

Accepted Manuscript

A relation-algebraic approach to the “Hoare logic” of functional dependencies

José N. Oliveira

PII: S1567-8326(14)00014-9
DOI: [10.1016/j.jlap.2014.02.013](http://dx.doi.org/10.1016/j.jlap.2014.02.013)
Reference: JLAP 401



To appear in: *The Journal of Logic and Algebraic Programming*

Please cite this article in press as: J.N. Oliveira, A relation-algebraic approach to the “Hoare logic” of functional dependencies, *The Journal of Logic and Algebraic Programming* (2014), <http://dx.doi.org/10.1016/j.jlap.2014.02.013>

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

Highlights

- This paper expresses functional dependencies using binary relation algebra.
- It shows that functional dependency calculus is similar to Hoare logic.
- This is based on reasoning about the degree of injectivity of a relation.
- Applications are given for query optimization and lexicographic sorting.
- The optimisations are obtained by point-free calculations.

A relation-algebraic approach to the “Hoare logic” of functional dependencies

José N. Oliveira^a

^a*High Assurance Software Lab
INESC TEC and Univ. Minho, 4710-070 Braga, Portugal*

Abstract

Abstract algebra has the power to *unify* seemingly disparate theories once they are encoded into the same abstract formalism. This paper shows how a relation-algebraic rendering of both database dependency theory and Hoare programming logic purports one such unification, in spite of the latter being an algorithmic theory and the former a data theory.

The approach equips relational data with *functional types* and an associated type system which is useful for database operation type checking and optimization.

The prospect of a generic, unified approach to both programming and data theories on top of libraries already available in automated deduction systems is envisaged.

Keywords: Relational mathematics, program calculi, data dependency theory

“Hardly anybody confronted with practical problems knows how to apply relational calculi [for which] there is almost no broadly available computer support [...]. We feel, however, that the situation is about to change dramatically as relational mathematics develops and computer power exceeds previous expectations.”

Gunther Schmidt (2010)

1. Introduction

In a paper addressing the influence of Alfred Tarski (1901-83) in computer science, Solomon Feferman (2006) quotes the following statement by his colleague John Etchemendy: “*You see those big shiny Oracle towers on Highway 101? They would never have been built without Tarski’s work on the recursive definitions of satisfaction and truth.*”

The ‘big shiny Oracle towers’ are nothing but the headquarters of Oracle Corporation, the giant database software provider sited in the San Francisco Peninsula. Still Feferman (2006): “*Does Larry Ellison know who Tarski is or anything about his work? [...] I learned subsequently from Jan Van den Bussche that [...] he marks the reading of Codd’s seminal paper as the starting point leading to the Oracle Corporation.*”

Bussche (2001) had in fact devoted attention to relating Codd and Tarski’s work: “*We conclude that Tarski produced two alternatives for Codd’s relational algebra: cylindric set algebra, and relational algebra with pairing [...]. For example, we can represent the ternary relation $\{(a, b, c), (d, e, f)\}$ as $\{(a, (b, c)), (d, (e, f))\}$ ”.*

Still Bussche (2001):

“Using such representations, we leave it as an exercise to the reader to simulate Codd’s relational algebra in RA^+ [relational algebra with pairing].”

To the best of the author’s knowledge, nobody has thus far addressed this *exercise* in a thorough and generic way. Instead, standard relational database theory (Maier, 1983; Abiteboul et al., 1995) includes a

Email address: jno@di.uminho.pt (José N. Oliveira)

well-known relation algebra but this is worked out in set theory and quantified logic, far from the objectives of Tarski’s life-long pursuit in developing methods for elimination of quantifiers from logic expressions. This effort ultimately lead to his *formalization of set theory without variables* (Tarski and Givant, 1987).

The topic has acquired recent interest with the advent of work on implementing extensions of Tarski’s algebra in automated deduction systems such as *Isabelle* (Struth, 2011) or *Prover9* and the associated counterexample generator *Mace4* (Höfner and Struth, 2007). This offers a potential for automation which has not been acknowledged by the database community. In this context, it is worth mentioning an early concern of the founding fathers of the standard theory (Beeri et al., 1977):

“[A] general theory that ties together dependencies, relations and operations on relations is still lacking”.

More than 30 years later, this concern is still justified, as database programming standards remain insensitive to techniques such as formal verification and extended static checking (Flanagan et al., 2002) which are regarded more and more essential to ensuring quality in complex software systems.

Contribution. The remainder of this paper will show how an algebraic treatment of standard data dependency theory along the *exercise* proposed by Bussche equips relational data with *functional types* and an associated type system which can be used to *type check* database operations.

Interestingly, such a *typed* approach to database programming will be shown to relate to other programming logics such as eg. *Hoare logic* (Hoare, 1969) or that of *strongest invariant functions* (Mili et al., 1985) which has been used in the analysis of while statements, for instance.

On the whole, the approach has a *unifying theories of programming* (Hoare and Jifeng, 1998) flavour, even though the exercise will not be carried out in “canonical” UTP.

Outline. Section 2 introduces functional dependencies (FD) and shows how to convert the standard definition into the Tarskian, quantifier-free style. The parallel between the *functions as types* approach which emerges from such a conversion and a similar treatment of Hoare logic starts in section 3. Section 4 shows that, in essence, *injectivity* is what matters in FDs and gives a correspondingly simpler definition of FD which is used in section 5 to re-factor the standard theory into a *type system of FDs*. Section 6 shows how to use this type system to type check database operations and section 7 shows how to calculate query optimizations from FDs. The last sections conclude and give an account of related and future work.

Some technical details are omitted from the current paper for conciseness. All can be found in a technical report available on-line (Oliveira, 2011).

2. Introducing functional dependencies

In standard relational data processing, real life objects or entities are recorded by assigning values to their observable properties or *attributes*. A database *table* is a collection of such attribute assignments, one per object, such that all values of a particular attribute (say i) are of the same type (say A_i). For n such attributes, a *relational database file* T can be regarded as a set of n -tuples, that is, $T \subseteq A_1 \times \dots \times A_n$. A *relational database* is just a collection of several such relations, or tables.

Attribute names normally replace natural numbers in the identification of attributes. The enumeration of all attribute names in a database table, for instance $S = \{\text{PILOT, FLIGHT, DATE, DEPARTS}\}$ concerning the airline scheduling system given as example in (Maier, 1983), is a finite set called the table’s *scheme*. This scheme captures the *syntax* of the data. What about *semantics*? Even non-experts in airline scheduling will accept “business rules” such as, for instance: *a single pilot is assigned to a given flight, on a given date*. This restriction is an example of a so-called *functional dependency* (FD) among attributes, which can be stated more formally by writing “FLIGHT DATE \rightarrow PILOT” to mean that *attribute PILOT is functionally dependent on FLIGHT and DATE*, or that FLIGHT, DATE *functionally determine* PILOT.

Data dependencies help in capturing the *meaning* of relational data. Data dependency theory involves not only functional dependencies (FD) but also multi-valued dependencies (MVD). Both are central to the standard theory, where they are addressed in an axiomatic way. Maier (1983) provides the following definition for FD-satisfiability:

Definition 1. Given subsets $x, y \subseteq S$ of the relation scheme S of a n -ary relation T , this relation is said to satisfy functional dependency $x \rightarrow y$ iff all pairs of tuples $t, t' \in T$ which “agree” on x also “agree” on y , that is,

$$\forall t, t' : t, t' \in T \Rightarrow (t[x] = t'[x] \Rightarrow t[y] = t'[y]) \quad (1)$$

(The notation $t[x]$ in (1) means “the values exhibited by tuple t for the attributes in x ”.)

□

How does one express formula (1) in Tarski’s relation algebra style, without the two-dimensional universal quantification and logical implications inside? For so doing we need to settle some notation. To begin with, $t[x]$ is better written as $x(t)$, where x is identified with the *projection function* associated to attribute set x . Regarding x and y in (1) as such functions we write:

$$\forall t, t' : t, t' \in T \Rightarrow (x(t) = x(t') \Rightarrow y(t) = y(t')) \quad (2)$$

Next, we observe that, given a function $f : A \rightarrow B$, the binary relation $R \subseteq A \times A$ which checks whether two values of A have the same image under f ¹ — that is, $a' R a \equiv f(a') = f(a)$ — can be written alternatively as $a'(f^\circ \cdot f)a$. Here, f° denotes the *converse* of f — that is, $a(f^\circ)b$ holds iff $b = f(a)$ — and the dot (\cdot) denotes the extension of function composition to binary relations²:

$$b(R \cdot S)c \equiv \exists a : b R a \wedge a S c \quad (3)$$

Using converse and composition the rightmost implication of (2) can be rewritten into $t(x^\circ \cdot x)t' \Rightarrow t(y^\circ \cdot y)t'$, for all $t, t' \in T$. Implications such as this can be expressed as relation inclusions, following the definition

$$R \subseteq S \equiv \forall b, a : b R a \Rightarrow b S a \quad (4)$$

However, just stating the inclusion $x^\circ \cdot x \subseteq y^\circ \cdot y$ would be a gross error, for the double scope of the quantification ($t \in T \wedge t' \in T$) would not be taken into account. To handle this, we first unnest the two implications of (2),

$$\forall t, t' : (t \in T \wedge t' \in T \wedge t(x^\circ \cdot x)t') \Rightarrow t(y^\circ \cdot y)t'$$

and treat the antecedent $t \in T \wedge t' \in T \wedge t(x^\circ \cdot x)t'$ independently, by replacing the set of tuples T by the binary relation $\llbracket T \rrbracket$ defined as follows³:

$$b\llbracket T \rrbracket a \equiv b = a \wedge a \in T \quad (5)$$

Note that $t \in T$ can be expressed in terms of $\llbracket T \rrbracket$ by $\exists u : t\llbracket T \rrbracket u$ and similarly for $t' \in T$. Then:

$$\begin{aligned} & (t \in T \wedge t' \in T \wedge t(x^\circ \cdot x)t') \\ \equiv & \{ \text{expansion of } t \in T \text{ and } t' \in T \} \\ & \exists u, u' : t\llbracket T \rrbracket u \wedge t'\llbracket T \rrbracket u' \wedge t(x^\circ \cdot x)t' \\ \equiv & \{ \wedge \text{ is commutative; } u = t \text{ and } u' = t'; \text{ converse} \} \\ & \exists u, u' : t\llbracket T \rrbracket u \wedge u(x^\circ \cdot x)u' \wedge u'\llbracket T \rrbracket^\circ t' \\ \equiv & \{ \text{composition (3) twice} \} \end{aligned}$$

¹This is known as the *nucleus* (Mili et al., 1985) or *kernel* (Oliveira, 2009) of a function f .

²Thus composition of both functions are relations should be read backwards. This is consistent with $b f a$ (function f regarded as a special case of relation) meaning $b = f(a)$ and not $a = f(b)$.

³This is a standard way of encoding a set T as a *partial identity* (Schmidt, 2010), thus called since $\llbracket T \rrbracket \subseteq id$. The set of all such relations forms a Boolean algebra which reproduces the usual algebra of sets. Moreover, partial identities are symmetric ($\llbracket T \rrbracket^\circ = \llbracket T \rrbracket$) and such that $\llbracket S \rrbracket \cdot \llbracket T \rrbracket = \llbracket S \rrbracket \cap \llbracket T \rrbracket$. Also known as *coreflexives* (Freyd and Scedrov, 1990) or as *monotypes* (Doornbos et al., 1997), partial identities are special cases of *tests* in Kleene algebras with tests (Kozen, 1997).

$$t(\llbracket T \rrbracket \cdot x^\circ \cdot x \cdot \llbracket T \rrbracket^\circ) t'$$

□

Finally, by putting this together with $t(y^\circ \cdot y) t'$ we obtain

$$\llbracket T \rrbracket \cdot x^\circ \cdot x \cdot \llbracket T \rrbracket^\circ \subseteq y^\circ \cdot y \quad (6)$$

as a quantifier-free relation algebra expression meaning the same as (1).

Generalization. To reassure the reader worried about the doubtful practicality of derivations such as the above, we would like to say that we don't need to do it over and over again: inequality (6), our Tarskian alternative to the original textbook definition (1), is all we need for calculating with functional dependencies. Moreover, we can start this by actually expanding the scope of the definition from sets of tuples $\llbracket T \rrbracket$ and attribute functions (x, y) to arbitrary binary relations R and suitably typed functions f and g :

$$R \cdot f^\circ \cdot f \cdot R^\circ \subseteq g^\circ \cdot g \quad (7)$$

In this wider setting, R can be regarded not only as a piece of data but also as the specification of a non-deterministic computation, or even the transition relation of a finite-state automaton; and f (resp. g) as a function which observes the input (resp. output) of R . Put back into quantified logic, such a wider notion of a functional dependency will expand as follows:

$$\forall a', a : f(a') = f(a) \Rightarrow (\forall b', b : b' R a' \wedge b R a \Rightarrow g(b') = g(b)) \quad (8)$$

In words: *inputs a, a' to R which are indistinguishable by f can only lead to outputs indistinguishable by g .*

Notationally, we will convey this interpretation by writing $R : f \rightarrow g$ or $f \xrightarrow{R} g$. We can still say that R satisfies FD $f \rightarrow g$, in particular wherever R is a piece of data. As can be easily checked, $f(a') = f(a)$ is an equivalence relation which, in the wider setting, can be regarded as the *semantics* of the datatype which R takes inputs from (think of $f : A \rightarrow B$ as a *semantic* function mapping a syntactic domain A into a semantic domain B), and similarly for g concerning the output type.

Summing up, the functions f and g in (7) can be regarded as *types* for R . Some type assertions of this kind will be very easy to check, for instance $id : f \rightarrow f$, just by replacing $R, f, g := id, f, f$ in (7) and simplifying. But type inference will be easier to calculate on top of the even simpler (re)statement of (7) which is given next.

3. Functions as types

Before proceeding let us record two properties of the relational operators *converse* and *composition*⁴:

$$(R \cdot S)^\circ = S^\circ \cdot R^\circ \quad (9)$$

$$(R^\circ)^\circ = R \quad (10)$$

Moreover, it will be convenient to have a name for the relation $R^\circ \cdot R$ which, for R a function f , is the equivalence relation “indistinguishable by f ” seen above. We define

$$\ker R \triangleq R^\circ \cdot R \quad (11)$$

and read $\ker R$ as “the *kernel* of R ”. Clearly, $a'(\ker R)a$ means $\exists b : b R a' \wedge b R a$ and therefore $\ker R$ measures the *injectivity* of R : the larger it is the larger the set of inputs which R is unable to distinguish (= the less *injective* R is).

⁴It may help to recall the same properties from elementary linear algebra, once converse is interpreted as *matrix transposition* and composition as *matrix-matrix multiplication* (Schmidt, 2010).

We capture this by introducing a preorder on relations which compares their *injectivity*:

$$R \leq S \quad \triangleq \quad \ker S \subseteq \ker R \quad (12)$$

As an example, take two list functions: *elems* computing the set of all elements of a list and *bagify* keeping the bag of such elements. The former loses more information (order and multiplicity) than the latter, which only forgets about order. Thus $\textit{elems} \leq \textit{bagify}$. A function f (relation in general) will be *injective* iff $\ker f \subseteq \textit{id}$ ($\textit{id} \leq f$), which easily converts to the usual definition: $f(a') = f(a) \Rightarrow a' = a$.

Summing up: for functions or any totally defined relations R and S ⁵, $R \leq S$ means that R is *less or as injective as* S ; for possibly partial R and S , it will mean that R is *less injective or more defined than* S . Therefore, for *total* relations R the preorder is universally bounded, $! \leq R \leq \textit{id}$, where the infimum is captured by constant function $!$ which maps every argument to a given (predefined) value, the choice of which is irrelevant⁶. The kernel of $!$ is therefore the largest possible, denoted by \top (for “top”). The other bound is trivial to check, since $\ker \textit{id} = \textit{id}$, this arising from the well-known fact that \textit{id} is the unit of composition. In general, $\textit{id} \leq R$ means that R is injective.

Equipped with this ordering, we may spruce up our relational characterization of the $f \xrightarrow{R} g$ type assertion, or functional dependency (FD):

$$\begin{aligned} & f \xrightarrow{R} g \\ \equiv & \quad \{ \text{definition (7)} \} \\ & R \cdot f^\circ \cdot f \cdot R^\circ \subseteq g^\circ \cdot g \\ \equiv & \quad \{ \text{converses (9,10); kernel (11)} \} \\ & \ker(f \cdot R^\circ) \subseteq \ker g \\ \equiv & \quad \{ (12): g \text{ is “less or as injective as } f \text{ wrt. } R” \} \\ & g \leq f \cdot R^\circ \end{aligned}$$

□

We thus reach a rather compact formula for expressing functional dependencies, whose layout invites us to actually swap the direction of the arrow notation (but, of course, this is optional and just a matter of taste):

Definition 2. Given an arbitrary binary relation $R \subseteq A \times B$ and functions $f : B \rightarrow D$ and $g : A \rightarrow C$, the “type assertion” $g \xleftarrow{R} f$ meaning that R satisfies FD $f \rightarrow g$ is captured by the equivalence:

$$g \xleftarrow{R} f \quad \equiv \quad g \leq f \cdot R^\circ \quad (13)$$

□

There are two main advantages in definition (13), besides saving ink. The most important is that it profits from the relational calculus of injectivity which will be addressed in the following section. The other is that it makes it easy to bridge with other programming logics, as is seen next.

Parallel with Hoare logic. As is widely known, Hoare logic is based on triples of the form $\{p\}R\{q\}$, with the standard interpretation: “if the assertion p is true before initiation of a program R , then the assertion q will be true on its completion” (Hoare, 1969).

⁵A relation R is totally defined (or *entire*) iff $\textit{id} \subseteq \ker R$.

⁶Thus $! \cdot f = !$, for all f . Also note that $R \leq S$ is a preorder, not a partial order, meaning that two relations indistinguishable with respect to their degree of injectivity can be different.

Let program R be identified with the relation which captures its state transition semantics and predicates p (and q) be identified with relation $s'[[p]]s \equiv s' = s \wedge p(s)$ (similarly for q) in which the reader identifies the earlier trick of converting sets to partial identities (section 2). Note how $[[p]]$ can be regarded as the semantics of a statement which checks $p(s)$ and does not change state, failing otherwise. In relation algebra the Hoare triple is captured by

$$\{p\}R\{q\} \equiv \text{rng}(R \cdot [[p]]) \subseteq [[q]] \quad (14)$$

meaning that the outputs of R (given by the range operator rng) for inputs pre-conditioned by p fall inside q ⁷; that is, q is *weaker* than the strongest (*liberal*) post-condition $\text{slp}(R, p)$, something we can express by writing

$$\{p\}R\{q\} \equiv q \leq p \cdot R^\circ \quad (15)$$

under a suitable preorder \leq expressing that q is less constrained than $p \cdot R^\circ$:⁸

$$R \leq S \equiv \text{dom } S \subseteq \text{dom } R \quad (16)$$

In spite of the different semantic context, there is a striking formal similarity between formulas (15) and (13) suggesting that Hoare logic and the logic we want to build for FDs share the same mathematics once expressed in relation algebra. Such similarities will become apparent in the sequel, particularly whereupon we write $p \xrightarrow{R} q$ (or the equivalent $q \xleftarrow{R} p$) for $\{p\}R\{q\}$ to put the two notations closer to each other. In this way, rules such as eg. that of composition, $\{p\}R\{q\} \wedge \{q\}S\{r\} \Rightarrow \{p\}R; S\{r\}$ become reminiscent of labelled transition systems:⁹

$$p \xrightarrow{R} q \wedge q \xrightarrow{S} r \Rightarrow p \xrightarrow{R;S} r \quad (17)$$

We will check the FD equivalent to composition rule (17) shortly.

4. A calculus of injectivity (\leq)

One of the advantages of relation algebra is its easy “tuning” to special needs, which we will illustrate below concerning the algebra of injectivity. We give just an example, taken from (Oliveira, 2011); the reader is referred to this report for technical details.

We start by considering two rules of relation algebra which prove very useful in program calculation:

$$f \cdot R \subseteq S \equiv R \subseteq f^\circ \cdot S \quad (18)$$

$$R \cdot f^\circ \subseteq S \equiv R \subseteq S \cdot f \quad (19)$$

In these equivalences¹⁰, which are widely known as *shunting rules* (Bird and de Moor, 1997; Schmidt, 2010), f is required to be a (total) function. In essence, they let one trade a function f from one side to the other of a \subseteq -equation just by taking converses. (This is akin to “changing sign” in trading terms in inequations of elementary algebra.)

It would be useful to have similar rules for the injectivity preorder, which we have chosen as support for our definition of a FD (13). Such rules turn out to be quite easy to infer, as is the case of the following Galois connection for trading a function f with respect to injectivity

$$R \cdot f \leq S \equiv R \leq S \cdot f^\circ \quad (20)$$

⁷See eg. (Kozen, 2000). Term $\text{rng}(R \cdot [[p]])$ instantiates the semiring *diamond* combinator of (Desharnais et al., 2006). Wehrman et al. (2009) give an even simpler semantics for Hoare triples: $P\{R\}Q \equiv R \cdot P \subseteq Q$, that is, P is at most the *weakest pre-specification* (residual relation) $R \setminus Q$, where $b(R \setminus Q)a$ means $\forall c : c R b \Rightarrow c Q a$ (Hoare and He, 1987).

⁸Details: by definition, $\text{dom } R = \ker R \cap \text{id}$ and $\text{rng } R = \text{dom}(R^\circ)$ (converse duality). Starting from (14), triple $\{p\}R\{q\}$ asserts $\text{rng}(R \cdot [[p]]) \subseteq [[q]]$, itself the same as $\text{dom}([[[p]] \cdot R^\circ]) \subseteq \text{dom}[[q]]$ (15) by converse duality and the fact that the domain of a partial identity is itself. Parentheses $[[_]]$ are dropped for improved readability.

⁹Forward composition $R; S$ means the same as $S \cdot R$.

¹⁰Technically, these equivalences should be regarded as (families of) Galois connections.

calculated as follows:

$$\begin{aligned}
& R \cdot f \leq S \\
\equiv & \quad \{ \text{definition (12); converses (9,10); kernel (11)} \} \\
& \ker S \subseteq f^\circ \cdot (\ker R) \cdot f \\
\equiv & \quad \{ \text{shunting rules (18,19)} \} \\
& f \cdot \ker S \cdot f^\circ \subseteq \ker R \\
\equiv & \quad \{ \text{converses, kernel and definition (12) again} \} \\
& R \leq S \cdot f^\circ
\end{aligned}$$

□

Below we put shunting rule (20) at work in the derivation of a *trading*-rule which will enable handling composite antecedent and consequent functions in FDs:

$$g \xleftarrow{h \cdot R \cdot k^\circ} f \equiv g \cdot h \xleftarrow{R} f \cdot k \quad (21)$$

Thanks to (20), the calculation of (21) is immediate:

$$\begin{aligned}
& g \xleftarrow{h \cdot R \cdot k^\circ} f \\
\equiv & \quad \{ \text{definition (13); converses} \} \\
& g \leq f \cdot k \cdot R^\circ \cdot h^\circ \\
\equiv & \quad \{ \text{shunting rule (20)} \} \\
& g \cdot h \leq (f \cdot k) \cdot R^\circ \\
\equiv & \quad \{ \text{definition (13)} \} \\
& g \cdot h \xleftarrow{R} f \cdot k
\end{aligned}$$

□

Another result about relational injectivity which will help in the sequel is

$$X \leq R \cup S \equiv X \leq R \wedge X \leq S \wedge R^\circ \cdot S \subseteq \ker X \quad (22)$$

where $R \cup S$ is the union of relations R and S . For $X := id$, (22) tells that $R \cup S$ is injective *iff* both R and S are injective and *don't "equivocate" each other*: wherever bSa and bRc hold, $c = a$. The calculation of (22) follows:

$$\begin{aligned}
& X \leq R \cup S \\
\equiv & \quad \{ \text{definitions of } \leq \text{ (12) and kernel (11)} \} \\
& (R \cup S)^\circ \cdot (R \cup S) \subseteq \ker X \\
\equiv & \quad \{ \text{distribution of converse and composition over union} \} \\
& (R^\circ \cdot R) \cup (R^\circ \cdot S) \cup (S^\circ \cdot R) \cup (S^\circ \cdot S) \subseteq \ker X \\
\equiv & \quad \{ \text{kernel (11)} \} \\
& \ker R \cup (R^\circ \cdot S) \cup (S^\circ \cdot R) \cup \ker S \subseteq \ker X \\
\equiv & \quad \{ \text{universal property: } R \cup S \subseteq X \equiv R \subseteq X \wedge S \subseteq X; \text{ (12)} \}
\end{aligned}$$

$$\begin{aligned}
& X \leq R \wedge R^\circ \cdot S \subseteq \ker X \wedge S^\circ \cdot R \subseteq \ker X \wedge X \leq S \\
\equiv & \quad \{ \text{the intermediate conjuncts are the same (taking converses)} \} \\
& X \leq R \wedge R^\circ \cdot S \subseteq \ker X \wedge X \leq S
\end{aligned}$$

□

Hoare logic counterparts. Galois connection (20) holds with no further change once \leq is replaced by the preorder adopted for Hoare triples (16), the reasoning being the same. Fact (22) is even simpler for such a preorder, as the third conjunct disappears.¹¹ Finally, the Hoare logic counterpart of (21) is

$$q \xleftarrow{h \cdot R \cdot k^\circ} p \equiv wp(h, q) \xleftarrow{R} wp(k, p) \quad (23)$$

where $wp(f, p) = \text{dom}(\llbracket p \rrbracket \cdot f)$ denotes the weakest-precondition for *function* f to ensure p on the output.¹² The first steps of the proof of (23) are the same as those of (21), leading to $q \cdot h \leq p \cdot k \cdot R^\circ$ (abbreviating $\llbracket p \rrbracket$, $\llbracket q \rrbracket$ to p , q). But the calculation requires further reasoning in this case because predicates p , q are (relationally) partial identities (tests) and therefore are not at the same level as functions. Since domain is idempotent, dom can be added to any of R or S in $R \leq S$ and thus dom is a self-adjoint concerning (16):

$$\text{dom } R \leq S \equiv R \leq \text{dom } S \quad (24)$$

Then:

$$\begin{aligned}
& q \cdot h \leq \text{dom}(p \cdot k \cdot R^\circ) \\
\equiv & \quad \{ \text{domain: } \text{dom}(R \cdot S) = \text{dom}(\text{dom } R \cdot S) \} \\
& q \cdot h \leq \text{dom}(\text{dom}(p \cdot k) \cdot R^\circ) \\
\equiv & \quad \{ (24) \} \\
& \text{dom}(q \cdot h) \leq \text{dom}(p \cdot k) \cdot R^\circ \\
\equiv & \quad \{ \text{weakest precondition for functions: } wp(f, p) = \text{dom}(p \cdot f) \} \\
& wp(h, q) \leq wp(k, p) \cdot R^\circ \\
\equiv & \quad \{ \text{Hoare triple (15)} \} \\
& wp(h, q) \xleftarrow{R} wp(k, p)
\end{aligned}$$

□

5. Building a type system of FDs

The machinery set up in the previous sections is enough for developing a type system whereby *dependencies, relations and operations on relations are tied together*, as envisaged by Beeri et al. (1977).

¹¹This happens because dom distributes through union, while \ker does not. Both versions of (20) are instances of a generic result concerning *Galois connection lifting*, see appendix D.6 of (Oliveira, 2011).

¹²Terms such as $h \cdot R \cdot k^\circ$ denote programs which begin by reversing a function, proceeding as R and then updating the state by another function; for instance, program $\{x := x-1; R; x := 2*x\}$ on a single-variable state x is denoted by relational term $(2*) \cdot R \cdot (1+)^\circ$, for some subprogram R .

Composition rule. FDs on relations which matching antecedent and consequent functions (as types) compose:

$$y \xleftarrow{S \cdot R} x \iff y \xleftarrow{S} z \quad \wedge \quad z \xleftarrow{R} x \quad (25)$$

Proof:

$$\begin{aligned} & h \xleftarrow{S} g \quad \wedge \quad g \xleftarrow{R} f \\ \equiv & \quad \{ (13) \text{ twice} \} \\ & h \leq g \cdot S^\circ \quad \wedge \quad g \leq f \cdot R^\circ \\ \Rightarrow & \quad \{ \leq\text{-monotonicity of } (\cdot S^\circ) ; \text{converse (9)} \} \\ & h \leq g \cdot S^\circ \quad \wedge \quad g \cdot S^\circ \leq f \cdot (S \cdot R)^\circ \\ \Rightarrow & \quad \{ \leq\text{-transitivity} \} \\ & h \leq f \cdot (S \cdot R)^\circ \\ \equiv & \quad \{ (13) \text{ again} \} \\ & h \xleftarrow{S \cdot R} f \end{aligned}$$

□

For R and S the same database table (tuple set), this rule subsumes Armstrong axiom F5 (Transitivity) in the standard FD theory (Maier, 1983). For R and S regarded as describing computations, rule (25) is the FD counterpart of the rule of composition in Hoare logic, recall (17).¹³

Consequence (weakening/strengthening) rule.

$$k \xleftarrow{R} h \iff k \leq g \quad \wedge \quad g \xleftarrow{R} f \quad \wedge \quad f \leq h \quad (26)$$

Proof: See (Oliveira, 2011), where this rule is shown to subsume and generalize standard Armstrong axioms F2 (*Augmentation*) and F4 (*Projectivity*). In the parallel with Hoare logic, it corresponds to the two *rules of consequence* (Hoare, 1969) which, put together and writing triples as arrows, becomes

$$q' \xleftarrow{R} p' \iff q' \Leftarrow q \quad \wedge \quad q \xleftarrow{R} p \quad \wedge \quad p \Leftarrow p'$$

for a program R and assertions p, q, p', q' . (Note that implications $q' \Leftarrow q$ etc correspond to $\llbracket q \rrbracket \subseteq \llbracket q' \rrbracket$ and therefore to $q' \leq q$ once brackets $\llbracket \rrbracket$ are dropped for simplicity.)

Reflexivity. We have seen already that

$$f \xleftarrow{id} f \quad (27)$$

holds trivially. This rule, which corresponds to the “*skip*” rule of Hoare logic, $p \xleftarrow{skip} p$, is easily shown to hold for any set T ,

$$f \xleftarrow{\llbracket T \rrbracket} f \quad (28)$$

as FDs are downward closed (that is, preserved by sub-relations, by monotonicity). Rule (28) is known as Armstrong axiom F1 (*Reflexivity*).

Note in passing that (25) and (27) together define a *category* whose objects are functions (types) and whose morphisms (arrows) are FDs.

¹³The proof is the same, as both (12) and (16) are preorders (thus transitive) compatible with relational composition. Recall that $R; S = S \cdot R$.

6. Type checking database operations

Merging (union). Let us proceed to an example of database operation *type checking*: we want to know what it means for the merging of two database files to satisfy a particular functional dependency $f \longrightarrow g$. That is, we want to find a *sufficient* condition for the union $R \cup S$ of two relations R and S to be of type $f \longrightarrow g$. The relational algebra of injectivity does most of the work:

$$\begin{aligned}
& g \xleftarrow{R \cup S} f \\
\equiv & \quad \{ \text{definition (13); converse distributes by union} \} \\
& g \leq f \cdot (R^\circ \cup S^\circ) \\
\equiv & \quad \{ \text{relational composition distributes through union} \} \\
& g \leq f \cdot R^\circ \cup f \cdot S^\circ \\
\equiv & \quad \{ \text{algebra of injectivity (22); definition (13) again, twice} \} \\
& g \xleftarrow{R} f \quad \wedge \quad g \xleftarrow{S} f \quad \wedge \quad R \cdot \ker f \cdot S^\circ \subseteq \ker g \\
\equiv & \quad \{ \text{introduce “mutual dependency” shorthand} \} \\
& g \xleftarrow{R} f \quad \wedge \quad g \xleftarrow{S} f \quad \wedge \quad g \xleftarrow{R,S} f
\end{aligned}$$

□

The “mutual dependency” shorthand $g \xleftarrow{R,S} f$ introduced in the last step for $R \cdot \ker f \cdot S^\circ \subseteq \ker g$ can be read as a generalization of the standard definition of FD to *two* relations instead of *one* — just generalize the second R in (8) to some S . For R and S two sets of tuples, it means that grabbing one tuple from one set and another tuple from the other set, if they cannot be distinguished by f then they will remain indistinguishable by g .

It should be stressed that the bottom line of the calculation expresses not only a *sufficient* but also a *necessary* condition for $g \xleftarrow{R \cup S} f$ to hold, as all steps are equivalences. Summing up, rule

$$g \xleftarrow{R \cup S} f \quad \equiv \quad g \xleftarrow{R} f \quad \wedge \quad g \xleftarrow{S} f \quad \wedge \quad g \xleftarrow{R,S} f \quad (29)$$

holds.¹⁴

Type checking other database operations will follow the same scheme. Below we handle in detail one particular such operation, *relational join* (Maier, 1983). This is justified not only for its relevance in data processing but also because it brings about other standard FD rules not yet addressed.

Joining (pairing). Recall from section 1 how Bussche (2001) explains the relevance of Tarski’s work on *pairing* in relation algebra by illustrating how a ternary (in general, n -ary) relation $\{(a, b, c), (d, e, f)\}$ gets represented by a binary one, $\{(a, (b, c)), (d, (e, f))\}$.

Pairing is not only useful for ensuring that sets of arbitrarily long (but finite) tuples are representable by binary relations but also for defining the *join* operator (\triangle) on such sets. This operator turns out to be particularly handy to formalize in case the two sets of tuples are already represented as relations R and S :

$$(a, b)(R \triangle S)c \quad \equiv \quad a R c \wedge b S c \quad (30)$$

¹⁴The counterpart of (29) in Hoare logic is $\{p\}R\{q\} \wedge \{p\}S\{q\} \equiv \{p\}(R \cup S)\{q\}$ written directly in the original triple notation, where $R \cup S$ denotes the non-deterministic choice between R and S .

Interestingly, relational join behaves as a *least upper bound* with respect to the injectivity preorder:

$$R \triangle S \leq T \quad \equiv \quad R \leq T \wedge S \leq T \quad (31)$$

This arises from fact¹⁵

$$\ker(R \triangle S) = \ker R \cap \ker S \quad (32)$$

as follows:

$$\begin{aligned} & R \triangle S \leq T \\ \equiv & \quad \{ (12) \text{ and } (32) \} \\ & \ker T \subseteq (\ker R) \cap (\ker S) \\ \equiv & \quad \{ \text{universal property: } X \subseteq R \cap S \equiv X \subseteq R \wedge X \subseteq S \} \\ & \ker T \subseteq \ker R \quad \wedge \quad \ker T \subseteq \ker S \\ \equiv & \quad \{ (12) \text{ twice } \} \\ & R \leq T \quad \wedge \quad S \leq T \end{aligned}$$

□

This combinator, termed *split* in (Bird and de Moor, 1997), *fork* in (Frias et al., 1997) and *strict fork* in Schmidt (2010), turns out to be more general than its use in data processing suggests. In particular, when R and S are functions f and g , $f \triangle g$ is the obvious function which pairs the outputs of f and g : $(f \triangle g)x = (f(x), g(x))$. Think for instance of the projection function f_x (resp. f_y) which, in the context of Definition 1 yields $t[x]$ (resp. $t[y]$) when applied to a tuple t . Then $(f_x \triangle f_y)t = (t[x], t[y]) = t[xy]$, where xy denotes the union of attributes x and y (Maier, 1983). So, attribute union corresponds to joining the corresponding projection functions. This gives us a quite uniform framework for handling both relational join and compound attributes. To make notation closer to what is common in data dependency theory we will abbreviate $f_x \triangle f_y$ to $f_x f_y$ and this even further to xy , identifying (as we did before) each attribute (say x) with the corresponding projection function (say f_x).

Keeping abbreviation fg of $f \triangle g$ (for functions), from (31) it is easy to derive facts $! \leq f \leq id$, $f \leq fg$ and $g \leq fg$. This is consistent with the use of juxtaposition to denote “sets of attributes”. Likewise, \leq can be regarded as expressing “attribute inclusion” in this context: the more attributes one observes the more injective the projection function corresponding to such attributes is¹⁶.

A first illustration of this unified framework is given below: the (generic) calculation of the so-called Armstrong axioms F3 (*Additivity*) and F4 (*Projectivity*)¹⁷. This is done in one go, for arbitrary (suitably typed) R, f, g, h ¹⁸:

$$gh \xleftarrow{R} f \quad \equiv \quad g \xleftarrow{R} f \quad \wedge \quad h \xleftarrow{R} f \quad (33)$$

Calculation:

$$\begin{aligned} & gh \xleftarrow{R} f \\ \equiv & \quad \{ (13); \text{ expansion of shorthand } gh \} \end{aligned}$$

¹⁵Fact (32) follows immediately from $(R \triangle S)^\circ \cdot (X \triangle Y) = R^\circ \cdot X \cap S^\circ \cdot Y$ (Bird and de Moor, 1997).

¹⁶Note how $!$ mimics the empty set and id mimics the whole set of attributes, enabling one to “see the whole thing” and thus discriminating as much as possible.

¹⁷See (Maier, 1983).

¹⁸In the Hoare logic counterpart of this rule, gh stands for the product $g \times h$ of predicates g and h defined by $(g \times h)(b, a) = g(b) \wedge h(a)$. The rule, derived in (Oliveira, 2009), ensures that the category of FDs has products.

$$\begin{aligned}
& g \triangle h \leq f \cdot R^\circ \\
\equiv & \quad \{ \text{universal property of } \triangle \text{ (31)} \} \\
& g \leq f \cdot R^\circ \quad \wedge \quad h \leq f \cdot R^\circ \\
\equiv & \quad \{ \text{(13) twice} \} \\
& g \xleftarrow{R} f \quad \wedge \quad h \xleftarrow{R} f
\end{aligned}$$

□

The typing rule for the join $R \triangle S$ of two relations R and S is calculated in the same way. The reasoning involves properties of the projection functions $\pi_1(x, y) = x$ and $\pi_2(x, y) = y$:

$$\begin{aligned}
& g \xleftarrow{R} f \quad \wedge \quad h \xleftarrow{S} f \\
\Rightarrow & \quad \{ \pi_1 \cdot (R \triangle S) \subseteq R \text{ and } \pi_2 \cdot (R \triangle S) \subseteq S; \text{ FDs are downward closed} \} \\
& g \xleftarrow{\pi_1 \cdot (R \triangle S)} f \quad \wedge \quad h \xleftarrow{\pi_2 \cdot (R \triangle S)} f \\
\equiv & \quad \{ \text{trading (21) twice} \} \\
& g \cdot \pi_1 \xleftarrow{R \triangle S} f \quad \wedge \quad h \cdot \pi_2 \xleftarrow{R \triangle S} f \\
\equiv & \quad \{ \text{F3+F4 (33)} \} \\
& (g \cdot \pi_1) \triangle (h \cdot \pi_2) \xleftarrow{R \triangle S} f \\
\equiv & \quad \{ \text{product of functions: } f \times g = (f \cdot \pi_1) \triangle (g \cdot \pi_2) \} \\
& g \times h \xleftarrow{R \triangle S} f
\end{aligned}$$

□

7. Beyond the type system: database operation optimization

As explained above, FD theory (resp. Hoare logic) can be regarded as a type system whose rules help in reasoning about data models (resp. programs) without going into the semantic intricacies of data business rules (resp. program meanings).

When compared to the quantified expression of Definition 1, quantifier-free equivalent (13) looks simpler and is therefore expected to be easier to use in practice. This section gives two illustrations of this, one concerned with query optimization and the other with optimizing lexicographic sorting of database files.

FDs for query optimization. This example, taken from Abiteboul et al. (1995), is also addressed by Wisnesky (2011): one wants to optimize the conjunctive query

$$\{(d, a') \mid (t, d, a) \in \text{Movies}, (t', d', a') \in \text{Movies}, t = t'\} \quad (34)$$

over a database file $\text{Movies}(\text{Title}, \text{Director}, \text{Actor})$ into a query accessing this file only once, knowing that FD $\text{Title} \rightarrow \text{Director}$ holds.

Using abbreviations M , t , d and a for (respectively) Movies , Title , Director and Actor , we want to solve for X the equation

$$d \cdot M \cdot (\ker t) \cdot M \cdot a^\circ = X \quad (35)$$

— whose left hand side is the relational equivalent of (34)¹⁹ — aiming at a solution X containing only one instance of M . The equation is solved by taking FD $d \xleftarrow{M} t$ itself as starting point and trying to re-write it into something one recognizes as an instance of (35):

$$\begin{aligned}
& d \xleftarrow{M} t \\
\equiv & \quad \{ (13) \} \\
& d \leq t \cdot M^\circ \\
\equiv & \quad \{ \text{expanding (11,12); } M^\circ = M \text{ since } M \text{ is a partial identity} \} \\
& M \cdot t^\circ \cdot t \cdot M \subseteq d^\circ \cdot d \\
\equiv & \quad \{ \text{composition } (\cdot M) \text{ with a partial identity (Oliveira, 2011)} \} \\
& M \cdot t^\circ \cdot t \cdot M \subseteq d^\circ \cdot d \cdot M \\
\Rightarrow & \quad \{ \text{shunting (18,19); monotonicity of } (\cdot a^\circ); \text{ kernel (11)} \} \\
& d \cdot M \cdot (\ker t) \cdot M \cdot a^\circ \subseteq d \cdot M \cdot a^\circ
\end{aligned}$$

□

We thus find $d \cdot M \cdot a^\circ$ as a candidate solution for X . To obtain $X = d \cdot M \cdot a^\circ$ it remains to check the converse inclusion:

$$\begin{aligned}
& d \cdot M \cdot a^\circ \subseteq d \cdot M \cdot (\ker t) \cdot M \cdot a^\circ \\
\Leftarrow & \quad \{ id \subseteq \ker t \text{ because kernels of functions are equivalence relations} \} \\
& d \cdot M \cdot a^\circ \subseteq d \cdot M \cdot M \cdot a^\circ \\
\equiv & \quad \{ M \cdot M = M \cap M = M \text{ because } M \text{ is a partial identity} \} \\
& d \cdot M \cdot a^\circ \subseteq d \cdot M \cdot a^\circ
\end{aligned}$$

□

Altogether, FD $d \xleftarrow{M} t$ grants the solution $X = d \cdot M \cdot a^\circ$ to equation (35) — that is

$$X = \{(d, a') \mid (t, d, a') \in \text{Movies}\}$$

— which optimizes the given query by only visiting the movies file once²⁰.

Optimizing lexicographic sorting. This example calculates an improvement in lexicographic sorting of database files subject to FDs. Let \leq and \sqsubseteq be two preorders of the same type. By the expression $\leq \triangleright \sqsubseteq$ we mean the lexicographic order

$$a(\leq \triangleright \sqsubseteq)a' \equiv a \leq a' \wedge (a \geq a' \Rightarrow a \sqsubseteq a')$$

which gives priority to \leq , that is,

$$\leq \triangleright \sqsubseteq = \leq \cap (\leq^\circ \Rightarrow \sqsubseteq) \tag{36}$$

¹⁹ As the interested reader may check by introducing the variables back. Note how $\ker t$ expresses $t = t'$ and projection functions d (for *Director*) and a (for *Actor*) work over tuple (t, d, a) and tuple (t', d', a') , respectively. The use of the same letters for data variables and the corresponding projection functions should help in comparing the two versions of the query.

²⁰ By the way: symmetry between a and d in calculation step $d \cdot M \cdot t^\circ \cdot t \cdot M \cdot a^\circ \subseteq d \cdot M \cdot a^\circ$ above immediately tells that FD $a \xleftarrow{M} t$ would also enable the proposed optimization.

where relational implication is the upper adjoint of intersection:

$$R \cap S \subseteq X \equiv R \subseteq (S \Rightarrow X) \quad (37)$$

Now suppose that T is a database file whose schema includes attribute x (resp. y) whose domain is ordered by partial order \leq_x (resp. \leq_y). Thus the tuples of T can be ordered not only by the preorders

$$t \leq_a^T t' \equiv t, t' \in T \wedge a(t) \leq_a a(t') \quad (38)$$

for $a \in \{x, y\}$, but also by lexicographic combinations thereof, eg. $\leq_x^T \triangleright \leq_y^T$. However, such lexicographic preorders can be simplified in presence of FDs. Below we calculate a sufficient condition for such a lexicographic preorder to reduce to one of its components, for instance:

$$\leq_x^T \triangleright \leq_y^T = \leq_x^T \iff y \xleftarrow{T} x \quad (39)$$

The relation-algebraic calculation of rule (39) goes in the same style as before:²¹

$$\begin{aligned} & \leq_x^T \triangleright \leq_y^T = \leq_x^T \\ \equiv & \quad \{ (36); X \cap Y = X \text{ equivalent to } X \subseteq Y \} \\ & \leq_x^T \subseteq ((\leq_x^T)^\circ \Rightarrow \leq_y^T) \\ \equiv & \quad \{ \text{Galois connection (37)} \} \\ & \leq_x^T \cap (\leq_x^T)^\circ \subseteq \leq_y^T \\ \equiv & \quad \{ \text{variable-free versions of (38) for } a \in \{x, y\}; \text{converses} \} \\ & T \cdot x^\circ \cdot \leq_x \cdot x \cdot T \cap T \cdot x^\circ \cdot \leq_x^\circ \cdot x \cdot T \subseteq T \cdot y^\circ \cdot \leq_y \cdot y \cdot T \\ \equiv & \quad \{ \text{distributions over intersection, as } x \cdot T \text{ is univalent; converses} \} \\ & T \cdot x^\circ \cdot (\leq_x \cap \leq_x^\circ) \cdot x \cdot T \subseteq T \cdot y^\circ \cdot \leq_y \cdot y \cdot T \\ \equiv & \quad \{ \leq_x \text{ is antisymmetric: } \leq_x \cap \leq_x^\circ = id \} \\ & T \cdot x^\circ \cdot x \cdot T \subseteq y^\circ \cdot \leq_y \cdot y \\ \leftarrow & \quad \{ \leq_y \text{ is reflexive} \} \\ & T \cdot x^\circ \cdot x \cdot T \subseteq y^\circ \cdot y \\ \equiv & \quad \{ (6); (13) \} \\ & y \xleftarrow{T} x \end{aligned}$$

□

²¹The fourth step in the reasoning relies on prop. 5.3 of (Schmidt, 2010): for univalent Q , distribution law $(R \cap S) \cdot Q = (R \cdot Q) \cap (S \cdot Q)$ holds — and therefore (taking converses) so does $Q^\circ \cdot (R \cap S) = (Q^\circ \cdot R) \cap (Q^\circ \cdot S)$.

8. Conclusions

“The great merit of algebra is as a powerful tool for exploring family relationships over a wide range of different theories. (...) It is only their algebraic properties that emphasise the family likenesses (...) Algebraic proofs by term rewriting are the most promising way in which computers can assist in the process of reliable design.”

Hoare and Jifeng (1998)

There is a growing interest in algebraic reasoning in computer science able to eventually promote calculational techniques in software engineering, hopefully unifying seemingly disparate theories once they are encoded into the same abstractions. Relation algebra (Schmidt, 2010) is particularly apt in this respect.

The current paper shows how a relation-algebraic rendering of both data dependency theory and Hoare logic purports one such unification, in spite of the latter being an algorithmic theory and the former a data theory, thanks to both algorithms and data structures being expressed in the *unified language of binary relations*.

In short (and informally), both logics rely on triples: *something* (a data set; a program) lies between an *antecedent* and a *consequent* observation (a data attribute; a state assertion); there is an ordering (injectivity; definition) on observations; triples express that antecedent observations are “enough” for consequents to hold “modally through” what is in between.

Triples are nicely captured by *arrows*, whose end-points can be regarded as *types*. On the data side, our approach equips relational data with *functional types* and an associated type system which can be used to *type check* database operations and optimize queries by calculation once they are written as Tarskian, quantifier-free formulas.

As formal verification is becoming more and more widespread to ensure quality in complex software systems, we believe our approach may contribute to *unified* formal verification tools blending in the same framework extended static checking and database programming.

Back to the opening story, surely Tarski’s work on satisfaction and truth is relevant to computer science. But Etchemendy’s answer could have been better tuned to the particular context of database technology suggested by the Oracle towers landscape:

[...] *“They would never have been built without Tarski’s work on the calculus of binary relations.”*

9. Related and future work

Functional dependencies have been characterized relationally by checking the determinism of the relations obtained by projecting tuple sets by antecedent and consequent attributes (Schmidt and Ströhlein, 1993; Jaoua et al., 1997). This alternative definition is equivalent to the one followed in the current paper.²² As a generalization, Jaoua et al. (1997) also study so-called *difunctional* dependencies.

Dependencies in relational databases have also been expressed using so-called *indiscernibility relations* (Okuma et al., 2003). Freyd and Scedrov (1990) develop a τ -category theory of relations based on *monic n-tuples*. Concepts such as *table*, *column*, *short column* etc. fit into the spirit of (pointfree) data dependency and database theory and should be carefully studied in the context of the current paper.

Wisnesky (2011) addresses the semantic optimization of monad comprehensions in functional programming by generalizing results from relational database theory. As this theory relies on the *powerset* monad, whose comprehensions correspond to database queries, by handling similar optimizations in relation algebra (as we did in section 7) we have followed the well-known shift towards the Kleisli adjoint category.

As is well-known, this shift can be generalized to any other monad. Wisnesky (2011) includes queries on probabilistic databases, this time relying on the (finite support) *distribution* monad. As shown by Oliveira (2012), the “Kleisli shift” wrt. this monad leads to typed *linear algebra*. Wong and Butz (2000) introduce Bayesian embedded multivalued dependencies as necessary and sufficient conditions for lossless

²²See section *Generic relational projections* in (Oliveira, 2011).

decomposition of probabilistic relations. Lossless decomposition and multivalued dependencies have been handled by Oliveira (2011) in the same way as FDs in the current paper. The prospect of calculating with data dependencies in probabilistic systems *directly* in matrix algebra is an interesting prospect for future work, in line with the consistent use of matrix notation by Schmidt (2010) in relation algebra. Whether this carries over to probabilistic Hoare logic (Barthe et al., 2012) remains to be seen.

Other ways of relating data dependency theory with algorithmic reasoning can be devised. For instance, Mili et al. (1985) reason about while-loops $w = (\mathbf{while} \ t \ \mathbf{do} \ b)$ in terms of so-called *strongest invariant functions*, where invariant functions f , ordered by injectivity, are such that $f \cdot \llbracket t \rrbracket = f \cdot b \cdot \llbracket t \rrbracket$ holds. A simple argument in relation algebra shows this equivalent to $f \cdot b \cdot \llbracket t \rrbracket \subseteq f$, thus entailing FD $f \xleftarrow{b \cdot \llbracket t \rrbracket} f$. How much of our FD relation-algebraic approach could be applied in this setting is open to research. This includes finding a meaningful counterpart of the *rule of iteration* (Hoare, 1969) at data level, a topic not addressed in the current paper.

Another law not considered in the correspondence between Hoare logic and functional dependencies is the *axiom of assignment*, $\{p[e/x]\}x := e\{p\}$ where $p[e/x]$ denotes the predicate which is obtained from p by replacing all occurrences of x by e . This axiom is interesting because it relies on the state structure of imperative programs: variables which hold data. Assignment $x := e$ means *selective updating*: program variable x is updated to e . A possible data-level counterpart to such selective updating is the SQL **update** command, which is of the form **UPDATE** R **SET** f **WHERE** x , meaning: *update all tuples in R which satisfy selection criterion x by tuple-transformation f , leaving the rest unchanged*. While the semantics of this operation is easy to encode in relation algebra, the parallel is somewhat artificial and needs further analysis. In general, future work should identify which generic properties of the \leq relation on types are common to both frameworks and derive a more general kernel theory which both are instances of.

Last but not least, another prospect for future work concerns automated reasoning. RelView (Berghammer, 2013) is a well-known system that calculates with relations “beyond toy size”. Many applications of relation algebra have been handled successfully in this tool. Algebraic structures such as idempotent semirings and Kleene algebras (which relation algebra is an instance of) have also been shown to be amenable to automation by eg. Höfner and Struth (2007) and Struth (2011). Möller et al. (2012) encode a database preference theory into idempotent semiring algebra and show how to use Prover9 to discharge proofs. Model checking of extended static checks in tools such as eg. the Alloy Analyser also blends well with algebraic-relational models (Oliveira and Ferreira, 2013).

The implementation of a generic, unified approach to both data and program theories on top of libraries already available in such automated deduction systems is a prospect for long term research.

Acknowledgments. This work is funded by ERDF - European Regional Development Fund through the COMPETE Programme (operational programme for competitiveness) and by National Funds through the *FCT - Fundação para a Ciência e a Tecnologia* (Portuguese Foundation for Science and Technology) within project FCOMP-01-0124-FEDER-028923.

The author would like to thank the anonymous referees for insightful comments which significantly improved the quality of the original submission. Feedback and exchange of ideas with Ryan Wisnesky about pointfree query reasoning are also gratefully acknowledged.

References

- G. Schmidt, Relational Mathematics, number 132 in Encyclopedia of Mathematics and its Applications, Cambridge University Press, 2010.
- S. Feferman, Tarski’s influence on computer science, Logical Methods in Computer Science 2 (2006) 1–13.
- J. Bussche, Applications of Alfred Tarski’s ideas in database theory, in: CSL’01: Proceedings of the 15th International Workshop on Computer Science Logic, Springer-Verlag, London, UK, 2001, pp. 20–37.
- D. Maier, The Theory of Relational Databases, Computer Science Press, 1983.
- S. Abiteboul, R. Hull, V. Vianu, Foundations of Databases, Addison-Wesley, 1995.
- A. Tarski, S. Givant, A Formalization of Set Theory without Variables, American Mathematical Society, 1987. AMS Colloquium Publications, volume 41, Providence, Rhode Island.
- G. Struth, Isabelle repository for relational and algebraic methods, 2011. URL: <http://staffwww.dcs.shef.ac.uk/people/G.Struth/isa>.

- P. Höfner, G. Struth, Automated reasoning in Kleene algebra, in: Proceedings, CADE-21, Springer-Verlag, 2007, pp. 279–294.
- C. Beeri, R. Fagin, J. Howard, A complete axiomatization for functional and multivalued dependencies in database relations, in: D. Smith (Ed.), Proc. 1977 ACM SIGMOD, Toronto, ACM, NY, USA, 1977, pp. 47–61.
- C. Flanagan, K. Leino, M. Lillibridge, G. Nelson, J. Saxe, R. Stata, Extended static checking for Java, in: PLDI, 2002, pp. 234–245.
- C. Hoare, An axiomatic basis for computer programming, *CACM* 12,10 (1969) 576–580, 583.
- A. Mili, J. Desharnais, J. Gagné, Strongest invariant functions: Their use in the systematic analysis of while statements, *Acta Inf.* 22 (1985) 47–66.
- C. Hoare, H. Jifeng, Unifying Theories of Programming, Series in Computer Science, Prentice-Hall International, 1998.
- J. Oliveira, Pointfree Foundations for (Generic) Lossless Decomposition, Technical Report TR-HASLab:3:2011, HASLab/INESC TEC & U. Minho, 2011. URL: <https://repositorium.sdum.uminho.pt/handle/1822/24648>.
- J. Oliveira, *Extended Static Checking by Calculation using the Pointfree Transform*, volume 5520 of *LNCS*, Springer-Verlag, 2009, pp. 195–251.
- P. Freyd, A. Scedrov, Categories, Allegories, volume 39 of *Mathematical Library*, North-Holland, 1990.
- H. Doornbos, R. Backhouse, J. van der Woude, A calculational approach to mathematical induction, *Theoretical Computer Science* 179 (1997) 103–135.
- D. Kozen, Kleene algebra with tests, *ACM Trans. Program. Lang. Syst.* 19 (1997) 427–443.
- D. Kozen, On Hoare logic and Kleene algebra with tests, *ACM Trans. Comput. Logic* 1 (2000) 60–76.
- J. Desharnais, B. Möller, G. Struth, Kleene algebra with domain, *ACM Trans. Comput. Logic* 7 (2006) 798–833.
- I. Wehrman, C. Hoare, P. W. O’Hearn, Graphical models of separation logic, *Information Processing Letters* 109 (2009) 1001–1004.
- C. Hoare, J. He, The weakest prespecification, *Information Processing Letters* 24 (1987) 127–132.
- R. Bird, O. de Moor, Algebra of Programming, Series in Computer Science, Prentice-Hall International, 1997.
- M. Frias, G. Baum, A. Haeberer, Fork algebras in algebra, logic and computer science, *Fundam. Inform.* (1997) 1–25.
- R. Wisnesky, Minimizing Monad Comprehensions, Technical Report TR-02-11, Harvard University, Cambridge, Massachusetts, 2011.
- G. Schmidt, T. Ströhlein, Relations and Graphs: Discrete Mathematics for Computer Scientists, EATCS Monographs on Theoretical Computer Science, Springer-Verlag, 1993.
- A. Jaoua, N. Belkhit, H. Ounalli, T. Moukam, Databases, in: C. Brink, W. Kahl, G. Schmidt (Eds.), Relational methods in computer science, Springer-Verlag New York, Inc., New York, NY, USA, 1997, pp. 197–210.
- H. Okuma, W. MacCaull, Y. Kawahara, Informational representability for contexts in Dedekind categories, Technical Report trcs208, Kyushu University, 2003. URL: <http://www.i.kyushu-u.ac.jp/doi/tr/trcs208.ps.gz>.
- J. Oliveira, Towards a linear algebra of programming, *Formal Aspects of Computing* 24 (2012) 433–458.
- S. Wong, C. Butz, The implication of probabilistic conditional independence and embedded multivalued dependency, in: IPMU00, 2000, pp. 876–881. 8th Conf. on Inf. Processing and Management of Uncertainty in K.-B. Systems.
- G. Barthe, B. Grégoire, S. Béguelin, Probabilistic relational Hoare logics for computer-aided security proofs, in: MPC’12, 2012, pp. 1–6.
- R. Berghammer, Computing and visualizing Banks sets of dominance relations using relation algebra and RelView, *JLAP* 82 (2013) 123–136.
- B. Möller, P. Roocks, M. Endres, An algebraic calculus of database preferences, in: MPC 2012, volume 7342 of *LNCS*, Springer Berlin Heidelberg, 2012, pp. 241–262.
- J. Oliveira, M. Ferreira, Alloy meets the algebra of programming: A case study, *IEEE Transactions on Software Engineering* 39 (2013) 305–326.