

Incomplete Information in Relational Databases

TOMASZ IMIELIŃSKI AND WITOLD LIPSKI, JR.

Polish Academy of Sciences, Warsaw, Poland

Abstract. This paper concerns the semantics of Codd's relational model of data. Formulated are precise conditions that should be satisfied in a semantically meaningful extension of the usual relational operators, such as projection, selection, union, and join, from operators on relations to operators on tables with "null values" of various kinds allowed. These conditions require that the system be safe in the sense that no incorrect conclusion is derivable by using a specified subset Ω of the relational operators; and that it be complete in the sense that all valid conclusions expressible by relational expressions using operators in Ω are in fact derivable in this system. Two such systems of practical interest are shown. The first, based on the usual Codd's null values, supports projection and selection. The second, based on many different ("marked") null values or variables allowed to appear in a table, is shown to correctly support projection, positive selection (with no negation occurring in the selection condition), union, and renaming of attributes, which allows for processing arbitrary conjunctive queries. A very desirable property enjoyed by this system is that all relational operators on tables are performed in exactly the same way as in the case of the usual relations. A third system, mainly of theoretical interest, supporting projection, selection, union, join, and renaming, is also discussed. Under a so-called closed world assumption, it can also handle the operator of difference. It is based on a device called a conditional table and is crucial to the proof of the correctness of the second system. All systems considered allow for relational expressions containing arbitrarily many different relation symbols, and no form of the universal relation assumption is required.

Categories and Subject Descriptors: H.2.3 [Database Management]: Languages—*query languages*; H.2.4 [Database Management]: Systems—*query processing*

General Terms: Theory

Additional Key Words and Phrases: Relational database, incomplete information, null values, marked nulls, relational algebra, query language semantics, query processing

1. Introduction

Attempts to represent incomplete information in the relational model of data have been made since the very beginning of the relational database theory [3, 7, 8, 10, 12, 13, 18–21, 27, 28, 30, 31, 33, 34]. (See also the extensive bibliography in [22].) The main device in this context has been the *null value*, a special symbol @

This work was supported in part by the Polish Academy of Sciences under Contract MR.I.3. The second author was also aided through a visiting appointment by Istituto di Analisi dei Sistemi ed Informatica, Consiglio Nazionale delle Ricerche, Roma, Italy.

Some of the results of this paper were presented at the Seventh International Conference on Very Large Data Bases, Cannes, France, September 1981.

Authors' present addresses: T. Imieliński, Rutgers University, Department of Computer Science, Hill Center for the Mathematical Sciences, Busch Campus, New Brunswick, NJ, 08903; W. Lipski, Jr., Laboratoire de Recherche en Informatique, Université de Paris-Sud, Centre d'Orsay, Bât. 490, 91405 Orsay Cedex, France.

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

© 1984 ACM 0004-5411/84/1000-0761 \$00.75

allowed as an entry in a table, with its intended meaning being *value at present unknown (but the attribute applicable)*.

Representing incomplete information in a database immediately raises the much more difficult and important problem of processing this information so that the user can obtain—perhaps partial, but correct—responses to his or her queries on the basis of the incomplete information. In the context of the relational model, this comes down to defining the relational operators on tables with null values in a semantically correct way. The intuitive condition that should be satisfied is the following. If a null value has a specified semantic interpretation, that is, if we assume that a specified relation exists between a table with nulls and the real world, then this relation should be similar for the tables that are arguments of a relational operator and for the table obtained as the result.

Solutions to this problem proposed so far have been in many respects partial and unsatisfactory. Codd's approach [7, 8] is based on a three-valued logic and a so-called null substitution principle, which has been criticized on semantic grounds by Grant [10] and Lipski [20]. Lipski [19–21] and Vassiliou [30] consider, essentially, only the relational operator of selection, where the selection condition can be an arbitrary Boolean combination of atomic conditions. Proposals involving a richer subset of the relational operators—and it is the join that seems to be the main source of difficulty—usually do not give a clear and precise explanation of the sense in which the definitions of relational operators on tables with nulls are semantically correct [7, 8, 18]. An exception is the work of Biskup [3]. However, his definition of, say, a join $T \bowtie U$ of two tables with nulls works correctly when T and U are “independent” but does not necessarily provide a method for correctly evaluating an expression $f(T) \bowtie g(T)$ where f and g are some relational expressions. This is because in the latter case there may be “additional information” implied by the fact that both arguments of the join come from the same table T , and his definition of join is not able to take this into account (this is discussed in more detail at the end of Section 4.) In fact, queries in a relational database are arbitrary relational expressions, so that we should be able to handle correctly relational expressions rather than just single relational operators.

In this paper we formulate the precise conditions that we believe should be satisfied in any semantically meaningful extension of the usual relational operators, such as projection, selection, union, and join, from operators on relations to operators on “incompletely specified relations”—typically, tables with “null values” of different kinds allowed as entries. These conditions are embodied into the definition of a representation system (see Section 3).

Roughly speaking, these conditions require that our system be safe in the sense that no incorrect conclusion is derivable by using a specified subset Ω of the relational operators, and that it be complete in the sense that all valid conclusions expressible by relational operators using operators in Ω are in fact derivable. The intuition behind the notion of a representation system can also be explained in the following way. By performing relational operators over tables with nulls, we may introduce some corruption of information. However, if the query language that the user has at his or her disposal is weak enough, that is, if only a specified subset Ω of the relational operators is available, then this corruption is not visible by the user.

Two representation systems of practical interest are shown. The first, described in Section 4, is based on tables with the usual Codd null values @ (we call them Codd tables). The following is an example of a Codd table.

SUPPLIER	LOCATION	PRODUCT	QUANTITY
Smith	London	Nails	@
Brown	@	Bolts	@
Jones	@	Nuts	40,000

This table represents the information that Smith, located in London, supplies (an unknown quantity of) nails; Brown supplies bolts; and Jones supplies 40,000 nuts. The system based on Codd tables supports projection and selection, and it is shown that no representation system based on Codd tables can support projection and join at a time.

Section 5 describes some algorithms for correctly evaluating a selection over a Codd table, including the case in which attribute domains are linearly ordered, and the selection condition is an arbitrary Boolean combination of elementary inequalities between attributes or between attributes and values.

The second representation system, presented in Section 6, supports projection, positive selection (with no negation occurring in the selection condition), union, join, and renaming of attributes—this allows for processing arbitrary conjunctive queries [4], as shown in Section 8. It is based on tables allowing many different (“marked”) null values, or variables (such tables are called *V*-tables). The following is an example of a *V*-table:

COURSE	TEACHER	WEEKDAY
Databases	<i>x</i>	Monday
Programming	<i>y</i>	Tuesday
Databases	<i>x</i>	Thursday
FORTRAN	Smith	<i>z</i>

This table contains three different variables, *x*, *y*, *z*. Notice that it represents the information that the teacher of the course on databases, though unknown, is the same on Monday and Thursday, a fact not representable by means of a Codd table.

A very desirable property enjoyed by our representation system based on *V*-tables is that all relational operators on *V*-tables are performed in exactly the same way as in the case of the usual relations, treating variables as if they were regular values in appropriate attribute domains (this may be referred to as the *naive evaluation*). It may also be noted that *V*-tables appear in a natural way in the context of updating relational views. Assume, for instance, that a database is given by

SL :	<table style="border-collapse: collapse;"> <tr> <th style="border: none;">SUPPLIER</th> <th style="border: none;">LOCATION</th> </tr> <tr> <td style="border: none;">Smith</td> <td style="border: none;">London</td> </tr> </table>	SUPPLIER	LOCATION	Smith	London
SUPPLIER	LOCATION				
Smith	London				

SP :	<table style="border-collapse: collapse;"> <tr> <th style="border: none;">SUPPLIER</th> <th style="border: none;">PRODUCT</th> </tr> <tr> <td style="border: none;">Smith</td> <td style="border: none;">Nails</td> </tr> </table>	SUPPLIER	PRODUCT	Smith	Nails
SUPPLIER	PRODUCT				
Smith	Nails				

and suppose that the relational view

$$\pi_{\text{LOCATION, PRODUCT}}(\text{SL} \bowtie \text{SP})$$

is updated by adding two tuples, (New York, Bolts) and (Los Angeles, Nuts). It

is then natural to define the effect of this update to be

SL :	<table style="border-collapse: collapse; width: 100%;"> <thead> <tr> <th style="padding: 5px;">SUPPLIER</th> <th style="padding: 5px;">LOCATION</th> </tr> </thead> <tbody> <tr> <td style="padding: 5px;">Smith</td> <td style="padding: 5px;">London</td> </tr> <tr> <td style="padding: 5px;">x</td> <td style="padding: 5px;">New York</td> </tr> <tr> <td style="padding: 5px;">y</td> <td style="padding: 5px;">Los Angeles</td> </tr> </tbody> </table>	SUPPLIER	LOCATION	Smith	London	x	New York	y	Los Angeles
SUPPLIER	LOCATION								
Smith	London								
x	New York								
y	Los Angeles								

SP :	<table style="border-collapse: collapse; width: 100%;"> <thead> <tr> <th style="padding: 5px;">SUPPLIER</th> <th style="padding: 5px;">PRODUCT</th> </tr> </thead> <tbody> <tr> <td style="padding: 5px;">Smith</td> <td style="padding: 5px;">Nails</td> </tr> <tr> <td style="padding: 5px;">x</td> <td style="padding: 5px;">Bolts</td> </tr> <tr> <td style="padding: 5px;">y</td> <td style="padding: 5px;">Nuts</td> </tr> </tbody> </table>	SUPPLIER	PRODUCT	Smith	Nails	x	Bolts	y	Nuts
SUPPLIER	PRODUCT								
Smith	Nails								
x	Bolts								
y	Nuts								

A negative result shown in Section 6 is that, rather surprisingly, no representation system based on V -tables can support projection and arbitrary selection (recall that these operators can be supported by Codd tables).

A third representation system, described in Section 7, supports projection, selection, union, join, and renaming and is mainly of a theoretical interest. It is based on an idea of a so-called conditional table and is crucial in the proof of the correctness of the second system. A conditional table is a V -table with an additional column, *con*, containing a condition; for example, the conditional table

SUPPLIER	LOCATION	PRODUCT	con
x	London	Nails	$x = \text{Smith}$
Brown	New York	Nails	$x \neq \text{Smith}$

represents the information that nails are supplied either by Smith in London or by Brown in New York, but not by both at a time.

All representation systems considered in this paper allow for relational expressions containing arbitrarily many different relation symbols, and no form of the universal relation assumption [24] is required.

Most of the results of the paper are developed under a so-called open world assumption (see [26]) that, roughly speaking, means that we are not able to represent negative information. In Section 9 we briefly show how the results can be extended in a straightforward manner to a modified interpretation of tables where negative information is representable. It turns out that under such a modified interpretation the system based on conditional tables can also handle the operator of difference, thus supporting the full strength of the relational algebra.

The last section contains conclusions and briefly describes related work, in particular, the work concerning the problem of handling dependencies in the context of V -tables, and the relation between V -tables and tableaux of Aho, Sagiv, and Ullman [1].

2. Basic Definitions

In this section, we give some basic definitions and notation concerning the relational data model ([5]; see also [29]).

Throughout the paper we consider a fixed, finite (unless otherwise stated) set \mathcal{U} of *attributes*. Attributes are usually denoted by A, B, C , and sets of attributes by X, Y, Z , with possible subscripts. A set of attributes, say $\{A, C\}$, is usually written as AC . Associated with every $A \in \mathcal{U}$ is an *attribute domain* $D(A)$. We always assume that $|D(A)| \geq 2$ and we denote $D = \bigcup_{A \in \mathcal{U}} D(A)$. Elements of D are sometimes called *constants*. Elements of $D(A), D(B), D(C)$ are usually denoted by a, b, c , respectively, with possible primes, etc. By a *tuple* on X we mean any mapping t that associates a value $t(A) \in D(A)$ with every $A \in X$. A tuple is usually denoted as a string of values associated with the attributes; for example, ac is a

tuple on AC . For a tuple t on X and for any $Y \subseteq X$, we denote by $t[Y]$ the restriction of t to Y ; for example, if $t = abc$ then $t[AC] = ac$. By a relation on X we mean any finite set of tuples on X . If t, r are a tuple and a relation on X , we then write

$$\alpha(t) = \alpha(r) = X,$$

and we call $\alpha(t)$ and $\alpha(r)$ the type of t and r , respectively. A multirelation of type $\langle X_1, \dots, X_n \rangle$ is any sequence $\langle r_1, \dots, r_n \rangle$ in which r_i is a relation on X_i , $1 \leq i \leq n$. Relations are usually denoted by r, s , with possible subscripts, and multirelations by \mathbf{r}, \mathbf{s} (conforming to the general convention that boldface indicates multiobjects). The type of multirelation \mathbf{r} is denoted by $\alpha(\mathbf{r})$. By an empty multirelation we mean any multirelation of the form $\mathbf{0} = \langle \emptyset, \dots, \emptyset \rangle$; formally, we shall assume that there is a different empty multirelation for each type.

For two multirelations $\mathbf{r} = \langle r_1, \dots, r_n \rangle, \mathbf{s} = \langle s_1, \dots, s_n \rangle$ of the same type, we write $\mathbf{r} \subseteq \mathbf{s}$ if $r_i \subseteq s_i, 1 \leq i \leq n$. If $\mathbf{r} = \langle r_1, \dots, r_n \rangle$ and $1 \leq i \leq n$, then we define

$$\text{pr}_i(\mathbf{r}) = r_i.$$

The set of all multirelations of type $\langle X_1, \dots, X_n \rangle$ is denoted by $\mathcal{R}(X_1, \dots, X_n)$; in particular, $\mathcal{R}(X)$ is the set of all relations of type X . The set of all multirelations is denoted by \mathcal{R} .

The class of all nonempty homogeneous sets of multirelations is denoted by \mathcal{S} (a set \mathcal{X} of multirelations is homogeneous if $\alpha(\mathbf{r}) = \alpha(\mathbf{s})$ for all $\mathbf{r}, \mathbf{s} \in \mathcal{X}$). Elements of \mathcal{S} are denoted by $\mathcal{X}, \mathcal{Y}, \mathcal{Z}$, and the definition of $\alpha(\mathcal{X})$, the type of \mathcal{X} , is the natural one.

We consider the usual relational operators:

Projection

$$\pi_Y(r) = \{t[Y] : t \in r\} \quad (Y \subseteq \alpha(r)).$$

Selection

$$\sigma_E(r) = \{t \in r : E(t) = \text{true}\}.$$

Here E is a selection condition; that is, any expression built up from atomic conditions of the form $(A = a), (A = B), A, B \in \mathcal{Z}, D(A) = D(B), a \in D(A)$ (and truth constants **true**, **false**) by means of the logical connectives \vee (or), \wedge (and), \neg (not). $E(t)$ denotes the truth value obtained by substituting the value $t(A)$ for any occurrence of A in E and evaluating the expression in the natural way. We always assume that $\alpha(r)$ contains all the attributes occurring in E . A selection condition E and the selection operator σ_E is called *positive* if it does not contain the \neg symbol.

We sometimes consider a more general form of selection, where E is any mapping from tuples to $\{\text{true}, \text{false}\}$. (Note: The operation $\sigma_{A=B}$ is also called *restriction*.)

Union

$$r \cup s, \quad \text{that is, the usual set-theoretical union.}$$

We always assume that both arguments of union are of the same type.

Join

$$r \bowtie s = \{t : \alpha(t) = X \cup Y \wedge t[X] \in r \wedge t[Y] \in s\},$$

where

$$X = \alpha(r), \quad Y = \alpha(s).$$

Difference

$$r - s = r \setminus s, \quad \text{that is, the usual set-theoretical difference.}$$

We always assume that both arguments of difference are of the same type. In Section 8 we also consider the following additional operation.

Renaming (of an attribute)

$$s_B^A(r) = \{s_B^A(t) : t \in r\}$$

where $A \in \alpha(r)$, $B \in \mathcal{Z} \setminus \alpha(r)$, $D(A) = D(B)$, and $s_B^A(t)$ is the tuple of type $(\alpha(r) \setminus \{A\}) \cup \{B\}$ with

$$(s_B^A(t))(C) = \begin{cases} t(C) & \text{if } C \in \alpha(r) \setminus \{A\}, \\ t(A) & \text{if } C = B. \end{cases}$$

Intuitively, $s_B^A(r)$ is the result of renaming column A or r to B .

We show in Section 8 how the results of this paper can be extended to another version of relational algebra, based on Cartesian product (see, e.g., [29]) instead of join.

By a *relation name* we mean a symbol R (the letter S is also used) with possible subscripts, with an associated *type* $\alpha(R) \subseteq \mathcal{Z}$.

An *instance* of R is any relation r such that $\alpha(r) = \alpha(R)$. An instance is usually denoted by a lowercase version of the letter denoting the relation name.

For any subset Ω of the relational operators of projection (P), selection (S), positive selection (S⁺), join (J), difference (D), and renaming (R), a *relational Ω -expression* means any well-formed expression built up from relation names and relational operators in Ω ; for example, $\sigma_{A=a}(R_1) \bowtie \pi_{AB}(R_2)$ is a relational PS⁺J-expression. A *multirelational Ω -expression* is any sequence $\mathbf{f} = \langle f_1, \dots, f_k \rangle$ of relational Ω -expressions. Relational and multirelational Ω -expressions are often called simply Ω -expressions; they are usually denoted by f, g , with possible subscripts, and by \mathbf{f}, \mathbf{g} , respectively. They are assumed to be *typed* in the following sense. We assume that associated with every \mathbf{f} is a sequence $\langle R_1, \dots, R_n \rangle$ (either clear from the context or given explicitly by writing $\mathbf{f}(R_1, \dots, R_n)$) containing all relation names occurring in \mathbf{f} (and possibly some other relation names). The *argument type* of \mathbf{f} is then defined as

$$\alpha(\mathbf{f}) = \langle \alpha(R_1), \dots, \alpha(R_n) \rangle.$$

For any $\mathbf{f} = \langle f_1, \dots, f_k \rangle$ and any multirelation $\mathbf{r} = \langle r_1, \dots, r_n \rangle$ of type $\alpha(\mathbf{r}) = \alpha(\mathbf{f})$, we define $\mathbf{f}(\mathbf{r}) = \langle f_1(\mathbf{r}), \dots, f_k(\mathbf{r}) \rangle$ where $f_i(\mathbf{r})$ denotes the relation obtained by substituting r_j for all occurrences of R_j in f_i , $j = 1, \dots, n$. We also define $\beta(\mathbf{f}) = \alpha(\mathbf{f}(\mathbf{r}))$ and call it the *result type* of \mathbf{f} .

It may be mentioned in this context that the relational algebra can be embedded in a natural way into a simple algebraic system called a *cylindric set algebra*, such that all relations considered are of type \mathcal{Z} and all relational operators are total (i.e., there are no restrictions on the arguments). This approach has several advantages (i.e., the “complicated” join operator becomes the usual set-theoretical intersection) but leads to infinite relations and creates some problems concerning the finite representability of these relations. This approach is treated in more detail in another paper by the authors [14].

3. Representation Systems

In this section we formulate some general conditions that we require to be satisfied in any semantically meaningful extension of the relational algebra to “tables with nulls.”

To this end we introduce the notion of a representation system. By a *representation system* we mean a triple $\langle \mathcal{T}, \text{Rep}, \Omega \rangle$ where \mathcal{T} is the set of *multitables*, Rep is a mapping, $\text{Rep}: \mathcal{T} \rightarrow \mathcal{F}$ (recall that \mathcal{F} is the set of all nonempty homogeneous sets of multirelations), Ω is a set of relational operators, and a certain natural condition, to be defined at the end of this section, is satisfied. Roughly speaking, this condition asserts that there is a way, consistent with respect to operators in Ω , to define $f(T)$ for any Ω -expression f and any multitable $T \in \mathcal{T}$.

First, however, let us give some intuitive explanations concerning the components of a representation system.

A multitable is usually a “generalized multirelation” with tuples allowed to contain not only values in appropriate attribute domains, but also some special symbols (“null values” of various kinds). Most of the notions and notation concerning relations and multirelations, such as a tuple, $\alpha(r)$, or $t[Y]$, carry over to multitables. The set $\text{Rep}(T)$ of multirelations defines the information contained in multitable T ; that is, it specifies the set of possibilities represented by T .

Suppose $\mathcal{X} \in \mathcal{F}$ represents our information about a certain unknown $r^* \in \mathcal{X}$ and let f be an Ω -expression with $\alpha(f) = \alpha(\mathcal{X})$. Clearly, our information about $f(r^*)$ is then represented by

$$f(\mathcal{X}) = \{f(r) : r \in \mathcal{X}\}.$$

In this way, any Ω -expression f can be treated in a natural way as a mapping $f: \mathcal{F} \rightarrow \mathcal{F}$ (more exactly, a partial mapping defined only for those $\mathcal{X} \in \mathcal{F}$ with $\alpha(\mathcal{X}) = \alpha(f)$).

Suppose now that we want to give a natural definition of multitable $f(T)$ where $T \in \mathcal{T}$ and f is an Ω -expression ($\alpha(T) = \alpha(f)$, where $\alpha(T)$ can be defined to be $\alpha(\text{Rep}(T))$). If we think of T as representing in an incomplete way some unknown $r^* \in \text{Rep}(T)$, then ideally we could expect that there exist the same relation between $f(T)$ and $f(r^*)$, so that $f(\text{Rep}(T)) = \text{Rep}(f(T))$; that is, the following diagram commutes:

$$\begin{array}{ccc}
 & \text{Rep} & \\
 \mathcal{T} & \xrightarrow{\quad} & \mathcal{F} \\
 \downarrow f & & \downarrow f \\
 \mathcal{T} & \xrightarrow{\quad} & \mathcal{F} \\
 & \text{Rep} &
 \end{array} \tag{3.1}$$

(for multitables of appropriate type). Unfortunately, this approach is usually not feasible, since in most practical situations the structure of the set $f(\text{Rep}(T))$ is not “regular” enough to be representable by any $U \in \mathcal{T}$ (see, however, an exception at the end of Section 9). In other words, there is no $U \in \mathcal{T}$ such that $\text{Rep}(U) = f(\text{Rep}(T))$.

Clearly, if we are to define $f(T)$ in a semantically meaningful way, then we should require that $\text{Rep}(f(T))$ approximate the information given by $f(\text{Rep}(T))$ in some natural sense (and we should precisely state in which sense). Before we describe what we consider to be a natural notion of “equivalence” between $\text{Rep}(f(T))$ and $f(\text{Rep}(T))$, we give some examples illustrating different kinds of equivalences between sets of multirelations.

Example 3.1. Let \mathcal{X} consists of all those relations that contain at least one of the following two relations:

$$\left[\begin{array}{cc} a & b \\ a' & b' \end{array} \right], \quad \left[\begin{array}{cc} a & b \\ a' & b'' \end{array} \right].$$

(Note: Here, and in all other examples in this paper, different symbols denoting values in attribute domains always stand for different values.)

Let \mathcal{Y} be the set of all relations containing tuple ab (clearly $\mathcal{X} \subsetneq \mathcal{Y}$). In a way, \mathcal{X} and \mathcal{Y} are equivalent: if we assume $\mathbf{r}^* \in \mathcal{X}$, then the only tuple that can be concluded to be in \mathbf{r}^* is ab , and the same is true for \mathcal{Y} .

However, \mathcal{X} provides the information that $a' \in \pi_A(\mathbf{r}^*)$, whereas \mathcal{Y} does not. Hence \mathcal{X} and \mathcal{Y} are not “equivalent with respect to projection.” \square

Example 3.2. Let \mathcal{X} consist of all those relations that contain at least one of the following two relations:

$$\left[\begin{array}{ccc} a & b & c \end{array} \right], \quad \left[\begin{array}{ccc} a & b' & c' \\ a' & b & c' \end{array} \right],$$

and let \mathcal{Y} be defined in the same way by relations

$$\left[\begin{array}{ccc} a & b & c \end{array} \right], \quad \left[\begin{array}{ccc} a & b' & c' \\ a' & b & c'' \end{array} \right].$$

It is easy to see that \mathcal{X} and \mathcal{Y} are equivalent with respect to projection, in the sense of the previous example. However, $\mathbf{r}^* \in \mathcal{X}$ implies

$$ab \in \pi_{AB}(\pi_{AC}(\mathbf{r}^*) \bowtie \pi_{BC}(\mathbf{r}^*)),$$

whereas this conclusion cannot be made under the assumption that $\mathbf{r}^* \in \mathcal{Y}$. Hence, \mathcal{X} and \mathcal{Y} are not equivalent “with respect to PJ-expressions.” \square

Example 3.3. Let $X = ABC$ and let

$$\mathcal{Z} = \{r \in \mathcal{R}(X) : abc \in r \vee ab'c \in r\}.$$

Define

$$\begin{aligned} \mathcal{X} &= \{\langle r, r \rangle : r \in \mathcal{Z}\}, \\ \mathcal{Y} &= \{\langle r, s \rangle : r, s \in \mathcal{Z}\}. \end{aligned}$$

It is not difficult to prove that \mathcal{X} and \mathcal{Y} are equivalent with respect to all PS-expressions (because PS-expressions involve only unary operators). However, they are not equivalent with respect to PJ-expressions. Indeed, $\langle r^*, s^* \rangle \in \mathcal{X}$ implies

$$ac \in \pi_{AC}(r^* \bowtie s^*)$$

(notice that $r^* \bowtie s^* = r^* \cap s^* = r^*$), whereas this conclusion cannot be made under the assumption that $\langle r^*, s^* \rangle \in \mathcal{Y}$. \square

The above examples suggest the following definition. For a relational multi-expression \mathbf{f} and $\mathcal{X} \in \mathcal{F}$ ($\alpha(\mathbf{f}) = \alpha(\mathcal{X})$) define the \mathbf{f} -information in \mathcal{X} , denoted by $\mathcal{X}^{\mathbf{f}}$, to be

$$\mathcal{X}^{\mathbf{f}} = \bigcap \mathbf{f}(\mathcal{X})$$

(if $\mathbf{f} = \langle f_1, \dots, f_k \rangle$, then $\bigcap \mathbf{f}(\mathcal{X})$ is understood as $\langle \bigcap f_1(\mathcal{X}), \dots, \bigcap f_k(\mathcal{X}) \rangle$). In other words, $\mathcal{X}^{\mathbf{f}}$ is the largest multirelation \mathbf{s} such that $\mathbf{s} \subseteq \mathbf{f}(\mathbf{r}^*)$ for all $\mathbf{r}^* \in \mathcal{X}$. Putting it still another way, if $\mathcal{X}^{\mathbf{f}} = \langle s_1, \dots, s_k \rangle$, then $t \in s_i$ means that from $\mathbf{r}^* \in \mathcal{X}$ we may conclude that $t \in f_i(\mathbf{r}^*)$.

Let $\mathcal{X}, \mathcal{Y} \in \mathcal{F}$, $\alpha(\mathcal{X}) = \alpha(\mathcal{Y})$. We say that \mathcal{X} and \mathcal{Y} are Ω -equivalent (in symbols $\mathcal{X} =_{\Omega} \mathcal{Y}$) if $\mathcal{X}^{\mathbf{f}} = \mathcal{Y}^{\mathbf{f}}$ for any Ω -expression \mathbf{f} with $\alpha(\mathbf{f}) = \alpha(\mathcal{X})$. Note that in Example 3.1, $\mathcal{X} =_{\Omega} \mathcal{Y}$ (we consider expressions of the form $f(R_1, \dots, R_n)$)

= R_i , to be the only relational \mathcal{O} -expressions), but $\mathcal{X} \not\equiv_P \mathcal{Y}$; in Example 3.2, $\mathcal{X} \equiv_P \mathcal{Y}$, but $\mathcal{X} \not\equiv_{PJ} \mathcal{Y}$; in Example 3.3, $\mathcal{X} \equiv_{PS} \mathcal{Y}$, but $\mathcal{X} \not\equiv_{PJ} \mathcal{Y}$.

A multitable \mathbf{T} is said to Ω -represent \mathcal{X} if $\text{Rep}(\mathbf{T}) \equiv_{\Omega} \mathcal{X}$.

We are now ready to complete the definition of a representation system. A triple $\langle \mathcal{T}, \text{Rep}, \Omega \rangle$, where $\mathcal{T}, \text{Rep}, \Omega$ are as described before, is a representation system if, for any Ω -expression f and for any multitable $\mathbf{T} \in \mathcal{T}$ ($\alpha(f) = \alpha(\mathbf{T})$), there is a multitable $\mathbf{U} \in \mathcal{T}$ that Ω -represents $f(\text{Rep}(\mathbf{T}))$ —in other words, if $f(\mathbf{T})$ can be defined for every Ω -expression f and multitable $\mathbf{T} \in \mathcal{T}$ ($\alpha(f) = \alpha(\mathbf{T})$) in such a way that

$$\text{Rep}(f(\mathbf{T})) \equiv_{\Omega} f(\text{Rep}(\mathbf{T})). \tag{3.2}$$

Two multitables \mathbf{T}, \mathbf{U} are called Ω -equivalent if $\text{Rep}(\mathbf{T}) \equiv_{\Omega} \text{Rep}(\mathbf{U})$ and *Rep-equivalent* if $\text{Rep}(\mathbf{T}) = \text{Rep}(\mathbf{U})$. Notice that for a given \mathbf{T} and f there may exist many Ω -equivalent multitables $\mathbf{U} \in \mathcal{T}$ that Ω -represent $f(\text{Rep}(\mathbf{T}))$.

There are two more conditions satisfied by most representation systems considered in this paper (these conditions are not part of the formal definition of a representation system). The first one is that any multirelation is—or at least can be identified with—a multitable, so that we may assume $\mathcal{R} \subseteq \mathcal{T}$. This reflects the fact that we are interested only in systems that extend the usual relational algebra. The second assumption (not valid in Section 9) is that for any multitable \mathbf{T} and for any $r, s \in \mathcal{R}$

$$r \in \text{Rep}(\mathbf{T}) \wedge s \supseteq r \Rightarrow s \in \text{Rep}(\mathbf{T}). \tag{3.3}$$

This condition—which is essentially the *open world assumption* of Reiter [26]—is equivalent to saying that we are not able to represent negative information; that is, the knowledge that a relationship expressed by a certain tuple t is definitely *not* true in the real world r^* (or in $f(r^*)$, where f is some relational expression). An intuitive consequence of this assumption is that for any $r \in \mathcal{R}$

$$\text{Rep}(r) = \{s \in \mathcal{R} : s \supseteq r\}.$$

We conclude this section by the following simple lemma.

LEMMA 3.1. *If Ω contains only unary operators, then $\langle \mathcal{T}, \text{Rep}, \Omega \rangle$ is a representation system iff $f(\mathbf{T})$ can be defined for every relational (rather than multirelational) Ω -expression f and every $\mathbf{T} \in \mathcal{T}$ with $\alpha(\mathbf{T}) = \alpha(f)$, in such a way that $\text{Rep}(f(\mathbf{T})) \equiv_{\Omega} f(\text{Rep}(\mathbf{T}))$.*

The easy proof is omitted.

4. Codd Tables

In this section we prove that the usual tables with null values @ (see [7]) provide a basis for a representation system with $\Omega = \text{PS}$. We also show that these tables can support neither $\Omega = \text{PSU}$ nor $\Omega = \text{PJ}$.

By a @-tuple on X we mean any function t that associates a value $t(A) \in D(A) \cup \{\text{@}\}$ with every $A \in X$. A *Codd table* on X is any finite set of @-tuples on X , and a *Codd multitable* of type $\langle X_1, \dots, X_n \rangle$ is any sequence $\mathbf{T} = \langle T_1, \dots, T_n \rangle$ where T_i is a Codd table on X_i , $1 \leq i \leq n$. The set of all Codd multitables is denoted by $\mathcal{T}_{\text{@}}$.

Let us define a partial order \leq on $D \cup \{\text{@}\}$ in such a way that $\text{@} < a$, $a \in D$ are the only nontrivial relationships. For two @-tuples t, t' on X , we write $t \leq t'$

if $t(A) \leq t'(A)$ for every $A \in X$, that is, if t and t' agree on every $A \in X$ such that $t(A) \neq @$. For any Codd table T on X we define

$$\text{Rep}(T) = \{r \in \mathcal{R}(X) : \text{for every } t \in T \text{ there is } t' \in r \text{ such that } t \leq t'\},$$

and for every Codd multitable $\mathbf{T} = \langle T_1, \dots, T_n \rangle$,

$$\text{Rep}(\mathbf{T}) = \text{Rep}(T_1) \times \dots \times \text{Rep}(T_n).$$

Intuitively, $r \in \text{Rep}(\mathbf{T})$ iff r contains a multirelation obtained from \mathbf{T} by replacing all occurrences of $@$ by some values in appropriate attribute domains (different occurrences may be replaced by different values, but equal values are allowed as well).

Let us now have a closer look at the notion of Ω -equivalence of homogeneous sets of multirelations for some simple cases of Ω .

We first note the following trivial fact:

$$\mathcal{X} \equiv_{\emptyset} \mathcal{Y} \Leftrightarrow \bigcap \mathcal{X} = \bigcap \mathcal{Y} \tag{4.1}$$

Let us now consider the case of $\Omega = P$.

THEOREM 4.1. *Any $\mathcal{X} \in \mathcal{F}$ can be P-represented by a Codd multitable.*

PROOF. Let us begin with the simple case where $\alpha(\mathcal{X}) = X$. For any $Y \subseteq X$, let $Q(Y)$ be obtained from $\mathcal{X}^{\ast Y}$ by extending with $@$'s every tuple $t \in \mathcal{X}^{\ast Y}$ to a $@$ -tuple \bar{t} on X (notice that $Q(\emptyset) = \emptyset$ if $\emptyset \in \mathcal{X}$ and $Q(\emptyset)$ consists of a tuple of $@$, otherwise). Define $T = \bigcup_{Y \subseteq X} Q(Y)$. We shall prove that $\text{Rep}(T) \equiv_P \mathcal{X}$. Indeed, consider an arbitrary $Y \subseteq X$. If $t \in \mathcal{X}^{\ast Y}$, then clearly T contains a $@$ -tuple t' that agrees with t on Y , and consequently every $r \in \text{Rep}(T)$ contains a tuple u that agrees with t on Y ; that is, $t \in \text{Rep}(T)^{\ast Y}$. Conversely, assume that $t \notin \mathcal{X}^{\ast Y}$, $\alpha(t) = Y$. Then no $@$ -tuple in T agrees with t on Y , and consequently there exists a relation $r \in \text{Rep}(T)$ not containing any tuple agreeing with t on Y (r may be obtained by replacing every null value in T by a value in the corresponding attribute domain, in such a way that if $A \in Y$ then the value replacing an occurrence of $@$ in column A of T is different from $t(A)$). Hence $t \notin \text{Rep}(T)^{\ast Y}$, which concludes the proof that $\text{Rep}(T) \equiv_P \mathcal{X}$.

In the general case, in which $\alpha(\mathcal{X}) = \langle X_1, \dots, X_n \rangle$, we construct tables T_i such that $\text{Rep}(T_i) \equiv_P \text{pr}_i(\mathcal{X})$, $1 \leq i \leq n$. It is then easily seen that $\mathbf{T} = \langle T_1, \dots, T_n \rangle$ P-represents \mathcal{X} . \square

Let us note the following simple fact concerning the P-equivalence of Codd multitables.

THEOREM 4.2. *Any two P-equivalent Codd multitables are Rep-equivalent.*

In other words, if $P \in \Omega$, then $\mathbf{T} \equiv_{\Omega} \mathbf{U}$ iff $\mathbf{T} \equiv \mathbf{U}$. The easy proof is omitted.

Clearly, two Codd multitables $\langle T_1, \dots, T_n \rangle, \langle U_1, \dots, U_n \rangle$ of the same type are Rep-equivalent iff $T_i \equiv U_i$, $1 \leq i \leq n$. The following structural characterization of Rep-equivalence of Codd tables was proved by Biskup [3]:

THEOREM 4.3. *For any Codd tables T, U of the same type $T \equiv U$ iff for every $t \in T$ there is $u \in U$ such that $t \leq u$ and for every $u \in U$ there is $t \in T$ such that $u \leq t$.*

The S-equivalence is even simpler than the P-equivalence, since it turns out to coincide with the \emptyset -equivalence (see (4.1)). Before we prove this, notice that $\sigma_E(r)$ may be expressed as $r \cap \sigma_E(1_Y)$, where $Y = \alpha(r)$ and $1_Y = \times_{A \in Y} D(A)$, and this is

true for arbitrary conditions $E: \mathbb{1}_Y \rightarrow \{\text{true}, \text{false}\}$, not only for those generated by atomic conditions of the form $(A = a), (A = B)$.

THEOREM 4.4. For any $\mathcal{X}, \mathcal{Y} \in \mathcal{T}$ of the same type,

$$\mathcal{X} \equiv_S \mathcal{Y} \Leftrightarrow \mathcal{X} \equiv_{\emptyset} \mathcal{Y}.$$

PROOF. Let $\mathcal{X}, \mathcal{Y} \in \mathcal{T}$, and let $f(R_1, \dots, R_n) = \sigma_E(R_i)$ be an arbitrary relational S-expression ($\alpha(\mathcal{X}) = \alpha(\mathcal{Y}) = \alpha(f)$). If $\mathcal{X} \equiv_{\emptyset} \mathcal{Y}$ then, by (4.1), $\cap \mathcal{X} = \cap \mathcal{Y}$, which implies $\cap \text{pr}_i(\mathcal{X}) = \cap \text{pr}_i(\mathcal{Y})$. Hence

$$\begin{aligned} \mathcal{X}^f &= \bigcap_{r \in \text{pr}_i(\mathcal{X})} \sigma_E(r) = \bigcap_{r \in \text{pr}_i(\mathcal{Y})} (r \cap \sigma_E(\mathbb{1}_{\alpha(r)})) = \sigma_E(\mathbb{1}_{\alpha(r)}) \cap \bigcap \text{pr}_i(\mathcal{X}) \\ &= \sigma_E(\mathbb{1}_{\alpha(r)}) \cap \bigcap \text{pr}_i(\mathcal{Y}) = \mathcal{Y}^f, \end{aligned}$$

which easily implies $\mathcal{X} \equiv_S \mathcal{Y}$. \square

Notice that the S-equivalence coincides with the \emptyset -equivalence even if we allow arbitrary functions $E: \mathbb{1}_Y \rightarrow \{\text{true}, \text{false}\}$ as selection conditions. Clearly, by Theorem 4.4, $\mathcal{X} \equiv_P \mathcal{Y}$ implies $\mathcal{X} \equiv_S \mathcal{Y}$. However, the PS-equivalence does not coincide with the P-equivalence on \mathcal{T} , as shown by the following example.

Example 4.1. Let

$$\mathcal{X} = \{\overline{abc}, \overline{ab'c}\}, \quad \mathcal{Y} = \{\overline{abc}, \overline{ab'c'}\}.$$

We have $\mathcal{X} \equiv_P \mathcal{Y}$ since both \mathcal{X} and \mathcal{Y} can be P-represented by $T = \overline{a@@}$. Consequently, also $\mathcal{X} \equiv_S \mathcal{Y}$. But $\mathcal{X} \not\equiv_{PS} \mathcal{Y}$, since for the PS-expression

$$f(R) = \pi_A(\sigma_{(B=b) \vee (C=c)}(R))$$

we have

$$\mathcal{X}^f = \{a\} \neq \emptyset = \mathcal{Y}^f.$$

In other words, \mathcal{X} and \mathcal{Y} are not distinguishable by either P-expressions or S-expressions, but are distinguishable by PS-expressions. \square

Before we prove that $(\mathcal{T}_{\emptyset}, \text{Rep}, \text{PS})$ is indeed a representation system, let us note the following two simple facts:

LEMMA 4.1. Every PSUJR-expression f is monotone; that is, $r \subseteq s \Rightarrow f(r) \subseteq f(s)$.

The easy proof, by induction on the number of operators in f , is omitted.

Let $\mathcal{X}, \mathcal{Y} \in \mathcal{T}, \alpha(\mathcal{X}) = \alpha(\mathcal{Y})$. We say that \mathcal{X} and \mathcal{Y} are *coinitial* (in symbols, $\mathcal{X} \approx \mathcal{Y}$) if for any $r \in \mathcal{X}$ there is $s \in \mathcal{Y}$ such that $s \subseteq r$, and for any $s \in \mathcal{Y}$ there is $r \in \mathcal{X}$ such that $r \subseteq s$. Obviously, for any monotone f

$$\mathcal{X} \approx \mathcal{Y} \Rightarrow f(\mathcal{X}) \approx f(\mathcal{Y}). \tag{4.2}$$

LEMMA 4.2. If \mathcal{X} and \mathcal{Y} are coinitial, then $\mathcal{X} \equiv_{PSUJR} \mathcal{Y}$.

PROOF. Let $\mathcal{X} \approx \mathcal{Y}$. Then for every $r \in \mathcal{X}$ there is a $\varphi(r) \in \mathcal{Y}$ such that $\varphi(r) \subseteq r$. By using the previous lemma, for any PSUJR-expression f we have

$$\mathcal{X}^f = \bigcap_{r \in \mathcal{X}} f(r) \supseteq \bigcap_{r \in \mathcal{X}} f(\varphi(r)) \supseteq \bigcap_{s \in \mathcal{Y}} f(s) = \mathcal{Y}^f,$$

and similarly $\mathcal{Y}^f \supseteq \mathcal{X}^f$. \square

We are now ready to prove the main result of this section.

THEOREM 4.5. *It is possible to correctly evaluate PS-expressions over Codd tables; more formally, $\langle \mathcal{T}_\otimes, \text{Rep}, \text{PS} \rangle$ is a representation system.*

PROOF. By Lemma 3.1, it is sufficient to show that for any $\mathbf{T} \in \mathcal{T}_\otimes$ and any relational PS-expression f with $\alpha(f) = \alpha(T)$ it is possible to define $f(\mathbf{T})$ in such a way that $\text{Rep}(f(\mathbf{T})) \equiv_{\text{PS}} f(\text{Rep}(\mathbf{T}))$.

We define $f(\mathbf{T})$ inductively, by using the following rules:

$$\begin{aligned} \text{If } \mathbf{T} &= \langle T_1, \dots, T_n \rangle, \text{ then } \text{pr}_i(\mathbf{T}) = T_i, \\ \pi_Y(T) &= \{t[Y] : t \in T\}, \\ \sigma_E(T) &= \{t \in T : E_*(t) = \text{true}\}, \end{aligned}$$

where

$$E_*(t) = \begin{cases} \text{true} & \text{if } E(u) \text{ for every } u \in \text{Compl}(t), \\ \text{false} & \text{otherwise,} \end{cases}$$

and $\text{Compl}(t)$ denotes the set of all tuples u (not containing null values) such that $t \leq u$.

By Lemma 4.2, it is now sufficient to prove that

$$\text{Rep}(f(\mathbf{T})) \approx f(\text{Rep}(\mathbf{T})) \tag{4.3}$$

for all f and \mathbf{T} . Let us first notice that

$$\text{if } \mathbf{T} = \langle T_1, \dots, T_n \rangle, \text{ then } \text{Rep}(\text{pr}_i(\mathbf{T})) = \text{pr}_i(\text{Rep}(\mathbf{T})), \tag{4.4}$$

and that for any Codd table T

$$\text{Rep}(\pi_Y(T)) = \pi_Y(\text{Rep}(T)), \tag{4.5}$$

$$\text{Rep}(\sigma_E(T)) \approx \sigma_E(\text{Rep}(T)). \tag{4.6}$$

We show (4.6) ((4.4) and (4.5) are obvious).

Let $r \in \text{Rep}(\sigma_E(T))$. Then r contains a relation $s \in \sigma_E(\text{Rep}(T))$ obtained by replacing nulls in the @-tuples $t \in T$ that belong to $\sigma_E(T)$ in such a way that the resulting tuples are in r , and by replacing nulls in every $t \in T \setminus \sigma_E(T)$ so as to obtain a tuple t^* with $E(t^*) = \text{false}$ (this is possible since for every $t \in T \setminus \sigma_E(T)$, $E_*(t) = \text{false}$).

Conversely, let $s \in \sigma_E(\text{Rep}(T))$. Then s contains a relation $r \in \text{Rep}(\sigma_E(T))$ obtained from T by replacing those @-tuples in T that belong to $\sigma_E(T)$ so as to obtain tuples in s , and by omitting the remaining @-tuples.

The desired formula (4.3) now easily follows by induction on the number of operators in f . We show, as an example, the inductive step in the case in which $f = \sigma_{Eg}$, under the inductive assumption $\text{Rep}(g(\mathbf{T})) \approx g(\text{Rep}(\mathbf{T}))$:

$$\begin{aligned} \text{Rep}(f(\mathbf{T})) &= \text{Rep}(\sigma_E(g(\mathbf{T}))) \approx \sigma_E(\text{Rep}(g(\mathbf{T}))) \\ &\approx \sigma_E(g(\text{Rep}(\mathbf{T}))) = f(\text{Rep}(\mathbf{T})) \end{aligned}$$

(we made use of (4.6) and then of (4.2)). \square

Note that in Theorem 4.5 selection can be arbitrary, on the basis of an arbitrary function from tuples to $\{\text{true}, \text{false}\}$. This includes as a special case selection based on arbitrary Boolean combinations generated by atomic conditions of the form $(A \leq a)$, $(A \leq B)$ ($A, B \in \mathcal{Z}$, $a \in D(A)$; we assume that the attribute domains are linearly ordered). These cases are treated in more detail in the next section.

The next two theorems say that, roughly speaking, our representation system based on Codd tables cannot, in addition, handle either union or join.

THEOREM 4.6. *It is not possible to correctly evaluate PSU-expressions over Codd tables; more formally $\langle \mathcal{T}_\otimes, \text{Rep}, \text{PSU} \rangle$ is not a representation system.*

PROOF. We now give an example of a PSU-expression f and a Codd table T such that $f(\text{Rep}(T))$ is not PSU-representable by any Codd multitable. Let

$$f(R) = \langle \sigma_{A=a}(R), \sigma_{A \neq a}(R) \rangle \quad (\alpha(f) = AB, a \in D(A)),$$

and let $T = \overline{\{a\}}$. Suppose that $U = \langle U_1, U_2 \rangle$ PSU-represents $f(\text{Rep}(T))$. Then obviously U P-represents $f(\text{Rep}(T))$. But the multitable $\langle \emptyset, \emptyset \rangle$ is easily seen to P-represent $f(\text{Rep}(T))$ (see the construction in the proof of Theorem 4.1), so that by Theorems 4.2 and 4.3, $U = \langle \emptyset, \emptyset \rangle$.

Consider the relational expression $g(S_1, S_2) = \pi_B(S_1 \cup S_2)$. We have

$$f(\text{Rep}(T))^g = \cap gf(\text{Rep}(T)) = \cap \pi_B(\text{Rep}(T)) = \{b\} \neq \emptyset = \text{Rep}(U)^g,$$

that is, U does not PSU-represent $f(\text{Rep}(T))$, contrary to our assumption. \square

It may be noted that if we restrict ourselves to the case in which all relational PSU-expressions considered involve only one relation symbol, then $f(\text{Rep}(T))$ can always be PSU-represented by a Codd table U (see [13]). This is a simple consequence of the fact that any PSU-expression f involving only one relation symbol can be transformed into an equivalent PS-expression of the form $\pi_Y(\sigma_E(R))$. The easy proof of this fact, by induction on the number of operators in f , is omitted.

THEOREM 4.7. *It is not possible to correctly evaluate PJ-expressions over Codd tables; more formally, $\langle \mathcal{T}_\otimes, \text{Rep}, \text{PJ} \rangle$ is not a representation system.*

PROOF. We give an example of a Codd table T and a PJ-expression f such that $f(\text{Rep}(T))$ is not PJ-representable by any Codd table. Let

$$T = \{a@c, a'@c'\} \quad (a \neq a', c \neq c') \\ f(R) = \pi_{AC}(R) \bowtie \pi_B(R).$$

It is easy to see that $f(\text{Rep}(T)) =_P \text{Rep}(T)$. By Theorem 4.2, if there is a Codd table U PJ-representing $f(\text{Rep}(T))$, then $T =_P U$, and consequently $f(\text{Rep}(T)) =_{PJ} \text{Rep}(T)$.

This last equivalence is, however, not true. Indeed, let

$$g(R) = \pi_{AC}(\pi_{AB}(R) \bowtie \pi_{BC}(R)).$$

We have

$$\text{Rep}(T)^g = \begin{bmatrix} a & c \\ a' & c' \end{bmatrix},$$

while

$$f(\text{Rep}(T))^g = \begin{bmatrix} a & c \\ a' & c' \\ a & c' \\ a' & c \end{bmatrix}. \quad \square$$

A representation system with $\Omega \supseteq \text{PJ}$ is considered in Section 6.

Now we compare our approach to that of Biskup [3]. Roughly speaking, Biskup defines the union of two Codd tables to be the usual set-theoretical union, so that

$$\text{Rep}(T \cup U) = \{r \cup s : r \in \text{Rep}(T) \wedge s \in \text{Rep}(U)\},$$

and he defines the join of two Codd tables in such a way that

$$\text{Rep}(T \bowtie U) \equiv_P \{r \bowtie s : r \in \text{Rep}(T) \wedge s \in \text{Rep}(U)\}. \quad (4.7)$$

His definition of join is based on the following informal matching rule:

$$@ \neq @, \quad @ \neq a.$$

(In fact, it may be shown that the sets on both sides of (4.7) are cointial so that they are PSUJ-equivalent.)

These definitions of union and join are semantically correct only if we assume that both arguments are independent; that is, every substitution of values in appropriate attribute domains for the null values in T and U is meaningful. This is clearly not the case when we want to correctly evaluate relational expressions such as

$$\begin{aligned} f(R) &= \sigma_{A=a}(R) \cup \sigma_{A \neq a}(R), \\ g(R) &= \pi_{AB}(R) \bowtie \pi_{AC}(R), \end{aligned}$$

over a Codd table. Consequently, Biskup's approach does not generalize to arbitrary relational expressions, though in a sense, it does give correct results for single relational operators. On the other hand, Biskup's approach is more general in that it allows "universal" null values in addition to the usual "existential" null values @.

5. Evaluating a Selection over a Codd Table

Recall that the selection of a Codd table is defined by

$$\sigma_E(T) = \{t \in T : E_*(t) = \text{true}\},$$

where $E_*(t) = \text{true}$ iff $E(t') = \text{true}$ for every $t' \in \text{Compl}(t)$. Of course, we should have a more efficient method of computing $E_*(t)$ than that involving the computation of $E(t')$ for every $t' \in \text{Compl}(t)$. ($\text{Compl}(t)$ may be infinite if attribute domains are infinite.)

In this section we always assume that the selection condition E is of the form of an arbitrary Boolean combination (using \neg , \vee , \wedge) of certain atomic conditions. The case in which these conditions are of the form $(A \text{ IN } F)$, $F \subseteq D(A)$ and in which subset entries of the form G , $G \subseteq D(A)$, meaning "an unknown value in G ," are allowed in the table, instead of just @, "an unknown value in $D(A)$," was extensively studied in [19], [20], and [21].

Here we sketch a method of computing $E_*(t)$ in the case in which the atomic conditions are of the form $(A = a)$, $(A = B)$, $(A \leq a)$, $(A \leq B)$ ($A, B \in \mathcal{U}$, $a \in D(A)$); in the third and fourth cases we assume a linear order on the corresponding attribute domains). The conditions $(A = a)$, $(A \leq a)$ are called *unary*, whereas $(A = B)$, $(A \leq B)$ are called *binary*. For simplicity, throughout this section we assume that every attribute domain is the same set D .

Note that we may "precompute" $E_*(t)$ by substituting in E the value $t(A)$ for every occurrence of A , for every A such that $t(A) \neq @$. This has the effect that some of the binary conditions become unary, and some conditions are reduced to $(a = b)$ or $(a \leq b)$, which can be replaced by **true** or **false** and then eliminated by using the absorption-type Boolean axioms. Let E' be the resulting condition. Clearly, the whole process can be carried out in time linear in the length of E . In this way our problem is reduced to evaluating $E'_*(t')$ where t' is a tuple of null values. This is obviously equivalent to testing whether E' is a tautology.

Let us begin with the case involving only conditions of the form $(A = a)$, $(A = B)$. We transform E' into a conjunctive normal form $\bigwedge_j W_j E_j^y$, where E_j^y is an atomic condition, and E_j^y denotes E_j or $\neg E_j$, depending on whether ϵ_{ij} is 1 or 0, respectively; we may also assume that no disjunct of this normal form contains both an equality and the negation to this equality. E' is a tautology iff for every i , $\bigvee_j W_j E_j^y$ is a tautology. A disjunction of the form $\bigvee_j W_j E_j^y$ is a tautology iff $\bigwedge_j E_j^{1-y}$ is not satisfiable. The satisfiability of this last conjunction can be tested in the following way. We construct a (nondirected) graph with vertices corresponding to attributes and constants occurring in $\bigwedge_j E_j^{1-y}$ (different occurrences of the same constant or attribute correspond to the same vertex), two vertices x, y joined by an edge iff the (nonnegated) condition $(x = y)$ appears in our conjunction.

It is clear that the conjunction is satisfiable iff every connected component of our graph contains at most one vertex corresponding to a constant. Obviously, finding the connected components can be done in time linear with respect to the number of vertices and edges of our graph, that is, $O(n)$, where n is the length of our conjunction (see, e.g., [11]). (Identifying *different* occurrences of attributes and constants may require $\Omega(n \log n)$ time.) The most complex part of this algorithm is, however, the transformation into a conjunctive normal form, since in the worst case it may involve an exponential growth of the length of our expression. We should not expect a substantially more efficient algorithm to evaluate $E_*(t)$, since it obviously contains as a special case the problem of deciding whether a Boolean expression is a tautology, a problem that is known to be NP-hard (see, e.g., [9]).

We now consider the case in which all four types of atomic conditions are allowed. Since

$$(x = y) \Leftrightarrow (x \leq y) \wedge (y \leq x), \tag{5.1}$$

our problem is reduced, in a similar way as before, to testing for the satisfiability of a conjunction $\bigwedge_j E_j$, where each E_j is either of the form $(x \leq y)$ or $(x < y)$. (It may be noted, however, that using (5.1) to eliminate equality may be not the most efficient approach.) Let us construct a directed graph with vertices corresponding to different attributes occurring in $\bigwedge_j E_j$, with an $<$ -edge or \leq -edge $\langle x, y \rangle$ from x to y iff $(x < y)$ or $(x \leq y)$, respectively, occurs in our conjunction. We associate with every vertex x an interval $I(x)$ with the beginning

$$b(x) = \max\{y \in D : (y < x) \text{ or } (y \leq x) \text{ appears in } \bigwedge_j E_j\}$$

and end

$$e(x) = \min\{y \in D : (x < y) \text{ or } (x \leq y) \text{ appears in } \bigwedge_j E_j\}$$

($b(x)$ is included into $I(x)$ iff $(b(x) < x)$ appears in the conjunction; similarly for $e(x)$).

We find strongly connected components of the graph consisting of \leq -edges (see, e.g., [25]; x and y are in the same strongly connected component iff there is a path from x to y and a path from y to x).

For every strongly connected component, we identify all vertices of the component and replace them with a single vertex x_C with $I(x_C) = \bigcap_{y \in C} I(y)$, where C is the set of vertices of the component. (In the resulting graph, there is an edge from $x_{C'}$ to $x_{C''}$, $C' \neq C''$, iff there was an edge from some $x' \in C'$ to some $x'' \in C''$ in the original graph.) The graph obtained is always acyclic and we test for the satisfiability of our conjunction in the following way. We try to assign a value $v(x) \in I(x)$ to every vertex x in such a way that $v(x) < v(y)$ ($v(x) \leq v(y)$) if there is a $<$ -edge (\leq -edge, respectively) $\langle x, y \rangle$. Let us assume, for simplicity, that our graph

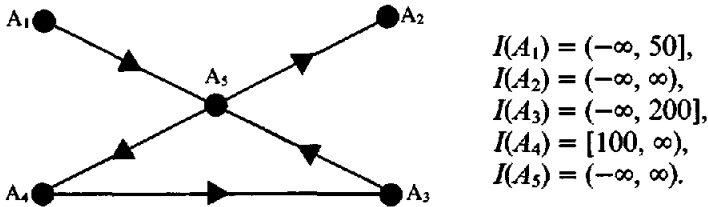
contains only \leq -edges, and that every $I(x)$ contains both $b(x)$ and $e(x)$ (the general case can be treated in a similar way). We set initially $v(x) = b(x)$ for any x , and we perform a breadth-first search (see, e.g., [11]) starting from the set of sources (i.e., vertices with no incoming edges). For any edge $\langle x, y \rangle$ processed by the search we set $v(y) := \max(v(y), v(x))$. If at any stage of the process this results in $v(y) > e(y)$, then clearly our conjunction is not satisfiable. Otherwise, the process determines an assignment of values to attributes that makes our conjunction true. We leave the details to the reader.

Example 5.1. Let our conjunction be

$$(A_1 \leq 50) \wedge (A_4 \geq 100) \wedge (A_3 \leq 200) \wedge (A_1 \leq A_5) \\ \wedge (A_5 \leq A_2) \wedge (A_5 \leq A_4) \wedge (A_4 \leq A_3) \wedge (A_3 \leq A_5)$$

(the domain of each attribute is the set of integers).

Our graph has the form



The strongly connected components have vertex sets $\{A_1\}$, $\{A_2\}$, $\{A_3, A_4, A_5\}$, and our graph is transformed into



The breadth-first search produces $v(A_1) = -\infty$, $v(A_{345}) = 100$, $v(A_2) = 100$, which means that our conjunction is satisfiable. \square

6. *V-Tables*

One of the reasons for the inability of Codd tables to correctly support the join operation appears to be the fact that we are not able to represent the information that two different occurrences of \emptyset represent the same value. In this section we consider tables where, for any $A \in \mathcal{Z}$, we have an infinite set of possible “null values.” The interpretation of such a table is that the values represented by two occurrences of the same null value, though unknown, are the same. Such tables will turn out to be suitable for a representation system with $\Omega = \text{PS}^* \text{UJ}$.

Formally, for every $A \in \mathcal{Z}$, let $V(A)$ be a countably infinite set of symbols called *variables*. We assume that $V(A) \cap D = \emptyset$, $V(A) \cap V(B) = \emptyset$ if $D(A) \neq D(B)$, and $V(A) = V(B)$ if $D(A) = D(B)$. By a *V-tuple* on X , we mean any mapping t that associates an element $t(A) \in D(A) \cup V(A)$ with every $A \in X$. A *V-table* on X is any finite set of *V-tuples* on X , and a *V-multitable* of type $\langle X_1, \dots, X_n \rangle$ is any sequence $\mathbf{T} = \langle T_1, \dots, T_n \rangle$ where T_i is a *V-table* on X_i , $1 \leq i \leq n$ (notice that the same variable may occur in several T_i 's). The set of all *V-multitables* is denoted by \mathcal{T}_V .

Let $V = \bigcup_{A \in \mathcal{Z}} V(A)$. By a *valuation*, we mean any mapping $v: V \rightarrow D$ such that $x \in V(A)$ implies $v(x) \in D(A)$. Any valuation can be extended to the set of constants by putting $v(c) = c$ for every $c \in D$; to *V-tuples*, by defining, for any *V-tuple* t on

X , $v(t)$ to be the tuple on X satisfying

$$(v(t))(A) = v(t(A)) \quad \text{for every } A \in X; \tag{6.1}$$

to V -tables,

$$v(T) = \{v(t) : t \in T\}; \tag{6.2}$$

and finally to V -multitables $\mathbf{T} = \langle T_1, \dots, T_n \rangle$,

$$v(\mathbf{T}) = \langle v(T_1), \dots, v(T_n) \rangle. \tag{6.3}$$

The intuition given at the beginning of this section can now be formalized by defining, for any V -multitable, $\mathbf{T} = \langle T_1, \dots, T_n \rangle$ of type $\langle X_1, \dots, X_n \rangle$:

$$\text{Rep}(\mathbf{T}) = \{r \in \mathcal{R}(X_1, \dots, X_n) : \text{there exists a valuation } v \text{ such that } v(\mathbf{T}) \subseteq r\}. \tag{6.4}$$

Note that in view of this definition, any Codd multitable can informally be identified with a V -multitable obtained by replacing any occurrence of @ by a different variable.

Let us associate with a V -multitable $\mathbf{T} = \langle T_1, \dots, T_n \rangle$ a first-order formula $\phi(\mathbf{T})$ (of a many-sorted predicate calculus) defined as

$$\exists x_1 \dots \exists x_m \bigwedge_t \phi(t)$$

where t ranges over all V -tuples appearing in \mathbf{T} and, for $t \in T_i$, $\phi(t)$ is $R_i(t)$, where R_i is a predicate symbol of a suitable type and x_1, \dots, x_m are all variables occurring in \mathbf{T} . It is easy to see that $\text{Rep}(\mathbf{T})$ is simply the set of all finite models of $\phi(\mathbf{T})$ (over a fixed universe given by $D(A)$, $A \in \mathcal{U}$).

In this section we always assume that all attribute domains are infinite.

The following is the main result of this paper.

THEOREM 6.1. *It is possible to correctly evaluate PS^+UJ -expressions over V -tables; more formally $\langle \mathcal{V}, \text{Rep}, PS^+UJ \rangle$ is a representation system.*

This theorem is proved in the next section, in which we develop a suitable formal notion, that of a *conditional table* (see also another proof in [23]). Here we only give the relevant definitions of the relational operators over V -tables.

For any V -multitable $\mathbf{T} = \langle T_1, \dots, T_n \rangle$ we put

$$\text{pr}_i(\mathbf{T}) = T_i, \quad 1 \leq i \leq n. \tag{6.5}$$

The definition of projection is the “natural one,”

$$\pi_Y(T) = \{t[Y] : t \in T\}, \tag{6.6}$$

and so is the definition of join of two V -tables T, W on X, Z , respectively:

$$T \bowtie W = \{t : \alpha(t) = X \cup Z \wedge t[X] \in T \wedge t[Z] \in W\}. \tag{6.7}$$

The union operator acts on V -tables as the usual set-theoretical union.

Finally,

$$\sigma_E(T) = \{t \in T : E_*(t) = \text{true}\} \tag{6.8}$$

where

$$E_*(t) = \begin{cases} \text{true} & \text{if } E(v(t)) = \text{true for every valuation } v, \\ \text{false} & \text{otherwise.} \end{cases}$$

Note that, if E is positive, then $E_*(t)$ can be computed in a very simple way by using the following rule: We evaluate every atomic condition $(A = B)$ in E to **true** if $t(A) = t(B)$ and to **false** otherwise; we evaluate every atomic condition $(A = a)$ to **true** if $t(A) = a$ and to **false** otherwise; and then we use \vee, \wedge in the natural way. In other words,

$$\begin{aligned} (A = B)_*(t) &= \begin{cases} \text{true} & \text{if } t(A) = t(B), \\ \text{false} & \text{otherwise,} \end{cases} \\ (A = a)_*(t) &= \begin{cases} \text{true} & \text{if } t(A) = a, \\ \text{false} & \text{otherwise,} \end{cases} \\ (E \vee E')_*(t) &= E_*(t) \vee E'_*(t), \\ (E \wedge E')_*(t) &= E_*(t) \wedge E'_*(t). \end{aligned}$$

The correctness of this rule follows easily from the fact that, if there is a valuation ν such that $E(\nu(t)) = \text{false}$, then $E(\nu'(t)) = \text{false}$ for any valuation ν' such that $\nu'(x) \neq \nu'(y)$ for $x \neq y$ and $\nu'(x)$ does not appear in E for any $x \in V$. Note that this rule does *not* evaluate $E_*(t)$ correctly if E contains negation, or if attribute domains are finite.

The above rules, together with the obvious rule

$$f(\mathbf{T}) = \langle f_1(\mathbf{T}), \dots, f_k(\mathbf{T}) \rangle \tag{6.9}$$

for any $\mathbf{f} = \langle f_1, \dots, f_k \rangle$, define inductively $f(\mathbf{T})$ for any multirelational PS+UJ-expression \mathbf{f} and any V -multitable \mathbf{T} ($\alpha(\mathbf{T}) = \alpha(\mathbf{f})$).

To sum up, we evaluate PS+UJ-expressions over V -tables in *exactly the same way* as if the variables were values in the attribute domains. Let us emphasize that the fact that this simple method of evaluating PS+UJ-expressions gives correct results is not quite trivial (see the proof in the next section).

Example 6.1. Let us evaluate

$$f(T) = \pi_{AC}(\pi_{AB}(T) \bowtie \sigma_{(B=C) \vee (C=c)}(\pi_{BC}(T))),$$

where T is the following V -table:

A	B	C
x	y	c
a	b	c
a'	b'	c'
a	y	z
x	d	d

We get

$$\begin{aligned} \pi_{AB}(T) &= \{xy, ab, a'b', xy, xd\}, \\ \pi_{BC}(T) &= \{yc, bc, b'c', yz, dd\}, \\ \sigma_{(B=C) \vee (C=c)}(\pi_{BC}(T)) &= \{yc, bc, dd\}, \\ \pi_{AB}(T) \bowtie \sigma_{(B=C) \vee (C=c)}(\pi_{BC}(T)) &= \{xyc, abc, xdd\}, \\ f(T) &= \{xc, ac, xd\}. \end{aligned}$$

□

Rather suprisingly, it turns out that V -multitables, which are “more powerful” than Codd multitables, cannot support a representation system with $\Omega = \text{PS}$.

THEOREM 6.2. *It is not possible to correctly evaluate PS-expressions over V -tables; more formally, $\langle \mathcal{T}_V, \text{Rep}, \text{PS} \rangle$ is not a representation system.*

PROOF. Let T be the following V -table:

A	B	C
a	y	c
a'	y	c

and let f be the following PS-expression—in fact, an S-expression:

$$f(R) = \sigma_{(A=a) \wedge (B=b) \vee (A=a') \wedge (B \neq b)}(R) \quad (a \neq a').$$

We claim that there is no V -table U PS-representing $f(\text{Rep}(T))$. Indeed, let

$$g(R) = \pi_C(\sigma_{(A=a) \vee (A=a')}(R)).$$

Notice that $f(\text{Rep}(T))$ is coinital with

$$\mathcal{X} = \{\overline{abc}\} \cup \{\overline{a'dc} : d \in D(B) \setminus \{b\}\},$$

so that $f(\text{Rep}(T))^s = \mathcal{X}^s = \{c\}$.

If there is a V -tuple $pqr \in U$, where $p \in D(A)$, then $p \in \text{Rep}(U)^{*A} \neq \mathcal{X}^{*A} = \emptyset$. If for every $pqr \in U$, p is a variable (in particular, if $U = \emptyset$), then clearly $\text{Rep}(U)^s = \emptyset$, although we have already noted that $f(\text{Rep}(T))^s = \{c\}$. Hence, in both cases, U does not PS-represent $f(\text{Rep}(T))$. \square

The apparent contradiction in the inability of V -multitables to support PS-expressions can be intuitively explained in the following way: $\langle \mathcal{T}, \text{Rep}, \Omega \rangle$ is a representation system if the class \mathcal{T} of multitables is a “fixpoint” with respect to Ω -expressions, in the sense that if we take any Ω -expression f and any $T \in \mathcal{T}$ then to Ω -represent $f(\text{Rep}(T))$ we do not need any more “representation power” than that available in \mathcal{T} . If we consider an arbitrary $\mathcal{T}' \supseteq \mathcal{T}$, there is no reason for \mathcal{T}' to be also such a fixpoint.

7. Conditional Tables

By a conditional table we mean a V -table extended by one additional special column, called *con*, which contains for any tuple a *condition* from a set \mathcal{L} .

More precisely, \mathcal{L} is the set of all expressions built up from atomic conditions of the form $(x = a)$, $(x = y)$, **false** and **true**, where for some $A \in \mathcal{U}$, $a \in D(A)$, $x, y \in V(A)$, by means of \neg , \vee , and \wedge . A condition is *positive*, if it does not contain \neg .

By a C -tuple on X we mean any mapping t defined on $X \cup \{\text{con}\}$ such that $t[X]$ is a V -tuple and $t(\text{con}) \in \mathcal{L}$. A *conditional table* (or *C-table*) on X is any finite set T of C -tuples on X . (Note: We assume that *con* is not part of $\alpha(t)$ or $\alpha(T)$.) A *conditional multitable* (or *C-multitable*) of type $\langle X_1, \dots, X_n \rangle$ is any sequence $\mathbf{T} = \langle T_1, \dots, T_n \rangle$ where T_i is a C -table on X_i , $1 \leq i \leq n$. The set of all C -multitables is denoted by \mathcal{T}_C .

Any valuation can be extended in the natural way to conditions.

For any valuation v , any C -table T on X , and any C -multitable $\mathbf{T} = \langle T_1, \dots, T_n \rangle$ of type $\langle X_1, \dots, X_n \rangle$, we define

$$v(T) = \{v(t[X]) : v(t(\text{con})) = \text{true}\}, \tag{7.1}$$

$$v(\mathbf{T}) = \langle v(T_1), \dots, v(T_n) \rangle, \tag{7.2}$$

and finally

$$\text{Rep}(\mathbf{T}) = \{r \in \mathcal{R}(X_1, \dots, X_n) : \text{there exists a valuation } v \text{ such that } v(\mathbf{T}) \subseteq r\}. \tag{7.3}$$

We say that conditions $\gamma, \delta \in \mathcal{L}$ are *equivalent*, and we write $\gamma \approx \delta$, if $v(\gamma) = v(\delta)$ for any valuation v . It can easily be shown that $\gamma \approx \delta$ iff γ can be transformed into δ by using the axioms of Boolean algebra and the axioms of equality. It should be clear that if we (a) replace each of the conditions in a C -table T by an equivalent one, (b) delete all $t \in T$ with $t(\text{con}) \approx \text{false}$, and (c) replace some $t_1, \dots, t_k \in T$ such that $t_1[X] = \dots = t_k[X]$ ($X = \alpha(T)$) by a single C -tuple t such that $t[X] = t_1[X]$ and $t(\text{con}) = \bigvee_{i=1}^k t_i(\text{con})$, then the resulting C -table U will be Rep-equivalent to T (i.e., $\text{Rep}(T) = \text{Rep}(U)$, see Section 3). We freely make use of the equivalent transformations of this kind; in particular, we use rule (c) to *normalize* a C -table T on X , that is, to replace T by a Rep-equivalent C -table T^0 on X not containing different C -tuples agreeing on X . In what follows we assume all C -tables to be normalized.

THEOREM 7.1. *It is possible to correctly evaluate PSUJ-expressions over C-tables; more formally, $\langle \mathcal{T}_C, \text{Rep}, \text{PSUJ} \rangle$ is a representation system.*

PROOF. We give definitions of the relational operators over C -tables that inductively define $f(T)$ for any C -multitable T and PSUJ-expression $f(\alpha(T) = \alpha(f))$, in such a way that

$$f(\text{Rep}(T)) =_{\text{PSUJ}} \text{Rep}(f(T)). \tag{7.4}$$

For any C -multitable $T = \langle T_1, \dots, T_n \rangle$ we put

$$\text{pr}_i(T) = T_i, \quad 1 \leq i \leq n, \tag{7.5}$$

and for any C -tables T, W on X, Z , respectively, we define

$$\pi_Y(T) = \{t[Y \cup \{\text{con}\}] : t \in T\}^0, \tag{7.6}$$

$$\sigma_E(T) = \{\sigma_E(t) : t \in T\}, \tag{7.7}$$

where $\sigma_E(t)$ is the C -tuple on X with

$$\begin{aligned} \sigma_E(t)[X] &= t[X], \\ \sigma_E(t)(\text{con}) &= t(\text{con}) \wedge E(t). \end{aligned}$$

($E(t)$ is the result of substituting $t(A)$ for A in E , for every $A \in X$), and

$$T \bowtie W = \{t \bowtie w : t \in T \wedge w \in W\}^0, \tag{7.8}$$

where $t \bowtie w$ is the C -tuple on $X \cup Z$ with

$$\begin{aligned} (t \bowtie w)(A) &= \begin{cases} t(A) & \text{if } A \in X, \\ w(A) & \text{if } A \in Z \setminus X, \end{cases} \\ (t \bowtie w)(\text{con}) &= t(\text{con}) \wedge w(\text{con}) \wedge \bigwedge_{A \in X \cap Z} (t(A) = w(A)). \end{aligned}$$

If $X = Z$, then

$$T \cup W = (T \cup W)^0, \tag{7.9}$$

where \cup on the left-hand side denotes the relational union operator that we define on C -tables, and \cup on the right-hand side is the usual set-theoretical union.

Finally, if $f = \langle f_1, \dots, f_k \rangle$, then

$$f(T) = \langle f_1(T), \dots, f_k(T) \rangle. \tag{7.10}$$

We now prove that under the definition of $f(T)$ given by (7.5)–(7.10),

$$v(f(T)) = f(v(T)) \tag{7.11}$$

for any valuation v . (On the right-hand side, $f(v(T))$ is understood as the result of performing f over the multirelation $v(T)$ in the usual way.) Clearly, (7.11) implies $\text{Rep}(f(T)) \approx f(\text{Rep}(T))$, which by Lemma 4.2, proves (7.4).

We show (7.11) by verifying that

$$\begin{aligned} v(\text{pr}_i(\mathbf{T})) &= \text{pr}_i(v(\mathbf{T})), \\ v(\pi_Y(T)) &= \pi_Y(v(T)), \\ v(T \bowtie W) &= v(T) \bowtie v(W), \\ v(T \cup W) &= v(T) \cup v(W) \quad (\alpha(T) = \alpha(W)), \end{aligned}$$

for any C -multitable \mathbf{T} and any C -tables T, W .

Let $\alpha(T) = X, \alpha(W) = Z$. In the case of projection, $t \in v(\pi_Y(T))$ iff there is a C -tuple $t' \in \pi_Y(T)$ such that $v(t'[Y]) = t, v(t'(\text{con})) = \text{true}$; that is, iff there exists a C -tuple $t'' \in T$ such that $\pi_Y(v(t''[X])) = t, v(t''(\text{con})) = \text{true}$, which is exactly the condition for t to be in $\pi_Y(v(T))$.

In the case of selection, $t \in v(\sigma_E(T))$ iff there is a $t' \in T$ such that $v(t'(\text{con})) = \text{true}, v(E(t')) = \text{true}$, and $t = v(t'[X])$, that is, iff $t \in \sigma_E(v(T))$.

In the case of join, a tuple t of type $X \cup Z$ is in $v(T \bowtie W)$ iff there are tuples $t' \in T, w' \in W$ such that $v((t' \bowtie w')(\text{con})) = \text{true}, v((t' \bowtie w')([X \cup Z])) = t$, and this is clearly equivalent to $t \in v(T) \bowtie v(W)$.

The case of union is obvious. \square

Example 7.1. Let us evaluate

$$f(R) = \sigma_{C=c}(\pi_{AC}(\pi_{AB}(R) \bowtie \pi_{BC}(R)))$$

over the following C -table:

$$T = \begin{array}{|c|c|c|c|} \hline A & B & C & \text{con} \\ \hline a & b & z & z \neq c \\ \hline a & y & c & y \neq b \\ \hline x & b & c & x \neq a \\ \hline \end{array}$$

We have

$$U = \pi_{AB}(T) \bowtie \pi_{BC}(T) = \begin{array}{|c|c|c|c|} \hline A & B & C & \text{con} \\ \hline a & b & z & (z \neq c) \wedge (z \neq c) \wedge (b = b) \\ \hline a & b & c & (z \neq c) \wedge (y \neq b) \wedge (b = y) \\ \hline a & b & c & (z \neq c) \wedge (x \neq a) \wedge (b = b) \\ \hline a & y & z & (y \neq b) \wedge (z \neq c) \wedge (y = b) \\ \hline a & y & c & (y \neq b) \wedge (y \neq b) \wedge (y = y) \\ \hline a & y & c & (y \neq b) \wedge (x \neq a) \wedge (y = b) \\ \hline x & b & z & (x \neq a) \wedge (z \neq c) \wedge (b = b) \\ \hline x & b & c & (x \neq a) \wedge (y \neq b) \wedge (b = y) \\ \hline x & b & c & (x \neq a) \wedge (x \neq a) \wedge (b = b) \\ \hline \end{array}$$

$$= \begin{array}{|c|c|c|c|} \hline A & B & C & \text{con} \\ \hline a & b & z & (z \neq c) \\ \hline a & b & c & (z \neq c) \wedge (y \neq b) \\ \hline a & y & c & (y \neq b) \\ \hline x & b & z & (x \neq a) \wedge (z \neq c) \\ \hline x & b & c & (x \neq a) \\ \hline \end{array}$$

$$W = \pi_{AC}(U) = \begin{array}{|c|c|c|} \hline A & C & \text{con} \\ \hline a & z & (z \neq c) \\ a & c & (y \neq b) \\ x & z & (x \neq a) \wedge (z \neq c) \\ x & c & (x \neq a) \\ \hline \end{array}$$

and finally

$$\sigma_{C=c}(W) = \begin{array}{|c|c|c|} \hline A & C & \text{con} \\ \hline a & z & (z \neq c) \wedge (z = c) \\ a & c & (y \neq b) \wedge (c = c) \\ x & z & (x \neq a) \wedge (z \neq c) \wedge (z = c) \\ x & c & (x \neq a) \wedge (c = c) \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline A & C & \text{con} \\ \hline a & c & (y \neq b) \\ x & c & (x \neq a) \\ \hline \end{array} \quad \square$$

Let T be a C -table on X . As usual we assume that no two different tuples in T agree on X . We define the *unconditional part* of T to be

$$T_* = \{t \in T : t(\text{con}) \approx \text{true}\}.$$

The unconditional part of a C -multitable $\mathbf{T} = \langle T_1, \dots, T_n \rangle$ is defined as $\mathbf{T}_* = \langle T_{1*}, \dots, T_{n*} \rangle$.

A C -multitable $\mathbf{T} = \langle T_1, \dots, T_n \rangle$ is called *positive* if for every T_i and every $t \in T_i$, $t(\text{con})$ is positive.

LEMMA 7.1. *Let all attribute domains be infinite, and let \mathbf{T} be a positive C -multitable. Then*

$$\bigcap \text{Rep}(\mathbf{T}) = \bigcap \text{Rep}(\mathbf{T}_*).$$

PROOF. Obviously, we may restrict ourselves to the case where \mathbf{T} is a single C -table, T .

Clearly $\text{Rep}(T) \subseteq \text{Rep}(T_*)$, and consequently $\bigcap \text{Rep}(T_*) \subseteq \bigcap \text{Rep}(T)$, so that it is sufficient to prove the converse inclusion.

Suppose $t \notin \bigcap \text{Rep}(T_*)$. It means that there is a valuation v such that $t \notin v(T_*)$. Let v' be a valuation such that $v'(x) = v'(y)$ for $x \neq y$ and for every variable x appearing in T , $v'(x)$ does not appear in T or t (the existence of such a valuation v' is an obvious consequence of the assumption that the attribute domains are infinite). Clearly, v' makes all atomic conditions (which are positive, by our assumption) false, and consequently it makes false all conditions $t(\text{con})$, $t \in T \setminus T_*$. Since $t \notin v(T_*)$, it follows that T_* does not contain any tuple agreeing with t on $\alpha(t)$. The valuation v' associates constants different from those appearing in t to variables in T , so that it cannot produce t ; that is, $t \notin v'(T)$. Consequently, $t \notin \bigcap \text{Rep}(T)$. \square

We are now going to prove, as we promised in the previous section, that $(\mathcal{R}_V, \text{Rep}, \text{PS}^+\text{UJ})$ is a representation system.

A C -multitable \mathbf{T} is called *unconditional* if $\mathbf{T}_* = \mathbf{T}$. Let $\mathbf{T} = \langle T_1, \dots, T_n \rangle$ be an unconditional C -multitable and let U be the V -multitable obtained from \mathbf{T} by

deleting the con column from every T_i , $1 \leq i \leq n$. Then

$$\text{Rep}(T) = \text{Rep}(U),$$

where $\text{Rep}(T)$ is computed according to (7.1)–(7.3) and $\text{Rep}(U)$ according to (6.1)–(6.4). In what follows we always identify an unconditional C -multitable with the corresponding V -multitable. To avoid confusion, we write $f^C(T)$ and $f^V(T)$ for an expression f evaluated for an unconditional C -multitable T (identified with a V -multitable) according to the rules given for C -tables (see (7.5)–(7.10)) and V -tables (see (6.5)–(6.9)), respectively.

LEMMA 7.2. *Let all attribute domains be infinite, let T be a positive C -multitable and let f be a PS^+UJ -expression with $\alpha(f) = \alpha(T)$. Then*

$$f^V(T_*) = f^C(T)_*.$$

PROOF. It is sufficient to prove that for any positive C -multitable $T = \langle T_1, \dots, T_n \rangle$

$$\text{pr}_i^V(T_*) = \text{pr}_i^C(T)_*, \quad 1 \leq i \leq n \tag{7.12}$$

and that for any positive C -tables T, W

$$\pi_Y^V(T_*) = \pi_Y^C(T)_*, \tag{7.13}$$

$$\sigma_E^V(T_*) = \sigma_E^C(T)_* \quad \text{for } E \text{ positive}, \tag{7.14}$$

$$T_* \bowtie^V W_* = (T \bowtie^C W)_*, \tag{7.15}$$

$$T_* \cup^V W_* = (T \cup^C W)_*, \tag{7.16}$$

($X = \alpha(T)$, $Z = \alpha(W)$).

Before we verify these equalities, let us note the following simple fact concerning positive conditions:

For any positive conditions γ, δ

$$\gamma \vee \delta \approx \text{true} \Rightarrow \gamma \approx \text{true} \quad \text{or} \quad \delta \approx \text{true}. \tag{7.17}$$

Indeed, it is easy to see that, if $\gamma \not\approx \text{true}$ and $\delta \not\approx \text{true}$, then $v(\gamma \vee \delta) = \text{false}$, where v is any valuation such that $v(x) \neq v(y)$ for $x \neq y$ and $v(x)$ does not appear in $\gamma \vee \delta$ for any $x \in V$.

Since (7.12) is obvious, we now consider the case of projection. A V -tuple t is in $\pi_Y^V(T_*)$ iff there is a $t' \in T_*$ such that $t'[Y] = t[Y]$, that is, if there is a $t'' \in T$ such that $t''[Y] = t[Y]$ and $t''(\text{con}) \approx \text{true}$. On the other hand, $t \in \pi_Y^C(T)_*$ iff there are tuples $t_1, \dots, t_k \in T$, $k \geq 1$ such that $t_i[Y] = t[Y]$, $1 \leq i \leq k$ and $\bigwedge_{i=1}^k t_i(\text{con}) \approx \text{true}$. But by (7.17) the last condition implies that $t_i(\text{con}) \approx \text{true}$ for some i , and we obtain exactly the same condition as before.

In the case of selection, $\sigma_E^V(T_*)$ consists of all tuples $t \in T$ such that for every valuation v , $v(t(\text{con})) = \text{true}$ and $v(E(t)) = \text{true}$, which is exactly the condition for t to be in $\sigma_E^C(T)_*$. Notice that (7.14) holds true for any, not necessarily positive, condition E . The requirement that E be positive is needed only to guarantee that $\sigma_E^C(T)_*$ is positive, so that we could inductively prove the lemma using (7.12)–(7.16).

We now consider the case of join. It is obvious that $T_* \bowtie^V W_* \subseteq (T \bowtie^C W)_*$, so it is sufficient to prove the converse inclusion. Suppose $t \in (T \bowtie^C W)_*$. Then there are C -tuples $t_1, \dots, t_k \in T$, $w_1, \dots, w_k \in W$, $k \geq 1$, such that $t_i[X] = t[X]$,

$w_i[Z \setminus X] = t[Z \setminus X]$, $1 \leq i \leq k$, and $W_{i-1}^k(t_i \bowtie w_i)(\text{con}) \approx \text{true}$. By (7.17), for some i

$$t_i(\text{con}) \wedge w_i(\text{con}) \wedge \bigwedge_{A \in X \cap Z} (t_i(A) = w_i(A)) \approx \text{true}.$$

This obviously means that $t_i \in T_*$, $w_i \in W_*$, and $t_i[X \cap Z] = w_i[X \cap Z]$; that is, $t \in T_* \bowtie^V W_*$.

Finally, in the case of union, $T_* \cup^V W_* \subseteq (T \cup^C W)_*$, and $t \in (T \cup^C W)_*$ implies that either $t \in T_*$ or $t \in W_*$ or there are $t' \in T$, $w' \in W$ with $t'[X] = w'[X] = t[X]$, $t'(\text{con}) \vee w'(\text{con}) \approx \text{true}$. The last case is reduced, however, by (7.17), to either the first or the second one. \square

We are now ready to give the promised proof of the fact that $\langle \mathcal{F}_V, \text{Rep}, \text{PS}^+\text{UJ} \rangle$ is a representation system.

PROOF OF THEOREM 6.1. We have to show that for any V -multitable T and any PS^+UJ -expression f with $\alpha(f) = \alpha(T)$,

$$\text{Rep}(f^V(T)) \equiv_{\text{PS}^+\text{UJ}} f(\text{Rep}(T));$$

that is, for any PS^+UJ -expression g with appropriate $\alpha(g)$,

$$\text{Rep}(f^V(T))^g = f(\text{Rep}(T))^g. \tag{7.18}$$

The left-hand side of (7.18) can be transformed as follows:

$$\text{Rep}(f^V(T))^g = \bigcap_{r \in \text{Rep}(f^V(T))} g(r) = \bigcap g(\text{Rep}(f^V(T))).$$

By Theorem 7.1,

$$g(\text{Rep}(f^V(T))) \equiv_{\text{PSUJ}} \text{Rep}(g^C(f^V(T))),$$

which implies the equality of the intersections

$$\bigcap g(\text{Rep}(f^V(T))) = \bigcap \text{Rep}(g^C(f^V(T))).$$

By applying Lemma 7.1 and then Lemma 7.2 to the right-hand side we get

$$\begin{aligned} \bigcap g(\text{Rep}(f^V(T))) &= \bigcap \text{Rep}(g^C(f^V(T))_*) \\ &= \bigcap \text{Rep}(gf^V(T)). \end{aligned}$$

Again using Lemmas 7.2 and 7.1 and Theorem 7.1 we obtain

$$\begin{aligned} \bigcap \text{Rep}(gf^V(T)) &= \bigcap \text{Rep}(gf^C(T)_*) \\ &= \bigcap \text{Rep}(gf^C(T)) \\ &= \bigcap gf(\text{Rep}(T)) \end{aligned}$$

which is the same as the right-hand side of (7.18):

$$f(\text{Rep}(T))^g = \bigcap_{r \in f(\text{Rep}(T))} g(r) = \bigcap gf(\text{Rep}(T)). \quad \square$$

To conclude this section, we note that

$$f(v(T)) \supseteq v(f(T)) \tag{7.19}$$

for any PS^+UJ -expression f , any V -multitable T ($\alpha(T) = \alpha(f)$), and any valuation v . This is proved easily by induction on the complexity of f . An immediate corollary is

$$\text{Rep}(f(T)) \supseteq f(\text{Rep}(T)). \tag{7.20}$$

An important fact about (7.20) is that it does not require any assumptions on the cardinality of attribute domains. Notice that, in the terminology of the Introduction, (7.20) says that the representation system based on V -multitables is always safe: $\text{Rep}(f(\mathbf{T}))$ includes all “possible” $f(\mathbf{r}^*)$, $\mathbf{r}^* \in \text{Rep}(\mathbf{T})$. The assumption that the attribute domains are infinite is only needed to prove its completeness.

8. Renaming of Attributes, Conjunctive Queries

It may be noted that the relational algebra based on $\Omega = \text{PSUJD}$ and a fixed finite set \mathcal{Z} of attributes is not relationally complete in the sense of Codd [6]. This is because our version of relational algebra is equivalent in expressive power to a many-sorted predicate calculus with only one variable available for each sort, and certain formulas of the usual (many-sorted) predicate calculus inherently need many “auxiliary variables” (cf. [14]). One way of making our algebra relationally complete is to add the operation of renaming attributes, and to consider an infinite set of attributes \mathcal{Z} , which contains an infinite number of attributes for each sort (i.e., for every $A \in \mathcal{Z}$, there exists different A_1, A_2, \dots with $D(A_i) = D(A)$). It is obvious that

$$s_B^A(\text{Rep}(T)) = \text{Rep}(s_B^A(T)), \tag{8.1}$$

where $s_B^A(T)$ (with T either a Codd table, V -table, or C -table) is defined in exactly the same way as $s_B^A(r)$. By (8.1), most of our results can easily be extended by adding renaming (R) to the set Ω of relational operators. In particular, by Theorem 6.1, we have the following corollary.

COROLLARY 8.1. $\langle \mathcal{T}_V, \text{Rep}, \text{PS}^* \text{UJR} \rangle$ is a representation system.

We now discuss another version of the relational algebra (see, e.g., [29]) where the type of a relation is defined as a *sequence*, rather than set, of attributes (repetitions are allowed). The relational operators are modified in the following way. In the projection operator, instead of a target attribute set we have a target sequence of positions (i.e., numbers of columns of the argument relation); in the selection condition, we refer to positions rather than to attributes; and instead of the natural join we consider the Cartesian product. The Cartesian product of two relations, r, s , of types $\langle A_1, \dots, A_n \rangle$ and $\langle B_1, \dots, B_m \rangle$, respectively, is a relation of type $\langle A_1, \dots, A_n, B_1, \dots, B_m \rangle$ defined as

$$r \times s = \{ \langle a_1, \dots, a_n, b_1, \dots, b_m \rangle : \langle a_1, \dots, a_n \rangle \in r \wedge \langle b_1, \dots, b_m \rangle \in s \}. \tag{8.2}$$

Note that any relation, in the modified sense, of type $\langle A_1, \dots, A_n \rangle$ can be uniquely represented by the usual relation of type $\{ \langle A_1, 1 \rangle, \dots, \langle A_n, n \rangle \}$, and that (8.2) can be expressed by first renaming the attributes of s so that it is of type $\{ \langle B_1, n + 1 \rangle, \dots, \langle B_m, n + m \rangle \}$ and then performing the natural join. By Corollary 8.1, this easily implies the following theorem (where X stands for the Cartesian product):

THEOREM 8.1. $\langle \mathcal{T}_V, \text{Rep}, \text{PS}^* \text{UX} \rangle$ is a representation system.

Conversely, the join can easily be expressed by the Cartesian product, restriction, and projection, so that, by Theorem 4.7, the Codd tables cannot support $\Omega = \text{PXE}$ (E stands for restriction).

THEOREM 8.2. $\langle \mathcal{T}_\emptyset, \text{Rep}, \text{PXE} \rangle$ is not a representation system.

(Clearly, in Theorems 8.1 and 8.2 the notions of a V -table and Codd table are modified into a "column-ordered" form, in the same way as it was done for relations.)

Another important fact implied by Corollary 8.1 is that the system based on V -tables can handle arbitrary conjunctive queries [4]. Roughly speaking, a conjunctive query is defined by a conjunction of atomic relational formulas of the form $R_i(x_1, \dots, x_p)$ (the R_i 's are predicates corresponding to database relations), where each of the x_i is either a constant or variable, preceded by existential quantifiers binding some of the variables appearing in the conjunction. The mapping defined by a conjunctive query q can be defined, for any database state \mathbf{r} , by

$$q(\mathbf{r}) = \{v(s) : v \text{ is a valuation such that } v(Q) \subseteq \mathbf{r}\} \quad (8.3)$$

where

Q is a suitable V -multitable (V -tuples appearing in Q correspond in the natural way to the atomic relational formulas in q ; $\alpha(Q) = \alpha(\mathbf{r}) = \alpha(q)$),
 s is a V -tuple such that all variables occurring in s occur in Q ($\alpha(s) = \beta(q)$).

LEMMA 8.1. *For any conjunctive query q there is a PS^+JR -expression f such that $q(\mathbf{r}) = f(\mathbf{r})$ for every \mathbf{r} ($\alpha(\mathbf{r}) = \alpha(f) = \alpha(q)$).*

SKETCH OF PROOF. Call a variable *repeated* if it occurs in at least two V -tuples in Q , s , and a *target* variable if it occurs in s (see (8.3)). For every repeated variable x , let A_x be an attribute not appearing in Q . Let $Q = \langle Q_1, \dots, Q_n \rangle$. For any V -tuple u in Q , let f_u be an expression consisting of a relation symbol R , preceded by (from right to left)

- (i) selections $\sigma_{A=a}$, for every $A \in \alpha(u)$ with $u(A) = a$, $a \in D$,
- (ii) restrictions $\sigma_{A=B}$, for all $A, B \in \alpha(u)$ with $u(A) = u(B) = x$, a variable,
- (iii) projection π_Y , where Y contains exactly one attribute A for any repeated variable x appearing in u ($u(A) = x$),
- (iv) renaming $s_{A_x}^u$, for every $A \in Y$, where $x = u(A)$.

Let $h = \pi_X(\bowtie_u f_u)$, where u ranges over all V -tuples in Q and $X = \{A_x : x \text{ is a target variable}\}$. The required expression f is obtained from h by suitably renaming the attributes in X . Some easy modifications are required if s contains constants or multiple occurrences of the same variable. (A similarity of this construction to "shallow expressions" of [32] may be noted.) \square

By Lemma 8.1, we obtain the following corollary.

COROLLARY 8.2. *The representation system $(\mathcal{F}_V, \text{Rep}, PS^+UJR)$ can support arbitrary conjunctive queries; that is, for any conjunctive query q and any V -multitable T with $\alpha(T) = \alpha(q)$ there is a V -table U such that $\text{Rep}(U) \equiv_{PS^+UJR} q(\text{Rep}(T))$.*

Here we denote $q(\text{Rep}(T)) = \{q(\mathbf{r}) : \mathbf{r} \in \text{Rep}(T)\}$. The V -table U is obtained by finding the PS^+JR -expression f equivalent to q (see Lemma 8.1) and putting $U = f(T)$. The V -table U will also be denoted by $q(T)$.

By using the fact that in computing $f(T)$ we treat variables and constants in exactly the same way, and that the equivalence between q and f does not depend on the nature of the attribute domains, one can easily prove the following fact generalizing (8.3):

THEOREM 8.3. For any conjunctive query q , given by a V -multitable \mathbf{Q} and V -tuple s , and for any V -multitable \mathbf{T} with $\alpha(\mathbf{T}) = \alpha(\mathbf{Q})$

$$q(\mathbf{T}) = \{\rho(s) : \rho \text{ is a generalized valuation such that } \rho(\mathbf{Q}) \subseteq \mathbf{T}\},$$

where by a generalized valuation we mean any mapping $\rho: V \rightarrow V \cup D$ such that $x \in V(A)$ implies $\rho(x) \in V(A) \cup D(A)$.

9. Closed World Interpretation of Tables

The interpretation of multitable considered so far has been based on the open-world assumption [26, 27], which was embodied in the property

$$r \in \text{Rep}(\mathbf{T}) \wedge s \supseteq r \Rightarrow s \in \text{Rep}(\mathbf{T})$$

(see (3.3)). We now consider a different interpretation, which we call the *closed world interpretation*.

Let \mathbf{T} be a V -multitable. The closed world interpretation of \mathbf{T} is defined as follows:

$$\text{rep}(\mathbf{T}) = \{r : \text{there exists a valuation } v \text{ such that } v(\mathbf{T}) = r\}, \quad (9.1)$$

where $v(\mathbf{T})$ is defined as before (see (6.1)–(6.3)). Similarly, for any C -multitable \mathbf{T} we define $\text{rep}(\mathbf{T})$ by formula (9.1), with $v(\mathbf{T})$ defined by (7.1) and (7.2). In the case of a Codd multitable \mathbf{T} , we define

$$\text{rep}(\mathbf{T}) = \text{rep}(\mathbf{U}),$$

where \mathbf{U} is any V -multitable obtained by replacing every occurrence of $@$ in \mathbf{T} by a different variable. It seems that what Codd meant in [6] was in fact the closed world interpretation of Codd tables.

Notice that if a tuple t cannot be obtained by replacing nulls by constants in any $@$ -tuple of a Codd table \mathbf{T} , then \mathbf{T} represents the “negative information” that the relationship expressed by t definitely does not hold, a fact not representable under the open world interpretation. The situation is similar in the case of V -tables and C -tables.

Clearly, for all three kinds of multitable,

$$\text{rep}(\mathbf{T}) \subseteq \text{Rep}(\mathbf{T});$$

moreover, $\text{rep}(\mathbf{T})$ and $\text{Rep}(\mathbf{T})$ are coinital.

Let $\Omega \subseteq \text{PSS}^+ \text{UJR}$, f be an Ω -expression, \mathbf{T} be a multitable of type $\alpha(f)$, and \mathbf{U} be a multitable of type $\beta(f)$. Since $\text{rep}(\mathbf{U}) \approx \text{Rep}(\mathbf{U})$, $\text{rep}(\mathbf{T}) \approx \text{Rep}(\mathbf{T})$, by Lemma 4.2 we obtain

$$\text{rep}(\mathbf{U}) \equiv_{\Omega} \text{Rep}(\mathbf{U})$$

and

$$f(\text{rep}(\mathbf{T})) \equiv_{\Omega} f(\text{Rep}(\mathbf{T})).$$

Hence,

$$\text{rep}(\mathbf{U}) \equiv_{\Omega} f(\text{rep}(\mathbf{T})) \quad \text{iff} \quad \text{Rep}(\mathbf{U}) \equiv_{\Omega} f(\text{Rep}(\mathbf{T})).$$

This equivalence immediately implies that all results concerning representation systems developed under the open-world assumption carry over to the closed-world interpretation:

THEOREM 9.1. *Let $\Omega \subseteq PSS^+UJR$ and let \mathcal{T} be either \mathcal{T}_\emptyset or \mathcal{T}_V or \mathcal{T}_C . Then*

$$\langle \mathcal{T}, Rep, \Omega \rangle \text{ is a representation system} \\ \text{iff } \langle \mathcal{T}, rep, \Omega \rangle \text{ is a representation system.}$$

Moreover, the correct operations on tables are the same under both interpretations.

It turns out that in the case of conditional tables we can handle all relational operations, including the difference, thus supporting the full strength of the relational algebra. Indeed, let T, U be two C -tables, $\alpha(T) = \alpha(U) = X$. For any $t \in T$, let us define t_U to be a C -tuple, $\alpha(t_U) = X$, such that

$$t_U[X] = t[X], \\ t_U(\text{con}) = t(\text{con}) \wedge \bigwedge_{u \in U} \bigvee_{A \in X} (t(A) \neq u(A)).$$

If we define

$$T - U = \{t_U : t \in T\},$$

then it is easily seen that

$$v(T - U) = v(T) - v(U)$$

for any valuation v , which, combined with (7.11), implies that

$$v(f(\mathbf{T})) = f(v(\mathbf{T}))$$

for any PSUJRD-expression f , any C -multitable \mathbf{T} ($\alpha(\mathbf{T}) = \alpha(f)$), and any valuation v . Consequently,

$$\text{rep}(f(\mathbf{T})) = f(\text{rep}(\mathbf{T})), \tag{9.2}$$

that is, diagram (3.1) commutes. Let us state this as a theorem.

THEOREM 9.2. *$\langle \mathcal{T}_C, rep, PSUJRD \rangle$ is a representation system. Moreover, all relational operators can be defined on C -tables in such a way that (9.2) holds.*

10. Conclusions

We have proposed a general condition that should be satisfied if we want to evaluate relational expressions over “tables with nulls” in a semantically correct way.

Suppose that an instance of an incomplete information relational database is given by a multitable \mathbf{T} , and let f be a relational query. While computing the response to f , we think of an unknown $r^* \in \text{Rep}(\mathbf{T})$ —corresponding to the true state of the real world—being transformed by f , and we require that $f(\mathbf{T})$ contain enough information to determine all tuples that surely (i.e., no matter which element of $\text{Rep}(\mathbf{T})$ the multirelation r^* is) appear in $f(r^*)$.

In this way our approach is similar in spirit to the *external interpretation* and *lower value* $\|\cdot\|_*$ in [20] (see also [22]).

Note that the tuples that surely appear in $f(r^*)$ are given by $\text{Rep}(\mathbf{T})^f$, and it is easy to see that $\text{Rep}(\mathbf{T})^f$ can be obtained by deleting from $f(\mathbf{T})$ all tuples containing nulls, both in the case of Codd tables and V -tables.

Another possibility is to treat $f(\mathbf{T})$, rather than $\text{Rep}(\mathbf{T})^f$, as the response to query f . This has the advantage of providing more information; for instance, a tuple $a@c$ (or ayc) appearing in the response can easily be shown to indicate that a tuple of the form abc , $b \in D(B)$ surely appears in $f(r^*)$. Moreover, this approach makes

it possible to treat a response to query f as an intermediate step in computing the correct response to a more complex query gf (note that in general it is *not* possible to determine $\text{Rep}(T)^{gf}$ on the basis of $\text{Rep}(T)^f$ and g).

It should be noted that $\text{Rep}(T)^f = \text{rep}(T)^f$, so that by Theorem 9.1 the process of computing the response to a query is exactly the same under both open and closed world assumptions.

The results of the paper show that the devices suitable for representing incomplete information heavily depend on what processing of the information we are going to correctly perform, or, in more concrete terms, what relational operators are allowed. Our approach stresses the requirement for a correct evaluation of relational expressions, rather than just single relational operators, as was usually the case in previous attempts in the literature to correctly handle null values.

From the practical point of view an especially appealing system seems to be one based on V -tables, supporting projection, positive selection, union, join, and renaming, which allows for processing arbitrary conjunctive queries. The reason is that all these relational operators and queries can be evaluated over V -tables in *exactly the same way* as in the case of the usual relations.

There are several important problems that have not been treated here. One is how to handle dependencies in our framework. It turns out that functional, multivalued, and join dependencies—more exactly, arbitrary implicational dependencies [2] (“generalized dependencies” in [29])—can be *represented* in any V -table, in the sense that for any V -table T and any set Σ of such dependencies, there is a V -table U such that $\text{Rep}(U) \equiv_{\text{PS}+\cup} \text{Rep}(T) \cap \text{Sat}(\Sigma)$, where $\text{Sat}(\Sigma)$ is the set of all relations satisfying the dependencies in Σ (see [12, 13, 16]). This is another argument for the usefulness of the representation system based on V -tables.

Another interesting topic is the relation between V -tables and the tableaux of Aho, Sagiv, and Ullman [1]. Roughly speaking, a tableau for a relational PSJ-expression f (with only selection of the form $\sigma_{A=a}$ allowed) can be treated as a V -table representing the inverse image of the summary of the tableau with respect to f (see [15]). It also turns out that for any PSJ-expression f (with only selection of the form $\sigma_{A=a}$) and every V -multitable T ($\alpha(T) = \alpha(f)$), there is a V -multitable U such that $\text{Rep}(U) = f^{-1}(\text{Rep}(T))$. This fact has several interesting applications (see [17]).

ACKNOWLEDGMENTS. Several comments and suggestions by J. Loś, W. M. Turski, and the referees are gratefully acknowledged.

REFERENCES

1. AHO, A. V., SAGIV, Y., AND ULLMAN, J. D. Equivalences among relational expressions. *SIAM J Comput* 8, 2 (May 1979), 218–246.
2. BEERI, C., AND VARDI, M. Y. Formal systems for tuple and equality generating dependencies. *SIAM J Comput* 13, 1 (Feb. 1984), 76–98.
3. BISKUP, J. A formal approach to null values in database relations. In *Advances in Database Theory*, H. Gallaire, J. Minker, and J. M. Nicolas, Eds. Plenum Press, New York, 1981, pp. 299–341.
4. CHANDRA, A. K., AND MERLIN, P. M. Optimal implementation of conjunctive queries in relational data bases. In *Proceedings of the 9th Annual ACM Symposium on Theory of Computing* (Boulder, Col., May 2–4). ACM, New York, 1977, pp. 77–90.
5. CODD, E. F. A relational model for large shared data banks. *Commun ACM* 13, 6 (June 1970), 377–387.
6. CODD, E. F. Relational completeness of data base sublanguages. In *Data Base Systems*, R. Rustin, Ed. Prentice-Hall, Englewood Cliffs, N.J., 1972, pp. 65–98.

7. CODD, E. F. Understanding relations (Installment #7). *FDT Bull. of ACM-SIGMOD* 7, 3-4 (Dec. 1975), 23-28.
8. CODD E. F. Extending the database relational model to capture more meaning. *ACM Trans Database Syst* 4, 4 (Dec. 1979), 397-434.
9. GAREY, M. R., AND JOHNSON, D. S. *Computers and Intractability: A Guide to the Theory of NP-completeness*. Freeman, San Francisco, Calif., 1979.
10. GRANT, J. Null values in a relational data base. *Inf. Process. Lett.* 6, 5 (Oct. 1977), 156-157.
11. HOROWITZ, E., AND SAHNI, S. *Fundamentals of Computer Algorithms*. Computer Science Press, Potomac, Md., 1979.
12. IMIELŃSKI, T. Problems of representing information in relational databases (in Polish). Ph.D. Thesis, Institute of Computer Science, Polish Academy of Sciences, 1981.
13. IMIELŃSKI, T., AND LIPSKI, W. On representing incomplete information in a relational data base. In *Proceedings of the 7th International Conference on Very Large Data Bases* (Cannes, France, Sept. 9-11) ACM, New York, 1981, pp. 388-397.
14. IMIELŃSKI, T., AND LIPSKI, W. The relational model of data and cylindric algebras. *J. Comput. System Sci.* 28, 1 (Feb. 1984), 80-102.
15. IMIELŃSKI, T., AND LIPSKI, W. A technique for translating states between database schemata. In *Proceedings of the ACM SIGMOD International Conference on Management of Data* (Orlando, Fla., June 2-4). ACM, New York, 1982, pp. 61-68.
16. IMIELŃSKI, T., AND LIPSKI, W. Incomplete information and dependencies in relational databases. In *Proceedings of the ACM SIGMOD International Conference on Management of Data* (San Jose, Calif., May 23-26). ACM, New York, 1983, pp. 178-184.
17. IMIELŃSKI, T., AND LIPSKI, W. Inverting relational expressions—a uniform and natural technique for various database problems. In *Proceedings of the ACM SIGACT-SIGMOD Symposium on Principles of Database Systems* (Atlanta, Ga., March 21-23). ACM, New York, 1983, pp. 305-311.
18. LACROIX, M., AND PIROTTE, A. Generalized joins. *ACM-SIGMOD Record* 8, 3 (Sept. 1976), 14-15.
19. LIPSKI, W. Informational systems with incomplete information. In *Proceedings of the 3rd International Colloquium on Automata, Languages and Programming* (Edinburgh, Scotland, July 20-23). Edinburgh University Press, Edinburgh, Scotland, 1976, pp. 120-130.
20. LIPSKI, W. On semantic issues connected with incomplete information databases. *ACM Trans Database Syst* 4, 3 (Sept. 1979), 262-296.
21. LIPSKI, W. On databases with incomplete information. *J. ACM* 28, 1 (Jan. 1981), 41-70.
22. LIPSKI, W. Logical problems related to incomplete information in databases. Tech. Rep. 138, Laboratoire de Recherche en Informatique, Université de Paris-Sud, Centre d'Orsay, Sept. 1983.
23. LIPSKI, W. On relational algebra with marked nulls. In *Proceedings of the 3rd ACM SIGACT-SIGMOD Symposium on Principles of Database Systems* (Waterloo, Ont., Canada, April 2-4). ACM, New York, 1984, pp. 201-203.
24. MAIER, D., ULLMAN, J. D., AND VARDI, M. Y. On the foundations of the universal relation model. *ACM Trans Database Syst.* 9, 2 (June 1984), 283-308.
25. REINGOLD, E. M., NIEVERGELT, J., AND DEO, N. *Combinational Algorithms: Theory and Practice*. Prentice Hall, Englewood Cliffs, N.J., 1977.
26. REITER, R. On closed world databases. In *Logic and Data Bases*, H. Gallaire and J. Minker, Eds. Plenum Press, New York, 1978, pp. 55-76.
27. REITER, R. Towards a logical reconstruction of relational database theory. In *Conceptual Modelling, Perspectives from Artificial Intelligence, Databases and Programming Languages*, M. L. Brodie, J. Mylopoulos, and J. Schmidt, Eds. Springer-Verlag, New York, 1984, pp. 191-233.
28. SIKLÓSSY, L. Efficient query evaluation in relational databases with missing values. *Inf. Process. Lett* 13, 4/5 (End 1981), 160-163.
29. ULLMAN, J. D. *Principles of Database Systems*. 2nd ed. Computer Science Press, Potomac, Md., 1982.
30. VASSILIOU, Y. Null values in data base management: A denotational semantics approach. In *Proceedings of the ACM-SIGMOD International Symposium on Management of Data* (Boston, Mass., May 30-June 1). ACM, New York, 1979, pp. 162-169.
31. VASSILIOU, Y. Functional dependencies and incomplete information. In *Proceedings of the 6th International Conference on Very Large Data Bases* (Montreal, Ont., Canada, Oct. 1-3). ACM, New York, 1980, pp. 260-269.
32. YANNAKAKIS, M., AND PAPADIMITRIOU, C. H. Algebraic dependencies. *J. Comput. Syst. Sci.* 25, 1 (Aug. 1982), 2-41.

33. ZANIOLO, C. Relational views in a data base system support for queries. In *Proceedings of the IEEE Computer Software and Applications Conference* (Chicago, Ill., Nov. 8–11). IEEE, New York, 1977, pp. 267–275.
34. ZANIOLO, C. Database relations with null values. In *Proceedings of the ACM SIGACT-SIGMOD Symposium on Principles of Database Systems* (Los Angeles, Calif., March 29–31). ACM, New York, 1982, pp. 27–33.

RECEIVED APRIL 1982; REVISED MARCH 1984; ACCEPTED MARCH 1984