F) \\
 = & \quad \{ F \text{ is an abstraction} \} \\
 & F \text{ id} \\
 = & \quad \{ \text{relators commute with id} \} \\
 & \text{id}
 \end{aligned}$$

Structural data refinement

b) **Representation:**

$$\begin{aligned}
 & F R \subseteq (F F)^\circ \\
 \equiv & \quad \{ \text{relators commute with converse} \} \\
 & F R \subseteq F(F^\circ) \\
 \Leftarrow & \quad \{ \text{relators are monotone} \} \\
 & R \subseteq F^\circ \\
 \equiv & \quad \{ R \text{ is a representation for } F \} \\
 & \text{TRUE}
 \end{aligned}$$

By the way

Datatypes such as $\mathcal{P}_f A$ are \leq postfix points [Bac00], cf.

$$\begin{array}{ccc}
 \mathcal{P}_f A & \xrightarrow{\text{ins}^\circ} & 1 + A \times \mathcal{P}_f A \\
 \downarrow \llbracket R, \text{ins}^\circ \rrbracket & \begin{array}{c} \leq \\ \text{ins} \end{array} & \downarrow \text{id} + \text{id} \times \llbracket R, \text{ins}^\circ \rrbracket \\
 B & \xleftarrow{R} & 1 + A \times B
 \end{array}$$

(hylomorphisms on finite sets).

Abstract database models

- A relational database is a tuple of **finite** relations
- Finite **simple** relations model **many-to-one** (M:1) relationships (inc. **primary key** relationships)
- Finite **simple+injective** relations model **one-to-one** (1:1) relationships

Notation:

- $B \leftarrow A$: all **simple** relations from A (the **key**) to B (the **data** of interest) —cf. (if also finite) `FiniteMap a b` in Haskell, `map A to B` in VDM-SL.
- $B \leftrightarrow A$: all **injective** relations from A to B —cf. (if also finite and simple) `inmap A to B` in VDM-SL,

“Maybe” transpose

Useful isomorphism

$$(B + 1)^A \xrightarrow{(\in \cdot)} \cong B \leftarrow A$$

Γ

converts **simple** relations into $(+ 1)$ -valued functions ($\in = i_1^\circ$):

$$f = \Gamma R \quad \equiv \quad (b R a \quad \equiv \quad (f a = i_1 b))$$

NB: generalizes to **generic transpose** [OR04]

... and exponentials

Multiple-key decomposition / synthesis:

$$\begin{aligned} & A \leftarrow B \times C \\ \cong & \quad \{ \Gamma \} \\ & (A + 1)^{B \times C} \\ \cong & \quad \{ \text{curry} \} \\ & ((A + 1)^C)^B \\ \cong & \quad \{ (\in \cdot)^B \} \\ & (A \leftarrow C)^B \end{aligned}$$

Calculating abs/reps

Altogether, the downwards isomorphism

$$(\in \cdot)^B \cdot \text{curry} \cdot \Gamma$$

is a convenient shorthand for a less readable pointwise **abstraction invariant**:

$$\begin{aligned} \overline{S} &= (\in \cdot) \cdot (\text{curry}(\Gamma S)) \\ &\equiv \{ \dots \text{relational calculus} \dots \} \\ (b, c) \in \text{dom } S &\equiv c \in \text{dom}(\overline{S} b) \wedge S(b, c) = \overline{S} b c \end{aligned}$$

NB: thanks to generic transpose, notation \overline{S} extends to other classes of relation.

Converse of simple is injective

$$\begin{aligned} B \times C &\leftrightarrow A \\ \text{IR} &\quad \{ (-^\circ) \text{ isomorphism} \} \\ A &\leftarrow B \times C \\ \text{IR} &\quad \{ \text{above} \} \\ (A \leftarrow C)^B & \\ \text{IR} &\quad \{ \text{isomorphism is } (-^\circ)^B \} \\ (C \leftarrow A)^B & \end{aligned}$$

Refinement by decomposition

Zip/unzip'ping simple relations:

$$B \times C \leftarrow A \stackrel{\text{unjoin}}{\leq} (B \leftarrow A) \times (C \leftarrow A)$$

where

$$\begin{aligned} \text{join} &= \langle -, - \rangle \\ \langle R, S \rangle &\stackrel{\text{def}}{=} (\pi_1^\circ \cdot R) \cap (\pi_2^\circ \cdot S) \\ \text{unjoin} &\stackrel{\text{def}}{=} \langle \pi_1 \leftarrow id, \pi_2 \leftarrow id \rangle \end{aligned}$$

where, for injective f ,

$$g \leftarrow f \stackrel{\text{def}}{=} (g \cdot) \cdot (\cdot f^\circ)$$

Refinement by decomposition

$$(B + C) \leftarrow A \stackrel{\text{uncojoin}}{\leq} (B \leftarrow A) \times (C \leftarrow A)$$

where

$$\begin{aligned} \text{uncojoin} &= \langle (i_1^\circ \cdot), (i_2^\circ \cdot) \rangle \\ \text{cojoin} &= \cup \cdot ((i_1 \cdot) \times (i_2 \cdot)) \end{aligned}$$

NB: $\text{cojoin} \cdot \text{uncojoin} = id$, since $\text{img } i_1 \cup \text{img } i_2 = id$.

Refinement by decomposition

Nested simplicity:

$$(D \times (C \leftarrow B)) \leftarrow A \stackrel{\text{unnjoin}}{\leq} (D \leftarrow A) \times (C \leftarrow (A \times B))$$

Definitions of $njoin$ and $unnjoin$ to follow from next slide's calculation

Calculation

$$\begin{aligned}
 & (D \times (C \leftarrow B)) \leftarrow A \\
 \cong & \quad \{ \text{Maybe transpose} \} \\
 & ((D \times (C \leftarrow B)) + 1)^A \\
 \leq & \quad \{ \text{Maybe-(right)strength is involved in the abstraction} \} \\
 & ((D + 1) \times (C \leftarrow B))^A \\
 \cong & \quad \{ \text{splitting} \} \\
 & (D + 1)^A \times (C \leftarrow B)^A \\
 \cong & \quad \{ \text{Maybe transpose and above} \} \\
 & (D \leftarrow A) \times (C \leftarrow A \times B)
 \end{aligned}$$

Getting away with finite lists

Several other \leq laws, eg.

$$\begin{array}{ccc}
 & \text{seq2index} & \\
 A^* & \begin{array}{c} \xrightarrow{\leq} \\ \xleftarrow{\text{list}} \end{array} & A \leftarrow \mathbb{N}
 \end{array}$$

such that, for instance,

$$\begin{aligned}
 \text{seq2index } [a, b, a] &= \{(a, 1), (b, 2), (a, 3)\} \\
 \text{list } \{(a, 11), (b, 12), (a, 33)\} &= [a, b, a]
 \end{aligned}$$

Getting away with recursion

Given

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

one has

$$\begin{array}{ccc}
 \mu F & \begin{array}{c} \xrightarrow{\leq} \\ \xleftarrow{F} \end{array} & \underbrace{(F K \leftarrow K)}_{\text{"heap"}} \times K
 \end{array}$$

for K a data type of "heap addresses", or "pointers", such that $K \cong \mathbb{N}$.

Abstraction function

- Main rôle in representation is played by simple **F-coalgebra** $F K \leftarrow K$, understood as a (finite) piece of “linear storage”, a “heap” or a “database” file.
- \bar{F} (recall \bar{F} notation from above), of type $(\mu F \leftarrow K)^{(F K \leftarrow K)}$, is nothing but the F-anamorphism combinator:

$$\begin{array}{ccc}
 \mu F & \xleftarrow{in} & F(\mu F) \\
 \bar{F}H \uparrow & & \uparrow F(\bar{F}H) \\
 K & \xrightarrow{H} & F K
 \end{array}
 \quad
 \bar{F} = \llbracket - \rrbracket_F
 \quad
 \bar{F} H = \mu X. in \cdot (F X) \cdot H$$

Partiality of implementation

Abstraction invariant $t = F(H, k)$ —that is, $t = (\bar{F}H)k$ —will hold only if

- $k \in \text{dom } H$, and
- the **accessibility** relation for H

$$\begin{array}{ccc}
 K & \xleftarrow{\prec_H} & K \\
 \prec_H & \stackrel{\text{def}}{=} & \in_F \cdot H
 \end{array}$$

is **well-founded** and **closed** ($K \xleftarrow{\in_F} F K$ is the membership of F .)

(Many details omitted here!)

Back to the *StringExp* example

Since

$$StringExp = \mu X. (String + String \times X^*)$$

we have:

$$\begin{aligned}
 &StringExp \\
 \leq &\quad \{ \text{remove recursion} \} \\
 &((String + String \times K^*) \leftarrow K) \times K \\
 \leq &\quad \{ \text{remove finite lists} \} \\
 &((String + String \times (K \leftarrow \mathbb{N})) \leftarrow K) \times K
 \end{aligned}$$

Example continued

$$\begin{aligned} &\leq \{ \text{recall } (B + C) \leftarrow A \leq (B \leftarrow A) \times (C \leftarrow A) \} \\ &\quad (String \leftarrow K) \times ((String \times (K \leftarrow IN)) \leftarrow K) \times K \\ &\leq \{ \text{remove nested } \leftarrow \} \\ &\quad (String \leftarrow K) \times (String \leftarrow K) \times (K \leftarrow (IN \times K)) \times K \\ &\cong \{ A \times A \cong A^2 \} \\ &\quad (String \leftarrow K)^2 \times (K \leftarrow (IN \times K)) \times K \\ &\cong \{ \text{recall } (A \leftarrow C)^B \cong A \leftarrow B \times C \} \\ &\quad \underbrace{(String \leftarrow (2 \times K))}_{SYMBOLS} \times \underbrace{(K \leftarrow (IN \times K)) \times K}_{EXPRESSIONS} \end{aligned}$$

Conclusions

- Database schema design as a special case of “do it by calculation” data refinement
- Computational alternative to state-of-the-art casuistic practice stemming from set-theoretic “normalization theory”
- Many more laws available, eg.

$$1 \leftarrow A \cong \mathcal{P}A$$

cf.

```
newtype Set a = MkSet (FiniteMap a ())
```

in the `FiniteMap / Set.lhs` Haskell libraries.