# DEPENDENCY SATISFACTION IN DATABASES WITH INCOMPLETE INFORMATION

Gösta Grahne

University of Helsinki, Department of Computer Science
Tukholmankatu 2, SF-00250 Helsinki 25, Finland

Abstract. Two of the major problems raised by information incompleteness in databases are how to evaluate queries and how to take data dependencies into account. We give a unified solution of these two intermingled problems for the relational model. Formal criteria for the correctness of the relational algebra and dependency satisfaction are presented. We give a correct redefinition of the complete relational algebra and present a method, called a chase, for enforcing a set of functional and full join dependencies on a relation with null-values of type "value exists, but is presently unknown". This novel chase can also be regarded as a generalization of previously known chase methods. The title of the paper reflects the emphasis of its contribution.

## 1. INTRODUCTION

The research in the field of information incompleteness has mainly followed three paths. The vast majority of the papers on the topic have sought a way to adapt the query language to handle null-values of different kinds (a.o. [Bis1, Bis2, IL1, IL2, Vas1]), while only a few papers (i.e. [IL3, Lie, Vas2]) have looked at the problems pointed at by the handling of data dependencies under information incompleteness. Even less

---

papers have orientated towards a solution to the problem of updating databases with incomplete information (i.e. [FUV, Sci]). All these directions are however only different aspects of the same problem: how to interpret and accordingly handle incomplete information. (For an overview of the field, see [Gra1] or [Lip].)

This paper offers a unified treatment of incomplete information of the type "value exists, but is presently unknown". Explicit solutions are given for the query problem and the data dependencies in the context of the relational model. Updates will be treated elsewhere [Gra2].

Following [IL2] and [Bis2] we consider a relation containing null-values of the aforementioned type as representing a set of relations of the ordinary kind (i.e. without nulls), one of which corresponds to the state of the real world. Since our information is incomplete we only know the set of possible states, and we store its representative. The query language should then be designed so that only information that holds in all of the possible states is inferred. In [IL2] it is shown how to extend a relational algebra consisting of projection, positive selection, union and join so that it is correct in this sence. A selection is positive if its formula only includes atomic expressions of the form $A=a$ combined with conjunction and disjunction ($A$ is an attribute name, $a$ some domain value). This algebra is quite restrictive, since it would for example be impossible to ask for a list of all em-

ployees earning more than 20000$ in some enterprise database. The authors of [IL2] show that the handling of the complete relational algebra requires a more complicated representative than just putting null-values in ordinary relations. Such a representative, called a c-table, and a correct redefinition of the complete relational algebra for these tables is given in [IL2].

Some authors have tried to solve the problem of data dependencies by incorporating the null-values in the definitions of the data dependencies. The dependencies are however not statements about the incomplete database, but statements about the real world. Thus their definitions remain unchanged, and the question of satisfaction is determined by the set of possible relations, in which each individual relation should satisfy the dependencies. But since we only store the representative relation with null-values, we need some syntactic method for checking dependency satisfaction. This syntactic satisfaction will, in absence of better names, be called "satisfaction".

The argumentation above also shows that the data dependencies affect query evaluation: the incomplete relation determines a set of possible relations, and this set is further restricted by the data dependencies. Since query evaluation is based on the set of possible relations, we must, if we want to maintain the correctness of the query language, restrict the set of possible relations to those relations that satisfy the data dependencies. In [IL3] it is shown how to take into account data dependencies in the system where only positive selections are allowed. Since we feel that this system is too restrictive, we present a system that supports the complete relational algebra and functional and full join dependencies. Chapter 2 of this paper contains the necessary preliminaries and our model that is based on so called tables, and Chapter 3 contains a formal criteria for the correctness of a relational algebra on tables, along with the correct

redefinition of the complete relational algebra. In Chapter 4 we give a formal criteria for dependency satisfaction in tables, and we present a method for transforming a table so that the corresponding set of relations will be restricted to those satifying a given set of functional and full join dependencies. The transformation procedure is called a chase, and it can also be regarded as a generalization of previously known chase methods.

2. PRELIMINARIES

For the definitions of the basics of the relational model we refer to [Mai, Ull]. Relation schemes will be denoted by $R_1$, $R_2$, ..., and their instances or relations by $r_1$, $r_2$, ..., where $r_i$ is a relation over $R_i$. The set of all relations over $R_i$ is Rel($R_i$). For notational convenience we assume one (infinite) common attribute domain D. A tuple is denoted by t, possibly indexed or primed, and t(X) is the restriction of t to the attributes X. The projection operator is denoted by π and the join operator by *. By a relational expression f we mean a well formed expression built up from projection, selection, union, join and difference. A PJ-expression is a relational expression involving only projection and join, for instance $\pi_{AB}(R_1)*\pi_{BC}(R_2)$. In general, if f is performed on relations $r_1$, $r_2$, ..., $r_n$ we can write $f(r_1, r_2, ..., r_n)$. The result of performing f is always a single relation of appropriate type.

X->A denotes that X functionally determines A, where A is a single attribute. The join dependency for the set $R = \{R_1, R_2, ..., R_p\}$ is denoted by *[R] or *[$R_1$, $R_2$, ..., $R_p$]. The join dependency *[R] is full for the scheme R, if $\bigcup_{i=1}^{p} R_i =$ R. A set of functional dependencies (FDs) and full join dependencies (JDs) is denoted by Σ. The set of JDs in Σ is denoted by $\Sigma_{JD}$. It is assumed that Σ is defined for some scheme R. Sat(Σ) is

Proceedings of the Tenth International
Conference on Very Large Data Bases.

Singapore, August, 1984

38

then the set of relations over R that satisfy the dependencies in $\Sigma$.

The device for representing a <u>set</u> of relations is called a <u>table</u>. Our tables are essentially extensions of so called c-tables of [IL2]. A table T consists of two parts, a set of <u>c-tuples</u> and a set of <u>general conditions</u>, denoted $T_C$ and $T_G$. For building up a table we need a set $G$ of all expressions with atoms of the form (x=a), (x=y), <u>true</u> and <u>false</u> combined by $\wedge$, $\vee$, $\neg$ and $\Rightarrow$ (conjunction, disjunction, negation and implication). In the atoms x and y are <u>variables</u> from a countably infinite set V, and a is a value from the domain D. The domain values will be called <u>constants</u>. The variables are used to express the null-values, and we assume that $V \cap D = \emptyset$. A c-tuple t over R is a mapping from $R \cup \{Con\}$ to $V \cup D \cup G$, such that $t(A) \in V$ or $t(A) \in D$ for $A \in R$, and $t(Con) \in G$. A c-table $T_C$ over R is a finite set of c-tuples over R. $T_G$ is a finite set of general conditions $\{g | g \in G\}$. By definition, <u>true</u> is always a member of $T_G$. The set of all tables over R is denoted by Tab(R). A <u>multitable</u> $T$ is a sequence $\langle T_1, T_2, \ldots, T_n \rangle$ of tables $T_i \in Tab(R_i)$, such that all tables have the same set of general conditions $T_G$. Tab($R$) is the set of all multitables over $R$.

A <u>valuation</u> v is a mapping from $V \cup D$ to D, such that for $x \in V$, $v(x) = a$ for some $a \in D$, and $v(a) = a$ for all $a \in D$. Valuations are extended to c-tuples by defining $(v(t))(A) = v(t(A))$ for $A \in R$ and $(v(t))(Con) = $ <u>true</u> if t(Con) comes to <u>true</u> when all variables x in t(Con) are substituted by v(x). Else $(v(t))(Con)$ is <u>false</u>. The valuation v applies to the general conditions in $T_G$ in the same way, and $v(T_G) = $ <u>true</u> if $v(g) = $ <u>true</u> for all $g \in T_G$.

The set of relations that a table T represents is denoted by Rep(T) and it is defined in the following way:

**Definition 1.** For $T \in Tab(R)$, Rep(T) = $\{r \mid r \in Rel(R)$ and there exists a valuation v such that r = $v(T_C)$ and $v(T_G) = $ <u>true</u>$\}$, where $v(T_C) = \{v(t(R)) \mid t \in T_C$ and $v(t(Con)) = $ <u>true</u>$\}$. For a multitable $T = \langle T_1, \ldots, T_n \rangle$, $T_i \in Tab(R_i)$ Rep($T$) = $\{\langle r_1, \ldots, r_n \rangle \mid r_i \in Rel(R_i)$ and there exists a valuation v such that $r_i = v(T_{i_C})$ and $v(T_G) = $ <u>true</u>$\}$.

Choosing $r = v(T_C)$ means that we make the <u>closed world assumption</u> [Rei]. The <u>open world assumption</u> would correspond to $r \supseteq v(T_C)$.

The next example shows a table and two of the relations that it represents. Note that the set Rep(T) in this and most other cases is infinite.

<u>Example 1.</u> $T \in Tab(SUPPLIER~PART~PROJECT)$

$T_C$(SUPPLIER PART PROJECT Con)

| Jones | x | steel | x=bolt∨x=nut |
| Smith | y | z | <u>true</u> |

$T_G = \{$<u>true</u>, x=y, z=concrete$\}$

$\{r_1, r_2\} \subseteq$ Rep(T)

$r_1$(SUPPLIER PART PROJECT)

| Jones | bolt | steel |
| Smith | bolt | concrete |

$r_2$(SUPPLIER PART PROJECT)

| Smith | nail | concrete |

Given two tables $T_1$ and $T_2$ over the same scheme, we say that $T_2$ <u>contains</u> $T_1$, den. $T_1 \sqsubseteq T_2$, if Rep($T_1$) $\subseteq$ Rep($T_2$), and that $T_1$ is <u>equivalent</u> to $T_2$, den. $T_1 \equiv T_2$, if $T_1 \sqsubseteq T_2$ and $T_2 \sqsubseteq T_1$. The equivalence of two tables means that they define the same set of relations, since we clearly have $T_1 \equiv T_2$ if Rep($T_1$) = Rep($T_2$). Also, $v(T_{1_C}) = v(T_{2_C})$ and $v(T_{1_G}) = v(T_{2_G})$ for all valuations v necessitates $T_1 \equiv T_2$.

A condition g (or a set of conditions $T_G$) implies a condition g', den. $g \Rightarrow g'$, if $v(g) = $

true implies $v(g') = \underline{true}$ (or $v(T_G) = \underline{true}$ implies $v(g') = \underline{true}$) for all valuations v. Mutual implication is denoted by $\simeq$ (equivalence). The following underline{normalization} underline{rules} can now be applied to tables.

(1) If for some $t \in T_C$, $\neg(T_G \wedge t(\text{Con}))$, then $T_C$ is replaced by $T_C - \{t\}$.

(2) If there exists c-tuples $t_1, \ldots, t_k \in T_C$, $T \in \text{Tab}(R)$, such that $t_1(R) = \ldots = t_k(R)$, then $T_C$ is replaced by $T_C - \{t_1, \ldots, t_k\} \cup \{t\}$, where $t(R) = t_1(R)$ and $t(\text{Con}) = \overset{k}{\underset{i=1}{\vee}} t_i(\text{Con})$.

A table T is said to be underline{normalized}, den. $T^o$, if none of these normalization rules can be applied to T. (There are also other normalization rules, but these two are sufficient for our purposes.) Normalization of T does not affect Rep(T). We have $T \equiv T^o$. Replacing conditions in c-tuples or in $T_G$ by equivalent ones also preserves equivalence.

For some purposes we need a form of set inclusion for tables. This underline{modified} underline{inclusion} is defined followingly:

underline{Definition 2}. Let $T_1$ and $T_2$ be two tables from Tab(R). $T_1 \subseteq_{\overline{\tau}} T_2$, $T_2$ underline{m-includes} $T_1$ if

(1) for each $g \in T_{1_G}$ there exists a $g' \in T_{2_G}$ such that $g \simeq g'$, and

(2) for each $t_1 \in T_{1_C}$ there exists a $t_2 \in T_{2_C}$ such that $t_1(R) = t_2(R)$ and $t_1(\text{Con}) \simeq t_2(\text{Con})$.

If both $T_1 \subseteq_{\overline{\tau}} T_2$ and $T_2 \subseteq_{\overline{\tau}} T_1$, we say that $T_1$ is underline{m-equal} to $T_2$ and denote it $T_1 \mp T_2$. Clearly m-equality implies equivalence, i.e. $T_1 \mp T_2$ implies $T_1 \equiv T_2$, but note that $T_1 \subseteq T_2$ does underline{not} imply $T_1 \subseteq_{\overline{\tau}} T_2$.

## 3.  EXTENSION OF THE RELATIONAL ALGEBRA

As mentioned in the introduction the correctness criteria for an extension of the algebra to a table T is that we only conclude information that holds in every possible state of the real world, i.e. in every relation in Rep(T). If we denote the extension of a relational expression f by $\hat{f}$ we can formalize the correctness criteria as

$$f(\text{Rep}(T)) = \text{Rep}(\hat{f}(T))$$

for all relational expressions f and multitables $T$. The lefthand side of the equality stands for $\{f(r_1, \ldots, r_n) | <r_1, \ldots, r_n> \in \text{Rep}(T)\}$, and the multitable $T$ is taken as $<T_1, \ldots, T_n>$ of the tables $T_i$ that $\hat{f}$ is applied on. The result $\hat{f}(T)$ is a table of appropriate type. Our notation reveals that we perform non-unary operations only between tables with the same set of general conditions, i.e. between individual tables of a multitable.

Imielinski and Lipski [IL2] have given a correct extension of the relational algebra for c-tables. The same extension can with a slight modification be used for tables also. Since we in the sequel only will need PJ-expressions we will be contended with giving only the definitions for the extension of the project and join operators.

underline{Definition 3}. For a table $T \in \text{Tab}(R)$ the underline{projection} $\hat{\pi}_X(T)$ on a set $X \subseteq R$ is a table $T'^o \in \text{Tab}(X)$ such that $T'_C = \{t(X \cup \text{Con}) | t \in T_C\}$ and $T'_G = T_G$.

underline{Definition 4}. For tables $T_1 \in \text{Tab}(R_1)$ and $T_2 \in \text{Tab}(R_2)$ with the same set of general conditions the underline{join} $T_1 \overset{\spadesuit}{} T_2$ of $T_1$ and $T_2$ is a table $T'^o \in \text{Tab}(R_1 \cup R_2)$ such that $T'_G = T_{1_G}$ and $T'_C = \{t_1 \overset{\spadesuit}{} t_2 | t_1 \in T_{1_C}$ and $t_2 \in T_{2_C}$, where $t_1 \overset{\spadesuit}{} t_2$ is the c-tuple over $R_1 \cup R_2$ with

$$(t_1 \overset{\spadesuit}{} t_2)(A) = \begin{cases} t_1(A), & \text{if } A \in R_1 \\ t_2(A), & \text{if } A \in R_2 - R_1 \end{cases}$$

Proceedings of the Tenth International
Conference on Very Large Data Bases.

Singapore, August, 1984

40

$$(t_1 \overset{\Phi}{} t_2)(Con) = t_1(Con) \wedge t_2(Con) \wedge \underset{A \in R_1 \cap R_2}{\wedge} (t_1(A) = t_2(A)).$$

As an illustration, consider

Example 2.

$$T_{1_C} \frac{(A\ B\ C\ Con)}{\begin{array}{cccc} a & x & c & \neg(x=b) \\ a' & b & y & \underline{true} \end{array}} \qquad T_{2_C} \frac{(C\ D\ E\ Con)}{\begin{array}{cccc} c & d & e & \underline{true} \\ z & d' & e' & \underline{\neg(z=c'')} \\ c' & d' & e & \underline{true} \end{array}}$$

$$T_{1_G} = T_{2_G} = \{\underline{true},\ \neg(y=c)\}$$

$$(T_1 \overset{\Phi}{} T_2)_C \frac{(A\ B\ C\ D\ E\ Con)}{\begin{array}{ccccccc} a & x & c & d & e & \neg(x=b) \wedge true \wedge c = c \\ a & x & c & d' & e' & \neg(x=b) \wedge \neg(z=c'') \wedge c = z \\ a' & b & y & d' & e' & true \wedge \neg(z=c'') \wedge y = z \\ a' & b & y & d' & e & \underline{true \wedge true \wedge y = c'} \end{array}}$$

$$(T_1 \overset{\Phi}{} T_2)_G = \{\underline{true},\ \neg(y=c)\}$$

The conditions of the first and last c-tuples of $(T_1 \overset{\Phi}{} T_2)_C$ can be replaced by $\neg(x=b)$ and $y=c$. Also, the un-normalized join of $T_1$ and $T_2$ contains the c-tuple $\langle a', b, y, d, e, \underline{true \wedge true \wedge y = c} \rangle$, but since $(\neg(y=c) \wedge y=c) \simeq \underline{false}$ the tuple is removed. For a similar reason the c-tuple $\langle a, x, c, d', e, \neg(x=b) \wedge true \wedge c = c' \rangle$ does not belong to the normalized result.

The rest of the relational operators can be extended to operate on tables along the same lines as the project and join operators, and the following theorem follows easily from Theorem 9.2 of [IL2].

Theorem 1. For any well formed expression f built up from projection, selection, union, join and difference, and multitables $T$, we have
$f(Rep(T)) = Rep(\overset{\wedge}{f}(T)).$ $\square$

In the ordinary relational algebra a relation r is always included in the result of certain PJ-expressions on r. A similar property will be needed for the extended algebra.

Lemma 1. Let $R = R_1, \ldots, R_p$, $\overset{p}{\underset{i=1}{\cup}} R_i = R$ and $T \in$ Tab(R). Then $T \underset{\top}{\subseteq} \overset{\wedge}{\pi}_{R_1}(T) \overset{\Phi}{} \ldots \overset{\Phi}{} \overset{\wedge}{\pi}_{R_p}(T).$

Proof. Property (1) of Definition 2 is immediate since $T_G$ is not changed. For property (2), take an arbitrary $t \in T_C$. By Definition 3 $t(R_i \cup Con) \in \overset{\wedge}{\pi}_{R_i}(T)$, for $i=1,\ldots,p$. By inductive use of Definition 4 $\overset{\wedge}{\pi}_{R_1}(T) \overset{\Phi}{} \ldots \overset{\Phi}{} \overset{\wedge}{\pi}_{R_p}(T)$ will include a c-tuple $t'$ with $t'(R) = t(R)$ and $t'(Con) = (t(Con) \wedge t(Con) \wedge \underset{A \in R_1 \cap R_2}{(t(A)=t(A))} \wedge \ldots \wedge \underset{A \in R_{p-1} \cap R_p}{t(Con)} \wedge t(A)=t(A)))$. Now $t'(Con)$ is clearly equivalent to $t(Con)$, so the result follows. $\square$

## 4. EXTENSION OF DEPENDENCY SATISFACTION

A set $\Sigma$ of dependencies restricts the relations that model the real world to those in Sat($\Sigma$). For a table T we make the natural interpretation that T "satisfies" $\Sigma$ if

$$Rep(T) \subseteq Sat(\Sigma).$$

Thus we need a method for transforming T in order to cut down Rep(T) to Rep(T)$\cap$Sat($\Sigma$) for any table T and set $\Sigma$ of FDs and JDs. This cutting down is necessary, since the extended algebra operates on the basis of Rep(T), and the set $\Sigma$ implies that not necessarily all members of Rep(T) are possible states of the real world.

The method is called a chase, and its idea will be clarifyed by concidering the algebraic counterparts of FDs and JDs. These counterparts are special cases of so called algebraic dependencies of [YP] (see also [Abi]). Stated with a momentary simplification, an algebraic dependency is a statement of the form $f(R) \subseteq R$, where f is a PJ-expression. A relation r over R satisfies such a dependency if $f(r) \subseteq r$. Since we demand the inclusion to hold for all members of Rep(T), where T$\in$Tab(R), we can using the definition of containment and Theorem 1 write the following deduction chain:

T "satisfies" $f(R) \subseteq R$ if and only if

$\quad f(Rep(T)) \subseteq Rep(T)$     if and only if

$\quad Rep(\hat{f}(T)) \subseteq Rep(T)$     if and only if

(*)    $\hat{f}(T) \sqsubseteq T$

Our chase is essentially a method for enforcing the containment (*) by applying the expression $\hat{f}$ on T. The enforcement will always be possible, even if $Rep(T) \cap Sat(\Sigma) = \emptyset$. There are two reasons for the intersection to be empty. First, it might be that every relation in $Rep(T)$ violates a FD in $\Sigma$. Then we must make $Rep(T)$ equal to the empty set. Second, the relations in $Rep(T)$ might violate some JDs of $\Sigma$. In this case we will add the necessary c-tuples to T, and thus we have in fact done more than just cutting down $Rep(T)$ to $Rep(T) \cap Sat(\Sigma)$. We shall be more precise after introducing some necessary definitions.

First we concider completions of relations: Let $r \in Rel(R)$. The underline{completion} of r with respect to $\Sigma$ is $Comp_\Sigma(r) = s$, such that $s \supseteq r$ and $s \in Sat(\Sigma_{JD})$, and there is no relation $s' \in Sat(\Sigma_{JD})$ with $s \not\supseteq s' \supseteq r$. $Comp_\Sigma(Rep(T)) = \{Comp_\Sigma(r) \mid r \in Rep(T)\}$. The completion means that the tuples that are necessary for making r satisfy $\Sigma_{JD}$ are added in s. The completion always exists and it is unique (see [Mai, exercise 8.41]). We are also free to choose any cover for $\Sigma$, i.e. if $Sat(\Sigma) = Sat(\Sigma^J)$ then $Comp_\Sigma(Rep(T)) = Comp_{\Sigma^J}(Rep(T))$.

For the definition of algebraic dependencies we need the notions of extended schemes and relations. An underline{extended relation scheme} $\overline{R}$ of R has two copies of every attribute. For instance, if $R = ABCD$ then $\overline{R} = ABCDABCD$. An underline{extended relation} $\overline{r}$ of r over R is accordingly $\{<t,t> \mid t \in r\}$. The two different copies of an attribute will be distinguished by the subscripts 1 and 2. We make the same kind of extensions to tables also, except that the conditions are not repeated. Thus an underline{extended table} for $T \in Tab(R)$ is a table $\overline{T} \in Tab(\overline{R})$ with $\overline{T}_C = \{<t(R),t(R),t(Con)> \mid t \in T_C\}$ and $\overline{T}_G = T_G$. An example will clarify the point.

Proceedings of the Tenth International
Conference on Very Large Data Bases.

42

Example 3.   $T_C(A\ B\ C\ Con)$     $T_G = \{\underline{true}, \neg(y=c)\}$

$\quad$
| a | x | c | $\neg(x=b)$ |
| a | b | y | true |

$\overline{T}_C(A_1 B_1 C_1 A_2 B_2 C_2 Con)$    $\overline{T}_G = \{\underline{true}, \neg(y=c)\}$

$\quad$
| a | x | c | a | x | c | $\neg(x=b)$ |
| a | b | y | a | b | y | true |

An underline{algebraic dependency} [YP] is a statement $f(\overline{R}) \subseteq f'(\overline{R})$, where f and f' are PJ-expressions. We will however only need the counterparts of FDs and JDs. The FD X->A for relations over R corresponds to the algebraic dependency

$$\pi_{A_1 A_2}(\pi_{X_1 A_1}(\overline{R}) * \pi_{X_1 A_2}(\overline{R})) \subseteq \pi_{A_1 A_2}(\overline{R}).$$

The JD $*[R_1,\ldots,R_p]$, $\overset{p}{\underset{i=1}{\cup}} R_i = R$, corresponds to to the algebraic dependency

$$\pi_{R_1}(R) * \ldots * \pi_{R_p}(R) \subseteq R.$$

Here correspondence means that the sets $\Sigma$ of FDs and JDs and $\Sigma'$ of the corresponding algebraic dependencies express the same constraints, i.e. that $Sat(\Sigma) = Sat(\Sigma')$. The algebraic counterpart of a JD should not raise anyones eyebrows. For the counterpart of a FD we give the following example of a relation that does underline{not} satisfy the dependency A->B.

Example 4.   $r(A\ B)$     $\overline{r}(A_1 B_1 A_2 B_2)$

| a | b |
| a | b' |

| a | b | a | b |
| a | b' | a | b' |

$\pi_{B_1 B_2}(\overline{r})(B_1 B_2)$

| b | b |
| b' | b' |

$\pi_{B_1 B_2}(\pi_{A_1 B_1}(\overline{r}) * \pi_{A_1 B_2}(\overline{r}))(B_1 B_2)$

| b | b |
| b | b' |
| b' | b' |
| b' | b |

The chase procedure can now be described with the following algorithm.

Algorithm Chase

Input: A Table $T \in Tab(R)$ and a set $\Sigma$ of FDs and JDs for R

Output: A table T' such that
$$Rep(T') = Comp_\Sigma(Rep(T)) \cap Sat(\Sigma)$$

Method: T' <- T

repeat foreach $*[R_1,\ldots,R_p] \in \Sigma$ do
$$T' \leftarrow \hat{\pi}_{R_1}(T') \overset{\Diamond}{\ast} \ldots \overset{\Diamond}{\ast} \hat{\pi}_{R_p}(T')$$
until the number of c-tuples in $T'_C$
no longer increases

foreach $X \rightarrow A \in \Sigma$ do begin
$$T'' \leftarrow \hat{\pi}_{A_1 A_2}(\hat{\pi}_{X_1 A_1}(T') \overset{\Diamond}{\ast} \hat{\pi}_{X_1 A_2}(T'))$$
foreach $t \in T''_C$ do
if $t(A_1) \neq t(A_2)$ then
  if $t(A_1)$ and $t(A_2)$ both
  are constants
  then $T'_G \leftarrow T'_G \cup \{\neg t(Con)\}$
  else $T'_G \leftarrow T'_G \cup \{t(Con) \Rightarrow$
  $t(A_1) = t(A_2)\}$
end

end of algorithm chase.

Some examples of the chase can be found in the end of this chapter. We now proceed to prove the correctness and some properties of the chase.

Theorem 2. For a given set of FDs and JDs and a table T, the chase algorithm only requires a finite number of steps and the resulting table T' is finite.

Proof. The algorithm does not introduce new variables or constants, so the resulting table T' is finite. Furthermore, each time we perform the repeat-until loop the number of c-tuples in $T'_C$ increases with at least one (except of course during the last loop). The loop is thus performed only a finite number of times. Since each FD is applied only once, the whole algorithm terminates after a finite number of steps. □

Different choices of the order for applying the JDs will result in different tables. In the full version of this paper we prove that these tables are m-equal, so the order is of no significance. Usually we will let T' denote the result of chasing T with some arbitrary order.

The next lemma gives the relationship between T and T'.

Lemma 2. Let T' be the result of chasing a table T with a set of FDs and JDs. Then $T \subseteq T'$.

Proof. The repeat-until loop is a repetitive application of PJ-expressions fullfilling the conditions for Lemma 1. Thus the result holds at this stage. The FDs only cause adding of conditions to $T_G$, and hence $T_G \subseteq T'_G$ and consequently $T \subseteq T'$. □

The main theorem states that the resulting table T' has the desired property. The proof of the theorem is given in the full version of this paper.

Theorem 3. Let T' be the result of performing the chase on a table $T \in Tab(R)$ with a set $\Sigma$ of FDs and JDs. Then $Rep(T') = Comp_\Sigma(Rep(T)) \cap Sat(\Sigma)$. □

In the chase of [MMS] any cover for the set of dependencies can be used. The results will be identical. In our chase we will get equivalent tables, which is sufficient for our purposes.

Corollary. If T' is the result of performing the chase with a set $\Sigma$ of FDs and JDs on a table T, and T'' is the result when using a set $\Sigma'$, where $Sat(\Sigma) = Sat(\Sigma')$, then $T' \equiv T''$.

Proof. By Theorem 3 and the definition of completions $Rep(T') = Comp_\Sigma(Rep(T)) \cap Sat(\Sigma) = Comp_{\Sigma'}(Rep(T)) \cap Sat(\Sigma') = Rep(T'')$. □

Before closing this section we will give two small examples of the chase. These examples show the kind of information that has not been deducable by previously known chase methods (i.e. [IL3, MMS, Vas2]). T' will as usual denote the result of chasing T with $\Sigma$.

Proceedings of the Tenth International
Conference on Very Large Data Bases.

43

Singapore, August, 1984

**Example 5.** $\Sigma = \{A\text{->}B\}$

$T_C$(A  B  Con)        $T_G = \{\underline{true}\}$

   a  b  <u>true</u>
   x  b' <u>true</u>

$T'_C$(A  B  Con)        $T'_G = \{\underline{true}, \neg(x\text{=}a)\}$

   a  b  <u>true</u>
   x  b' <u>true</u>

**Example 6.** $\Sigma = \{*[AB,AC]\}$

$T_C$(A  B  C  Con)        $T_G = \{\underline{true}\}$

   a  b  c  <u>true</u>
   x  b' c' <u>true</u>

$T'_C$(A  B  C  Con)        $T'_G = \{\underline{true}\}$

   a  b  c  <u>true</u>
   x  b' c' <u>true</u>
   x  b  c' <u>x=a</u>
   x  b' c  x=a

In Example 5 we are able to express the fact
that x cannot equal a. Example 6 is perhaps
more interesting. Any relation in $\text{Comp}_\Sigma(\text{Rep}(T))$
$\cap$ Sat$(\Sigma)$ should include two or four tuples, de-
pending on the value of x. The table T' ex-
presses exactly this and the content and condi-
tions of the extra tuples. The conclusion is
that the expressive power of our tables is re-
quired for properly handling data dependencies.
These two examples can in fact be used to show
that the usual device for representing null-
values, in [IL2, IL3] called v-tables (ordinary
relations with variables), is not capable of
fully supporting data dependencies. That is,
there are v-tables T and dependencies $\Sigma$ such
that there exists no v-table T' for which
Rep(T') = $\text{Comp}_\Sigma(\text{Rep}(T)) \cap \text{Sat}(\Sigma)$. The problem is
also noted by [IL3].

## 5. CONCLUSIONS

We have presented a relational system for
handling null-values of type "value exists, but
is unknown". The key idea is that a relation with
null-values, here modelled by a so called table,
represents a set of relations, one of which cor-
responds to the incompletely known state of the
real world. The formal criteria for the correct-
ness of an algebra that operates on tables is
that only information that holds in every rela-
tion in the represented set is inferred. We have
correctly extended the complete relational alge-
bra to operate on tables.

The main contribution of this paper however
lies in the capability of the system to support
functional and full join dependencies. The formal
criteria for a table to "satisfy" a set of depen-
dencies is that the dependencies are satisfied in
every relation in the set that the table repre-
sents. This strong form of satisfaction is re-
quired for maintaining the correctness of the
algebra in the precence of dependencies. We have
then presented a transformation algorithm, called
a chase, that enforces a given set of dependen-
cies on a table in such a way that the dependen-
cies are satisfied in every relation in the set
represented by the transformed table. This chase
algorithm can also be regarded as a generaliza-
tion of previously known chase methods, and we
have given some examples of the kind of infor-
mation that only our novel chase is able to de-
duce. Our results also show that the expressive
power of our tables is required for fully sup-
porting data dependencies in databases with in-
complete information.

# REFERENCES

Abi    S. Abiteboul, Algebraic analogues to fundamental notions of query and dependency theory. Rapports de Recherche 201. Institut National de Recherche en Informatique et en Automatique. Roquencourt, April 1983.

Bis1    J. Biskup, A formal approach to null values in database relations. In: Advances in Data Base Theory - Vol. 1 J. Minker & J. M. Nicolas (eds.) Plenum Press, New York London 1981, 299-341.

Bis2    J. Biskup, A foundation of Codd's relational maybe operations. ACM Transactions on Database Systems 8,4 (Dec. 1983), 608-636.

FUV    R. Fagin, J. D. Ullman & M. Y. Vardi, On the semantics of updates in databases. Proc. of the Second ACM SIGACT-SIGMOD Symposium on Principles of Database Systems 1983, 352-365.

Gra1    G. Grahne, Information incompleteness in databases. In: Proc. of the Winter School on Theoretical Computer Science. R. Back et al (eds.). Finnish Society of Information Processing Science, Lammi, January 1984, 65-98.

Gra2    G. Grahne, Updates in databases with incomplete information. In preparation.

IL1    T. Imielinski & W. Lipski, On representing incomplete information in a relational database. Proc. of the Sixth International Conference on Very Large Data Bases 1981, 388-397.

IL2    T. Imielinski & W. Lipski, Incomplete information in relational databases. To appear in J. ACM.

IL3    T. Imielinski & W. Lipski, Incomplete information and dependencies in relational databases. Proc. of the ACM SIGMOD Conference on Management of Data 1983, 178-184.

Lie    E. Lien, Multivalued dependencies with null values in relational databases. Proc. of the Fifth International Conference on Very Large Data Bases 1979, 61-66.

Lip    W. Lipski, Logical problems related to incomplete information in databases. Rapport de Recherche 138. Laboratoire de Recherche en Informatique, Universite de Paris-Sud, September 1983.

Mai    D. Maier, The Theory of Relational Databases. Pitman, London 1983.

MMS    D. Maier, A. O. Mendelson & Y. Sagiv, Testing implications of data dependencies. ACM Transactions on Database Systems 4,4 (Dec. 1979), 455-469.

Rei    R. Reiter, On closed world databases. In: Logic and Databases, H. Gallaire & J. Minker (eds.). Plenum Press, New York 1978, 55-76.

Sci    E. Sciore. The universal instance and database design. Ph.D. Thesis. Princeton University 1980.

Ull    J. D. Ullman, Principles of Database Systems (Second Edition). Computer Science Press, Potomac, Md. 1982.

Vas1    Y. Vassiliou, Null values in database management: a denotational semantics approach. Proc. ACM SIGMOD 1979 Conference on Management of Data, 162-169.

Vas2    Y. Vassiliou, Functional dependencies and incomplete information. Proc. of the Sixth International Conference on Very Large Data Bases 1981, 260-269.

YP    M. Yannakakis & C. H. Papadimitriou, Algebraic dependencies. Journal of Comp. and Syst. Sci. 25,1 (1982), 2-41.