# Handling Inconsistencies in Tables with Nulls and Functional Dependencies

**Dominique Laurent · Nicolas Spyratos**

arXiv:2108.02581v2 [cs.DB] 27 Nov 2021

**Abstract** In this paper we address the problem of handling inconsistencies in tables with missing values (also called nulls) and functional dependencies. Although the traditional view is that table instances must respect all functional dependencies imposed on them, it is nevertheless relevant to develop theories about how to handle instances that violate some dependencies. Regarding missing values, we make no assumptions on their existence: a missing value exists only if it is inferred from the functional dependencies of the table.

We propose a formal framework in which each tuple of a table is associated with a truth value among the following: true, false, inconsistent or unknown; and we show that our framework can be used to study important problems such as consistent query answering, table merging, and data quality measures - to mention just a few. In this paper, however, we focus mainly on consistent query answering, a problem that has received considerable attention during the last decades.

The main contributions of the paper are the following: (a) we introduce a new approach to handle inconsistencies in a table with nulls and functional dependencies, (b) we give algorithms for computing all true, inconsistent and false tuples, (c) we investigate the relationship between our approach and Four-valued logic in the context of data merging, and (d) we give a novel solution to the consistent query answering problem and compare our solution to that of table repairs.

**Keywords** Inconsistent database . Functional dependency . Null value . Data merging . Consistent query answering

Dominique Laurent
ETIS Laboratory - ENSEA, CY Cergy Paris University, CNRS
F-95000 Cergy-Pontoise, France
`dominique.laurent@u-cergy.fr`

Nicolas Spyratos
LISN Laboratory - University Paris-Saclay, CNRS
F-91405 Orsay, France
`nicolas.spyratos@lri.fr`

## 1 Introduction

In several applications today we encounter tables with missing values and functional dependencies. Such a table is often the result of merging two or more other tables coming from different sources. Typical examples include recording the results of collaborative work, merging of tables during data staging in data warehouses or checking the consistency of a relational database.

As an example of collaborative work consider two groups of researchers each studying three objects found in an archaeological site. The researchers of each group record in a table data regarding the following attributes of each object:

- Identifier (here of the form $i_n$ where $n$ is an integer, distinct objects being associated with distinct identifiers)
- Kind (such as statue, weapon, ...)
- Material from which the object is made (such as iron, bronze, marble, ...)
- Century in which the object is believed to have been made.

At the end of their work each group submits their findings to the site coordinator in the form of a table as shown in Figure 1 (tables $D_1$ and $D_2$). Each row of a table contains data recorded for a single object. For example, the row $(i_1, statue, marble, 1.BC)$ means that object $i_1$ is a statue made of marble and believed to have been made in the first century before Christ. Similarly the row $(i_2, statue, , 2.BC)$ means that object $i_2$ is a statue of unknown material, believed to have been made during the second century before Christ. Note that, in this tuple, there is a missing value, meaning that the material from which object $i_2$ is made could not be determined.

| $D_1$ | $Id$ | $K$ | $M$ | $C$ |
|---|---|---|---|---|
| | $i_1$ | $k$ | $m$ | $c$ |
| | $i_1$ | | $m'$ | |
| | $i_2$ | $k'$ | $m'$ | $c$ |
| | $i_2$ | $k'$ | $m''$ | |
| | $i_3$ | | $m$ | |

| $D_2$ | $Id$ | $K$ | $M$ | $C$ |
|---|---|---|---|---|
| | $i_1$ | $k$ | | $c$ |
| | $i_2$ | $k'$ | | $c'$ |
| | $i_2$ | $k'$ | $m''$ | |
| | $i_3$ | $k'$ | | |

| $D$ | $Id$ | $K$ | $M$ | $C$ |
|---|---|---|---|---|
| | $i_1$ | $k$ | $m$ | $c$ |
| | $i_1$ | | $m'$ | |
| | $i_1$ | $k$ | | $c$ |
| | $i_2$ | $k'$ | $m'$ | $c$ |
| | $i_2$ | $k'$ | $m''$ | |
| | $i_2$ | $k'$ | | $c'$ |
| | $i_3$ | | $m$ | |
| | $i_3$ | $k'$ | | |

**Fig. 1** The tables prepared by the two groups and the merged table

Now, the data contained in the two tables can be merged into a single table $D$ containing all tuples from the two tables, without duplicates, as shown in Figure 1 (table $D$). In doing this merging, we may have discrepancies between tuples of $D$. For example, object $i_1$ appears in $D$ as being made from two different materials; and object $i_2$ appears as made from two different materials and in two different

centuries. This kind of discrepancies may lead to 'inconsistencies' that should be identified by the site coordinator and resolved in cooperation with the researchers of the two groups.

It should be obvious from this example that the merging of two or more tables into a single table more often than not results in inconsistencies even if the individual tables are each consistent. For example, although each of the tables $D_1$ and $D_2$ shown in Figure 1 satisfies the functional dependencies $Id \rightarrow K$ and $Id \rightarrow C$, the merged table $D$ does not satisfy $Id \rightarrow C$.

A similar situation arises in data warehouses where one tries to merge views of the underlying sources into a single materialized view to be stored in the data warehouse.

As a last example, in a relational database, although each table may satisfy its functional dependencies, the database as a whole may violate some dependencies. To determine whether the database is consistent with its dependencies, one proceeds as follows: all tables are merged by placing their tuples into a single universal table $D$ possibly with missing values (under certain assumptions discussed in [25]); then all functional dependencies are applied on $D$ through the well known chase algorithm [12,24]. If the algorithm terminates successfully (*i.e.,* no inconsistency is detected) then the database is consistent; otherwise the algorithm stops when a first inconsistency is detected and the database is declared inconsistent.

So in general the question is: what should we do when a table is inconsistent? There are roughly three approaches: (*a*) reject the table, (*b*) try to correct or 'repair' it so that to make it consistent (and therefore be able to work with the repaired table) and (*c*) keep the table as is but make sure you know which part is consistent and which is not.

The first approach is followed by database theorists when checking database consistency, as explained above. This approach is clearly not acceptable in practice as the universal table might contain a consistent set of tuples that can be useful to users (*e.g.,* users can still query the consistent part of the table).

The second approach tries to alleviate the impact of inconsistent data on the answers to a query by introducing the notion of repair: a repair is a minimally different consistent instance of the table and an answer is consistent if it is present in every repair. This approach, referred to as 'consistent query answering', has motivated important research efforts during the past two decades and is still the subject of current research. The reader is referred to Section 6 for a brief overview of the related literature. However, this approach is always difficult to implement due to important issues related to computational complexity and/or to semantics (there is still no consensus regarding the definition of 'consistent answer').

In our work we follow the third approach that is, we keep inconsistencies in the table but we determine which part of the table is consistent and which is not. More specifically, we use set theoretic semantics for tuples and functional dependencies that allow us to associate each tuple of the table with one truth value among the following: true, false, inconsistent or unknown. By doing so we can study a number of important problems including in particular the problem of consistent query answering, and the definition of data quality measures.

Regarding consistent query answering, our model offers a fundamentally different and direct solution to the problem: the consistent answer is obtained by simply retrieving true tuples that fit the query requirements.

Moreover our approach offers the possibility of defining meaningful data quality measures. For example if a table contains a hundred tuples of which only five are true while the remaining ones are inconsistent, then the quality of data contained in the table is five percent. Since we have polynomial algorithms for computing all true, false and inconsistent tuples, we can define several quality measures of the data contained in a table, inspired by the work in [17]. We can then use such measures to accompany query answers so that users are informed of the quality of the answer they receive (*e.g.,* getting an answer from a table with ninety five per cent of true tuples is more reliable than if the table contained only five per cent of true tuples). However, defining and studying such measures lies outside the goals of the present paper. In this paper we focus on one important application of our approach, namely consistent query answering. A complete account of data quality measures will be reported in a future paper.

The main contributions of the present paper can be summarized as follows:

1. We introduce a new approach to handle inconsistencies in a table with nulls and functional dependencies; we do so by adapting the set theoretic semantics of [21] to our context and by extending the chase algorithm so that all inconsistencies are accounted for in the table.
2. We give polynomial algorithms in the size of the table for computing all true and all inconsistent tuples in the table.
3. We investigate the relationship of our approach with Four-valued logic in the context of data merging.
4. We propose a novel approach for consistent query answering and we investigate how our approach relates to existing approaches.

The paper is organized as follows: In Section 2 we recall basic definitions and notations regarding tables and we introduce the set theoretic semantics that we use in our work. In Section 3 we give definitions and properties regarding the truth values that we associate with tuples. In Section 4 we study computational issues and give algorithms for computing the truth values of tuples. In Section 5, we show how our approach relates to Four-value logic when merging two or more tables. In Section 6 we present a novel solution to the problem of consistent query answering and compare it to existing approaches. Section 7 contains concluding remarks and suggestions for further research.

## 2 The Model

In this section we present the basic definitions regarding tuples and tables as well as the set theoretic semantics that we use for tuples and functional dependencies. Our approach builds upon earlier work on the partition model [21].

### 2.1 The Partition Model Revisited

Following [21], we consider a universe $U = \{A_1, \ldots, A_n\}$ in which every attribute $A_i$ is associated with a set of atomic values called the domain of $A_i$ and denoted by $dom(A_i)$. An element of $\bigcup_{A \in U} dom(A)$ is called a *domain constant* or a *constant*. We call *relation schema* (or simply *schema*) any nonempty subset of $U$ and we

denote it by the concatenation of its elements; for example $\{A_1, A_2\}$ is simply denoted by $A_1 A_2$. Similarly, the union of schemas $S_1$ and $S_2$ is denoted as $S_1 S_2$ instead of $S_1 \cup S_2$.

We define a *tuple* $t$ to be a partial function from $U$ to $\bigcup_{A \in U} dom(A)$ such that, for every $A$ in $U$, if $t$ is defined over $A$ then $t(A)$ belongs to $dom(A)$. The domain of definition of $t$ is called the *schema* of $t$, denoted by $sch(t)$. We note that tuples in our approach satisfy the *First Normal Form* [24] in the sense that each tuple component is an atomic value from an attribute domain.

Regarding notation, we follow the usual convention that, whenever possible, lower-case characters denote domain constants and upper-case characters denote the corresponding attributes. Following this convention the schema of a tuple $t = ab$ is $AB$ and more generally, we denote the schema of $t$ as $T$.

Assuming that the schema of a tuple $t$ is understood, $t$ is denoted by the concatenation of its values, that is: $t = a_{i_1} \ldots a_{i_k}$ means that for every $j = 1, \ldots, k$, $t(A_{i_j}) = a_{i_j}$, $a_{i_j}$ is in $dom(A_{i_j})$, and $sch(t) = A_{i_1} \ldots A_{i_k}$.

We assume that for any distinct attributes $A$ and $B$, we have either $dom(A) = dom(B)$ or $dom(A) \cap dom(B) = \emptyset$. However, this may lead to ambiguity when two attributes have the same domain. Ambiguity can be avoided by prefixing each value of an attribute domain with the attribute name. For example, if $dom(A) = dom(B)$ we can say 'an $A$-value $a$' to mean that $a$ belongs to $dom(A)$, and 'a $B$-value $a$' to mean that $a$ belongs to $dom(B)$. In order to keep the notation simple we shall omit prefixes whenever no ambiguity is possible.

Denoting by $\mathcal{T}$ the set of all tuples that can be built up given a universe $U$ and the corresponding attribute domains, a *table* $D$ is a *finite* sub-set of $\mathcal{T}$ where duplicates are *not allowed*.

Given a tuple $t$, for every $A$ in $sch(t)$, $t(A)$ is also denoted by $t.A$ and more generally, for every subset $S$ of $sch(t)$ the restriction of $t$ to $S$, also called *sub-tuple* of $t$, is denoted by $t.S$. In other words, if $S \subseteq sch(t)$, $t.S$ is the tuple such that $sch(t.S) = S$ and for every $A$ in $S$, $(t.S).A = t.A$.

Moreover, $\sqsubseteq$ denotes the 'sub-tuple' relation, defined over $\mathcal{T}$ as follows: for any tuples $t_1$ and $t_2$, $t_1 \sqsubseteq t_2$ holds if $t_1$ is a sub-tuple of $t_2$. It is thus important to keep in mind that whenever $t_1 \sqsubseteq t_2$ holds, it is understood that $sch(t_1) \subseteq sch(t_2)$ also holds.

The relation $\sqsubseteq$ is clearly a partial order over $\mathcal{T}$. Given a table $D$, the set of all sub-tuples of the tuples in $D$ is called the *lower closure* of $D$ and it is defined by: $\mathsf{LoCl}(D) = \{q \in \mathcal{T} \mid (\exists t \in D)(q \sqsubseteq t)\}$. We shall call a table *reduced* if it contains only maximal tuples (*i.e.,* if no tuple in the set is sub-tuple of some other tuple in the set).

The notion of $\mathcal{T}$-*mapping*, as defined below, generalizes that of interpretation defined in [21].

**Definition 1** Let $U$ be a universe. A $\mathcal{T}$-*mapping* is a mapping $\mu$ defined from $\bigcup_{A \in U} dom(A)$ to $2^{\mathbb{N}}$. A $\mathcal{T}$-*mapping* $\mu$ can be extended to the set $\mathcal{T}$ as follows: for every $t = a_{i_1} \ldots a_{i_k}$ in $\mathcal{T}$, $\mu(t) = \mu(a_{i_1}) \cap \ldots \cap \mu(a_{i_k})$.

A $\mathcal{T}$-mapping $\mu$ is an *interpretation* if $\mu$ satisfies the *partition constraint* stating that for every $A$ in $U$, and for all distinct $a$ and $a'$ in $dom(A)$, $\mu(a) \cap \mu(a') = \emptyset$.

We emphasize that in [21] interpretations provide the basic tool for defining true tuples: a tuple $t$ is said to be true in an interpretation $\mu$ if $\mu(t)$ is nonempty.

5

To see the intuition behind this definition consider a relational table $D$ over $U$ and suppose that each tuple is associated with a unique identifier, say an integer. Now, for every $A$ in $U$ and every $a$ in $dom(A)$, define $\mu(a)$ to be the set of all identifiers of the tuples in $D$ containing $a$. Then $\mu$ is an interpretation as it satisfies the partition constraint. Indeed, due to the fact that, for every attribute $A$ in $U$, a tuple $t$ can not have more than one $A$-value, it is then impossible that $\mu(a) \cap \mu(a')$ be nonempty for any distinct values $a$, $a'$ in $dom(A)$.

Incidentally, if for every $A$ in $U$ we denote by $dom^*(A)$ the set of all $A$-values such that $\mu(a) \neq \emptyset$, then the set $\{\mu(a) \mid a \in dom^*(A)\}$ is a *partition* of $\bigcup_{a \in dom^*(A)} \mu(a)$ (whence the name "partition model"). The following example illustrates this important feature.

*Example 1* Considering $U = \{A, B, C\}$ and $D = \{ab, bc, ac, a'b', b'c', abc\}$, the tuples in $D$ can be respectively assigned the identifiers 1, 2, 3, 4, 5 and 6. In that case, we have $\mu(a) = \{1, 3, 6\}$, $\mu(a') = \{4\}$, $\mu(b) = \{1, 2, 6\}$, $\mu(b') = \{4, 5\}$, $\mu(c) = \{2, 3, 6\}$, $\mu(c') = \{5\}$, and $\mu(\alpha) = \emptyset$ for any constant $\alpha$ different than $a$, $a'$, $b$, $b'$, $c$ and $c'$.

It is clear that the $\mathcal{T}$-mapping $\mu$ is an interpretation and, since $dom^*(A)$, $dom^*(B)$ and $dom^*(C)$ are respectively equal to $\{a, a'\}$, $\{b, b'\}$ and $\{c, c'\}$, it is easy to see that $\{\mu(\alpha) \mid \alpha \in dom^*(A)\}$ is a partition of $\{1, 3, 4, 6\}$, $\{\mu(\beta) \mid \beta \in dom^*(B)\}$ is a partition of $\{1, 2, 4, 5, 6\}$, and $\{\mu(\gamma) \mid \gamma \in dom^*(C)\}$ is a partition of $\{2, 3, 5, 6\}$.

Moreover, extending $\mu$ to non unary tuples yields the following regarding the tuples in $D$: $\mu(ab) = \{1, 6\}$, $\mu(bc) = \{2, 6\}$, $\mu(ac) = \{3, 6\}$, $\mu(a'b') = \{4\}$, $\mu(b'c') = \{5\}$, and $\mu(abc) = \{6\}$. $\qquad\square$

Summarizing our discussion, when dealing with consistent tables in [21], only interpretations are relevant. In the present work, we follow the same idea, but we also extend the work of [21] so that we can deal with inconsistencies. As we shall see, non satisfaction of the partition constraint in Definition 1 is the key criterion to characterize inconsistent tuples.

## 2.2 Functional Dependencies

The notion of functional dependency in our approach is defined as in [21].

**Definition 2** Let $U$ be a universe. A *functional dependency* is an expression of the form $X \rightarrow Y$ where $X$ and $Y$ are nonempty sub-sets of $U$.

A $\mathcal{T}$-mapping $\mu$ *satisfies* $X \rightarrow Y$, denoted by $\mu \models X \rightarrow Y$, if for all tuples $x$ and $y$, respectively over $X$ and $Y$, the following holds: if $\mu(x) \cap \mu(y) \neq \emptyset$ then $\mu(x) \subseteq \mu(y)$.

Based on Definition 2, for all $X$ and $Y$ such that $X \cap Y = \emptyset$, and for every $\mathcal{T}$-mapping $\mu$, the following holds:

$$\mu \models X \rightarrow Y \text{ if and only if } \mu \models X \rightarrow A \text{ for every } A \text{ in } Y.$$

This is so because, for every $x$ and $y$ such that $\mu(x) \cap \mu(y) \neq \emptyset$, $\mu(x) \subseteq \mu(y)$ holds if and only if $\mu(x) \subseteq \mu(a)$ holds for every constant $a$ in $y$.

Therefore without loss of generality we can assume that all functional dependencies are of the form $X \rightarrow A$ where $A$ is an attribute not in $X$. Under this assumption, we consider pairs $\Delta = (D, \mathcal{FD})$ where $D$ is a table over $U$ and $\mathcal{FD}$

a set of functional dependencies over $U$, and we say that a $\mathcal{T}$-mapping $\mu$ satisfies $\Delta$, denoted by $\mu \models \Delta$, if $(i)$ for every $t$ in $D$, $\mu(t) \neq \emptyset$, and $(ii)$ $\mu$ satisfies every $X \to A$ in $\mathcal{FD}$.

To see how our notion of functional dependency relates to the standard one in relational databases [24], recall first that a relation $r$ over universe $U$ satisfies $X \to A$ if for all tuples $t$ and $t'$ in $r$ such that $t.X = t'.X$, we have $t.A = t'.A$.

In our approach, let $\Delta = (D, \mathcal{FD})$ and consider two tuples $t$ and $t'$ in $D$ such that $XA$ is a subset of $sch(t)$ and of $sch(t')$ and let $t.X = t'.X = x$. Then for every $\mathcal{T}$-mapping $\mu$ such that $\mu \models \Delta$, $\mu(t)$ and $\mu(t')$ are nonempty, implying that $\mu(x) \cap \mu(t.A)$ and $\mu(x) \cap \mu(t'.A)$ are also nonempty. By Definition 2, this implies that $\mu(x)$ is a sub-set of $\mu(t.A)$ and of $\mu(t'.A)$. As a consequence, assuming that $t.A \neq t'.A$ (i.e., that $X \to A$ is not satisfied in the sense of the relational model), means that $\mu(t.A) \cap \mu(t'.A)$ is nonempty, and therefore $\mu$ *can not be an interpretation*.

Therefore if we restrict $\mathcal{T}$-mappings to be interpretations then the notion of functional dependency satisfaction in our approach is *the same* as that of relational databases. As we shall see, this observation supports the notion of consistency for $\Delta$, to be given later (in Definition 4).

Given $\Delta = (D, \mathcal{FD})$ and tuples $t$, $t'$, $t''$, the following notations are extensively used in the remainder of the paper.

− $\Delta \vdash t$, denotes that if $\mu \models \Delta$ then $\mu(t) \neq \emptyset$.
− $\Delta \vdash (t \sqcap t')$, denotes that if $\mu \models \Delta$ then $\mu(t) \cap \mu(t') \neq \emptyset$.
− $\Delta \vdash (t \preceq t')$ denotes that if $\mu \models \Delta$ then $\mu(t) \subseteq \mu(t')$.
− $\Delta \vdash (t \preceq t' \sqcap t'')$ denotes that if $\mu \models \Delta$ then $\mu(t) \subseteq \mu(t') \cap \mu(t'')$.

Given $\Delta = (D, \mathcal{FD})$, we now build a particular $\mathcal{T}$-mapping $\mu$ such that $\mu \models \Delta$ as follows: Let $(\mu_i)_{i \geq 0}$ be the sequence defined by the steps below:

1. Associate each tuple $t$ with an identifier, $id(t)$, called the *tuple identifier* of $t$ (this can be an integer that identifies $t$ uniquely).
2. Let $\mu_0$ be the mapping defined for every domain constant $a$ by:
   $\mu_0(a) = \{id(t) \mid t \in D \text{ and } a \sqsubseteq t\}$.
3. While there exists $X \to A$ in $\mathcal{FD}$, $x$ over $X$ and $a$ in $dom(A)$ such that $\mu_i(xa) \neq \emptyset$ and $\mu_i(x) \not\subseteq \mu_i(a)$, define $\mu_{i+1}$ by: $\mu_{i+1}(a) = \mu_i(a) \cup \mu_i(x)$ and $\mu_{i+1}(\alpha) = \mu_i(\alpha)$ for any other constant $\alpha$.

**Lemma 1** *For every $\Delta = (D, \mathcal{FD})$, the sequence $(\mu_i)_{i \geq 0}$ has a unique limit $\mu^*$ such that $\mu^* \models \Delta$. Moreover:*

1. *For all $a_1$ and $a_2$ in the same attribute domain $dom(A)$, if $\mu^*(a_1) \cap \mu^*(a_2) \neq \emptyset$ then there exist $X \to A$ in $\mathcal{FD}$ and $x$ over $X$ such that $\mu^*(x) \neq \emptyset$ and $\mu^*(x) \subseteq \mu^*(a_1) \cap \mu^*(a_2)$.*
2. *For all $\alpha$ and $\beta$, $\Delta \vdash (\alpha \sqcap \beta)$ holds if and only if $\mu^*(\alpha) \cap \mu^*(\beta) \neq \emptyset$ holds.*

*Proof* See Appendix A. □

Given $\Delta = (D, \mathcal{FD})$, Lemma 1 shows the following:

1. There always exists a $\mathcal{T}$-mapping $\mu$ such that $\mu \models \Delta$.
2. When two constants from the same domain have common identifiers with respect to $\mu^*$ then this is due to a functional dependency.

3. For every tuple $t$, $\Delta \vdash t$ if and only if $\mu^*(t) \neq \emptyset$.

It is important to note that the $\mathcal{T}$-mapping $\mu^*$ as defined in Lemma 1 is not necessarily an interpretation as the following example shows.

*Example 2* Let $U = \{A, B, C\}$ and $\Delta = (D, \mathcal{FD})$ where $D = \{ab, bc, abc'\}$ and $\mathcal{FD} = \{B \to C\}$. Associating $ab$, $bc$ and $abc'$ respectively with 1, 2 and 3, $\mu^*$ is obtained as follows:
• First, we have $\mu_0(a) = \{1, 3\}$, $\mu_0(b) = \{1, 2, 3\}$, $\mu_0(c) = \{2\}$ and $\mu_0(c') = \{3\}$ and $\mu_0(\alpha) = \emptyset$ for any other domain constant $\alpha$.
• Then, considering $B \to C$, we have $\mu_1(a) = \{1, 3\}$, $\mu_1(b) = \mu_1(c) = \mu_1(c') = \{1, 2, 3\}$ and $\mu_1(\alpha) = \emptyset$ for any other domain constant $\alpha$.

Hence, $\mu^* = \mu_1$ and we remark that $\mu^*(c) \cap \mu^*(c') \neq \emptyset$, thus that $\mu^*$ is not an interpretation. Nevertheless, as stated by Lemma 1, it is easy to see that $\mu^* \models \Delta$.
$\square$

We note here that the authors of [22] use a construction similar to that of Lemma 1 to define a minimal model of $\Delta$, called 'query model', assuming that $D$ is consistent with $\mathcal{FD}$.

Now, in order to characterize when $\Delta \vdash (t \preceq a)$ holds, we introduce the notion of *closure of a tuple $t$ in $\Delta$* inspired by the well known relational notion of closure of a relation scheme with respect to a set of functional dependencies [24].

**Definition 3** Given a database $\Delta = (D, \mathcal{FD})$ and a tuple $t$, the *closure of $t$ in $\Delta$* (or *closure of $t$* for short, when $\Delta$ is understood), denoted by $t^+$, is the set of all domain constants $a$ such that $\Delta \vdash (t \preceq a)$ holds.

We notice that, based on Definition 3, for every constant $a$ occurring in a tuple $t$ (*i.e.*, if $a \sqsubseteq t$ holds) then $a$ is in $t^+$, because, in this case, $\mu(t) \subseteq \mu(a)$ holds for every $\mathcal{T}$-mapping $\mu$. However constants not occurring in $t$ may also appear in $t^+$ due to functional dependencies, as shown in the following example.

*Example 3* Continuing Example 2 where $U = \{A, B, C\}$ and $\Delta = (D, \mathcal{FD})$ with $D = \{ab, bc, abc'\}$ and $\mathcal{FD} = \{B \to C\}$, we show that $c$ belongs to $(ab)^+$.

Indeed, for every $\mu$ such that $\mu \models \Delta$, we have $\mu(ab) \subseteq \mu(b)$ (since $b \sqsubseteq ab$) and $\mu(b) \subseteq \mu(c)$ (due to $B \to C$ and the fact that $\mu(b) \cap \mu(c) \neq \emptyset$ must hold). Hence, by transitivity, $\mu(ab) \subseteq \mu(c)$ holds, implying that $\Delta \vdash (ab \preceq c)$ holds, which by Definition 3, means that $c$ belongs to $(ab)^+$. It should also be noticed that a similar argument shows that $c'$ also belongs to $(ab)^+$.
$\square$

Clearly computing the closure directly from its definition is inefficient. Algorithm 1 gives a method for computing the closure, since the following lemma states that this algorithm correctly computes the closure.

**Lemma 2** *Let $\Delta = (D, \mathcal{FD})$ and $t$ a tuple. Then Algorithm 1 computes correctly the closure $t^+$ of $t$.*

*Proof* See Appendix B.
$\square$

We draw attention on the fact that the database involved in Algorithm 1 is not $\Delta$ but the database $\Delta_t$ that can be seen as $\Delta$ in which the tuple $t$ has been added.

It should however be noticed that in case $\Delta \vdash t$, this distinction is not necessary because in this case, for every tuple $q$, $\Delta \vdash q$ holds if and only if $\Delta_t \vdash q$ holds.

**Algorithm 1** Closure of $t$

---
**Input:** $\Delta = (D, \mathcal{FD})$ and a tuple $t$.
**Output:** The closure $t^+$ of $t$
 1: $\Delta_t := (D_t, \mathcal{FD})$ where $D_t = D \cup \{t\}$
 2: $t^+ := \{a \mid a \sqsubseteq t\}$
 3: **while** $t^+$ changes **do**
 4:     **for all** $X \rightarrow A \in \mathcal{FD}$ **do**
 5:         **for all** $x$ such that for every $b$ in $x$, $b \in t^+$ and $\Delta_t \vdash xa$ **do**
 6:             $t^+ := t^+ \cup \{a\}$
 7: **return** $t^+$

---

This is a consequence of the fact that, as seen in Appendix B, if $\Delta \vdash t$ then for every $\mathcal{T}$-mapping $\mu$, $\mu \models \Delta$ holds if and only if $\mu \models \Delta_t$.

On the other hand, the following example shows that when $\Delta \nvdash t$, the introduction of $\Delta_t$ instead of $\Delta$ is necessary for correctly computing $t^+$.

*Example 4* Let $U = \{A, B, C\}$ and $\Delta = (D, \mathcal{FD})$ where $D = \{ac, b\}$ and $\mathcal{FD} = \{A \rightarrow B, B \rightarrow C\}$.

It is easy to see that when numbering the tuples in $D$ by 1 for $ac$ and 2 for $b$, the $\mathcal{T}$-mapping $\mu^*$ for $\Delta$ is defined by: $\mu^*(a) = \mu^*(c) = \{1\}$, $\mu^*(b) = \{2\}$ and $\mu^*(\alpha) = \emptyset$ for any other constant $\alpha$.

For $t = ab$, we argue that $c$ is in $t^+$, that is, for every $\mu$ such that $\mu \models \Delta$, $\mu(ab) \subseteq \mu(c)$ holds. Indeed, this trivially holds if $\mu(ab) = \emptyset$ (as is the case with $\mu^*$), and otherwise the following proof can be done:
- As $\mu(ab) \neq \emptyset$, $A \rightarrow B$ implies that $\mu(a) \subseteq \mu(b)$.
- As $\mu(ac) \neq \emptyset$, $\mu(a) \subseteq \mu(b)$ implies $\mu(bc) \neq \emptyset$. Thus, $\mu(b) \subseteq \mu(c)$, due to $B \rightarrow C$.
- Therefore, $\mu(a) \subseteq \mu(c)$, and since $\mu(ab) \subseteq \mu(a)$, we have $\mu(ab) \subseteq \mu(c)$.

On the other hand, computing $(ab)^+$ using a modified version of Algorithm 1 where $\Delta_t$ is replaced by $\Delta$ would output $a$ and $b$ in the closure. It should also be noticed that computing $(ab)^+$ using Algorithm 1 is as follows: by the statement on line 2, $a$ and $b$ are inserted into the closure, and then, since for $t = ab$, $\Delta_t \vdash ab$ the above reasoning shows that $\Delta_t \vdash bc$ as well. Therefore, $c$ is inserted into the closure because the test line 5 succeeds. $\square$

The following example shows a case where the tuple $t$ of which the closure is computed is such that $\Delta \vdash t$.

*Example 5* As seen in Example 2, if $U = \{A, B, C\}$ and $\Delta = (D, \mathcal{FD})$ where $D = \{ab, bc, abc'\}$ and $\mathcal{FD} = \{B \rightarrow C\}$, $\mu^*$ is defined by: $\mu^*(a) = \{1, 3\}$, $\mu^*(b) = \mu^*(c) = \mu^*(c') = \{1, 2, 3\}$ and $\mu^*(\alpha) = \emptyset$ for any other domain constant $\alpha$.

In this case, the computation of $(ab)^+$ according to Algorithm 1 is as follows:
- As $ab$ is in $D$, $\Delta_t = \Delta$. We thus run Algorithm 1 with $\Delta$ instead of $\Delta_t$.
- $(ab)^+$ is first set to $\{a, b\}$.
- Considering $B \rightarrow C$, since $b$ is in $(ab)^+$, and since $\Delta \vdash bc$ and $\Delta \vdash bc'$ (this holds because $\mu^*(c)$ and $\mu^*(c')$ are nonempty), $c$ and $c'$ are inserted in $(ab)^+$.

As no further step is processed, $(ab)^+ = \{a, b, c, c'\}$, as seen in Example 3. Thus $\Delta \vdash (ab \preceq c)$ and $\Delta \vdash (ab \preceq c')$ hold, implying $\Delta \vdash (ab \preceq c \sqcap c')$. $\square$

## 3 Semantics

In this section we provide basic definitions and properties regarding the truth value associated with a tuple. The following definition is borrowed from [21].

**Definition 4** $\Delta$ is said to be *consistent* if there exists an *interpretation* $\mu$ such that $\mu \models \Delta$.

Since in our approach, inconsistent tables are *not* discarded, it is crucial to be able to provide semantics to any $\Delta = (D, \mathcal{FD})$, being it consistent or not. To this end, inspired by Belnap's Four-valued logic [5], we consider *four* possible truth values for a given tuple $t$ in $\Delta$. The notations of truth values for tuples in our approach and their intuitive meaning are as follows, for a given tuple $t$:

- Truth value `true`: $t$ is true in $\Delta$.
- Truth value `false`: $t$ is false in $\Delta$. This means that we do *not* follow the Closed World Assumption (CWA), according to which any non true tuple is false [20].
- Truth value `inc` (*i.e.,* inconsistent): $t$ is true *and* false in $\Delta$. This truth value is necessary for 'safely' dealing with inconsistent tuples.
- Truth value `unkn` (*i.e.,* unknown): $t$ is not true, not false and not inconsistent in $\Delta$. This truth value is necessary for dealing with tuples not falling in one of the previous three categories.

In order to formalize the exact meaning of these truth values in our approach, we introduce the following terminology and notation for a given tuple $t$:

- If $\Delta \vdash t$ holds, $t$ is said to be *potentially true* in $\Delta$. Notice here that by Lemma 1, $t$ is potentially true if and only if $\mu^*(t) \neq \emptyset$.
- If $\Delta \vdash (t \preceq a \sqcap a')$ holds for some distinct $a$ and $a'$ in the same attribute domain, then we use the notation $\Delta \mathrel{|\!\!\sim} t$, and in this case, $t$ is said to be *potentially false* to reflect that $\mu(a) \cap \mu(a')$ must be empty for $\mu$ to be an interpretation. By Definition 3, $\Delta \mathrel{|\!\!\sim} t$ holds if and only if there exist $a$ and $a'$ in the same attribute domain such that $a$ and $a'$ are in $t^+$.

Consequently, if a tuple $t$ is such that $\Delta \vdash t$ and $\Delta \mathrel{|\!\!\sim} t$, then for $\mu$ to be an interpretation, $\mu$ must associate $t$ with a set expected to be empty and nonempty, which is of course a case of inconsistency! This explains why, in our approach, 'potentially true' and 'potentially false', should respectively be understood as 'true or inconsistent' and 'false or inconsistent'.

Based on this intuition, each tuple is assigned one of the four truth values according to the following definition.

**Definition 5** Given $\Delta = (D, \mathcal{FD})$ and a tuple $t$, the truth value of $t$ in $\Delta$, denoted by $v_\Delta(t)$, is defined as follows:

- $v_\Delta(t) = \texttt{true}$      if $\Delta \vdash t$ and $\Delta \mathrel{|\!\!\not\sim} t$; $t$ is said to be *true in $\Delta$*.
- $v_\Delta(t) = \texttt{false}$      if $\Delta \not\vdash t$ and $\Delta \mathrel{|\!\!\sim} t$; $t$ is said to be *false in $\Delta$*.
- $v_\Delta(t) = \texttt{inc}$      if $\Delta \vdash t$ and $\Delta \mathrel{|\!\!\sim} t$; $t$ is said to be *inconsistent in $\Delta$*.
- $v_\Delta(t) = \texttt{unkn}$      if $\Delta \not\vdash t$ and $\Delta \mathrel{|\!\!\not\sim} t$; $t$ is said to be *unknown in $\Delta$*.

We point out that the four truth values as defined above correspond exactly to the four truth values defined in the Four-valued logic [5]. The reader is referred to Section 5 for more details on this point. We illustrate Definition 5 through the following example.

*Example 6* As in Example 2, let $U = \{A, B, C\}$ and $\Delta = (D, \mathcal{FD})$ where $D = \{ab, bc, abc'\}$ and $\mathcal{FD} = \{B \to C\}$.

It has been seen in Example 5 that $(ab)^+ = \{a, b, c, c'\}$. Thus $\Delta \mathrel{|\!\sim} ab$ holds. Moreover, it is easy to see from Example 2 that $\mu^*(ab) \neq \emptyset$, implying that $\Delta \vdash ab$ holds as well. As a consequence, by Definition 5, $v_\Delta(ab) = \mathtt{inc}$, meaning that $ab$ is inconsistent in $\Delta$. We notice that similar arguments hold for $abc$, $abc'$, $bc$, $bc'$ and $b$, showing that these tuples are also inconsistent in $\Delta$.

Moreover, based on Definition 4, we also argue that $\Delta$ is *not* consistent, because every $\mu$ such that $\mu \models \Delta$ cannot be an interpretation. This is so because $\mu^*(c) \cap \mu^*(c') \neq \emptyset$ and Lemma 1 imply that for $\mu$ such that $\mu \models \Delta$, $\mu(c) \cap \mu(c') \neq \emptyset$.

Now, consider the tuple $bc''$ where $c''$ is a constant in $dom(C)$ distinct from $c$ and $c'$. To compute $(bc'')^+$ using Algorithm 1, the database $\Delta_t = (D_t, \mathcal{FD})$ where $D_t = \{ab, bc, abc', bc''\}$ is first defined and then, the closure is first set to $\{b, c''\}$. The subsequent computation steps rely on $B \to C$ and on that $\Delta_t \vdash bc$ and $\Delta_t \vdash bc'$ to insert $c$ and $c'$ in the closure.

It therefore follows that $(bc'')^+ = \{b, c'', c, c'\}$, thus that $\Delta \mathrel{|\!\sim} bc''$ holds. Since $\Delta \not\vdash bc''$ (because $\mu^*(bc'') = \emptyset$ and $\mu^* \models \Delta$), it follows that $v_\Delta(bc'') = \mathtt{false}$. Hence $bc''$ and all its super-tuples are false in $\Delta$.

As an example of unknown tuple in $\Delta$, let $a'$ be in $dom(A)$ such $a' \neq a$, and consider $a'c$. Since $\mu^*(a'c) = \emptyset$, $\Delta \not\vdash a'c$. On the other hand, it can be seen that $(a'c)^+ = \{a', c\}$, because $D \cup \{a'c\}$ does not allow any specific tuple derivation using $B \to C$. Hence, $\Delta \mathrel{\not|\!\sim} a'c$, which shows that $v_\Delta(a'c) = \mathtt{unkn}$. $\qquad\square$

The following example shows that computing *all* inconsistent tuples in $\Delta$ is not an easy task.

*Example 7* Let $\Delta = (D, \mathcal{FD})$ be defined over $U = \{A, B, C\}$ by $D = \{abc, ac'\}$ and $\mathcal{FD} = \{A \to B, B \to C\}$.

Here again, the tuples in $D$ along with the functional dependencies in $\mathcal{FD}$ show no explicit inconsistency. However computing $\mu^*$ yields the following:
- To define $\mu_0$, we associate the tuples $abc$ and $ac'$ with the integers 1 and 2, respectively. It follows that $\mu_0(a) = \{1, 2\}$, $\mu_0(b) = \{1\}$, $\mu_0(c) = \{1\}$, $\mu_0(c') = \{2\}$ and $\mu_0(\alpha) = \emptyset$ for any other domain constant $\alpha$.
- The next steps modify $\mu_0$ so as to satisfy $A \to B$ and $B \to C$ as follows:

1. Due to $A \to B$, $\mu_1$ is defined by: $\mu_1(a) = \{1, 2\}$, $\mu_1(b) = \{1, 2\}$, $\mu_1(c) = \{1\}$ and $\mu_1(c') = \{2\}$;
2. Due to $B \to C$, $\mu_2$ is defined by: $\mu_2(a) = \{1, 2\}$, $\mu_2(b) = \{1, 2\}$, $\mu_2(c) = \{1, 2\}$ and $\mu_2(c') = \{1, 2\}$.

As $\mu_2 \models \mathcal{FD}$, $\mu^* = \mu_2$. Moreover, we have $a^+ = \{a, b, c, c'\}$ and $b^+ = \{b, c, c'\}$ showing that, by Lemma 2, $\Delta \vdash a \preceq (c \sqcap c')$ and $\Delta \vdash (b \preceq c \sqcap c')$, thus that $a$ and $b$ are inconsistent in $\Delta$. It can then be seen that, for example, $abc$, $bc'$ and $ac$ are also inconsistent in $\Delta$.

Now, let $\Delta_1 = (D_1, \mathcal{FD})$ such that $D_1 = \{ac, ac'\}$. In this case, $\mu^*$ is defined by $\mu^*(a) = \{1, 2\}$, $\mu^*(c) = \{1\}$, $\mu^*(c') = \{2\}$ and $\mu^*(\alpha) = \emptyset$ for any other domain constant $\alpha$. Therefore, $a^+ = \{a\}$, showing that $a$ is *not* inconsistent in $\Delta_1$. As a consequence, $ac$, $ac'$ along with all their sub-tuples are true in $\Delta_1$ and all other tuples are unknown in $\Delta_1$. $\qquad\square$

The following proposition shows that our notion of inconsistent tuple complies with Definition 4.

**Proposition 1** $\Delta = (D, \mathcal{FD})$ *is consistent if and only if there exists no tuple* $t$ *such that* $v_\Delta(t) = \mathtt{inc}$.

*Proof* We first note that if there exists a tuple $t$ such that $v_\Delta(t) = \mathtt{inc}$, then $\Delta \vdash t$ and $\Delta \not\hspace{-0.3em}\sim t$. Hence there exist $a$ and $a'$ in the same attribute domain $dom(A)$ such that $\Delta \vdash (t \preceq a \sqcap a')$. Thus every $\mathcal{T}$-mapping $\mu$ such that $\mu \models \Delta$ satisfies that $\mu(t) \neq \emptyset$ and $\mu(t) \subseteq \mu(a) \cap \mu(a')$, implying that $\mu(a) \cap \mu(a') \neq \emptyset$. Hence, $\mu$ is not an interpretation, showing that, by Definition 4, $\Delta$ is not consistent.

Conversely, assuming that there is no tuple $t$ such that $v_\Delta(t) = \mathtt{inc}$, that is such that $\Delta \vdash t$ and $\Delta \not\hspace{-0.3em}\sim t$, we prove that $\mu^*$ is an interpretation of $\Delta$. Indeed, if $a_1$ and $a_2$ are two constants in the same attribute domain $A$ such that $\mu^*(a_1) \cap \mu^*(a_2) \neq \emptyset$, then by Lemma 1(1), there exist $X \to A$ in $\mathcal{FD}$ and $x$ over $X$ such that $\mu^*(x) \neq \emptyset$ and $\mu^*(x) \subseteq \mu^*(a_1) \cap \mu^*(a_2)$. Thus by Lemma 1(2), for every $\mu$ such that $\mu \models \Delta$, $\mu(x) \subseteq \mu(a_1) \cap \mu(a_2)$ and $\mu(x) \neq \emptyset$. We therefore obtain that $\Delta \not\hspace{-0.3em}\sim x$ and $\Delta \vdash x$, thus that $v_\Delta(x) = \mathtt{inc}$, which is a contradiction. Therefore $\mu^*$ is an interpretation, and the proof is complete. □

Based on Definition 5, we stress the following important remarks about potentially true and potentially false tuples in a given $\Delta = (D, \mathcal{FD})$:

- Let $t$ be a potentially true tuple. Since $\Delta \vdash t$ holds, as a consequence of Lemma 1, we have that $\mu^*(t) \neq \emptyset$. Therefore true or inconsistent tuples are those tuples that are associated with a nonempty set by *every* $\mathcal{T}$-mapping $\mu$ such that $\mu \models \Delta$. This implies that potentially true tuples in $\Delta$ are built up with constants occurring in $D$, and thus are in finite number. We provide in this paper effective algorithms for computing the sets of true tuples and inconsistent tuples.
- As potentially false tuples $t$ are such that $\Delta \not\hspace{-0.3em}\sim t$, they may not satisfy $\Delta \vdash t$. Hence, Lemma 1 cannot be used to characterize them. Moreover, if $\Delta \not\hspace{-0.3em}\sim t$, then every tuple $t'$ such that $t \sqsubseteq t'$ also satisfies $\Delta \not\hspace{-0.3em}\sim t'$. This is so because in this case, if $\Delta \vdash (t \preceq a \sqcap a')$, then $\Delta \vdash (t' \preceq a \sqcap a')$ holds as well. Thus, the number of potentially false tuples may be infinite in case some of the attribute domains are infinite.
- Moreover, since every false tuple $t$ is potentially false and does not satisfy $\Delta \vdash t$, it also follows as above that every tuple $t'$ such that $t \sqsubseteq t'$ is also false. Thus, the number of false tuples may be infinite in case some of the attribute domains are infinite. However, the following proposition allows to characterize when a given tuple $t$ is false.

**Proposition 2** *Given* $\Delta = (D, \mathcal{FD})$ *and a tuple* $t$, $v_\Delta(t) = \mathtt{false}$ *if and only if* $\Delta \not\vdash t$ *and* $v_{\Delta_t}(t) = \mathtt{inc}$, *where* $\Delta_t = (D_t, \mathcal{FD})$ *and* $D_t = D \cup \{t\}$.

*Proof* Assuming that $v_\Delta(t) = \mathtt{false}$ indeed entails that $\Delta \not\vdash t$ by Definition 5. Moreover, as $\Delta \not\hspace{-0.3em}\sim t$, there exist $A$ in $U$ and $a$ and $a'$ in $dom(A)$ such that $\Delta \vdash (t \preceq a \sqcap a')$. Now, given $\mu$ such that $\mu \models \Delta_t$, it has been shown that $\mu$ also satisfies that $\mu \models \Delta$ (see Appendix B). Hence, $\mu(t) \subseteq \mu(a) \cap \mu(a')$ holds, which implies that $\Delta_t \vdash (t \preceq a \sqcap a')$, that is $\Delta_t \not\hspace{-0.3em}\sim t$. Since it holds that $\Delta_t \vdash t$, we obtain that $v_{\Delta_t}(t) = \mathtt{inc}$.

**Algorithm 2** Chasing a table

**Input:** $\Delta = (D, \mathcal{FD})$
**Output:** The chased table $\Delta^* = (D^*, \mathcal{FD})$ and a set $inc(\mathcal{FD})$ containing sets of tuples
    associated with each $X \to A$ in $\mathcal{FD}$
1: $D^* := D$
2: **for all** $X \to A$ in $\mathcal{FD}$ **do**
3:     $inc(X \to A) := \emptyset$
4: **while** $D^*$ changes **do**
5:     **for all** $X \to A \in \mathcal{FD}$ **do**
6:         **for all** $t_1$ in $D^*$ such that $XA \subseteq sch(t_1)$ **do**
7:             **for all** $t_2$ in $D^*$ such that $X \subseteq sch(t_2)$ and $t_1.X = t_2.X$ **do**
8:                 **if** $A \notin sch(t_2)$ **then**
9:                     $D^* := D^* \cup \{t_2 a_1\}$ where $a_1 = t_1.A$
10:                **if** $A \in sch(t_2)$ and $t_1.A \neq t_2.A$ **then**
11:                   Let $y_i = t_i.(sch(t_i) \setminus A)$ and $a_i = t_i.A$, for $i = 1, 2$
12:                   $D^* := D^* \cup \{y_2 a_1\}$
                  // $y_1 a_2$ is inserted into $D^*$ when processing $t_2$ in place of $t_1$
                  // and $t_1$ in place of $t_2$
13:                   $inc(X \to A) := inc(X \to A) \cup \{x\}$ where $x = t_1.X = t_2.X$
    // Reduction: keep in $D^*$ only maximal tuples
14: **for all** $t_1$ in $D^*$ **do**
15:     **for all** $t_2$ in $D^*$ **do**
16:         **if** $t_2 \sqsubseteq t_1$ and $t_1 \neq t_2$ **then**
17:             $D^* := D^* \setminus \{t_2\}$
18: $inc(\mathcal{FD}) := \{inc(X \to A) \mid inc(X \to A) \neq \emptyset\}$
19: **return** $\Delta^* = (D^*, \mathcal{FD})$ and $inc(\mathcal{FD})$

Conversely, if $v_{\Delta_t}(t) = \texttt{inc}$ then, by Definition 5, $\Delta_t \vdash t$ and $\Delta_t \not\vdash t$ hold. Therefore, there exist $A$ in $U$ and $a$ and $a'$ in $dom(A)$ such that $\Delta_t \vdash (t \preceq a \sqcap a')$, which by Algorithm 1 and Lemma 2, implies that $a$ and $a'$ are in $t^+$. We thus obtain that $\Delta \not\vdash t$, which combined with our hypothesis that $\Delta \nvdash t$, implies that $v_\Delta(t) = \texttt{false}$. The proof is therefore complete.     $\square$

## 4 Computing the Semantics

Similarly to standard two valued logic, where computing the semantics of $\Delta$ means computing the set of all tuples true in $\Delta$, in our approach, computing the semantics amounts to compute all true, inconsistent or false tuples, knowing that unknown tuples are the remaining ones.

However, as mentioned above, the set of false tuples may be infinite, making it impossible to compute them all. In this work, the case of false tuples is only partially addressed, and we rather concentrate on potentially true tuples, with the goal of investigating consistent query answering in our approach (see Section 6).

### 4.1 The Chase Procedure in our Approach

We first propose an effective algorithm for the computation of all potentially true tuples in a given $\Delta$. This algorithm is in fact inspired by the standard chase algorithm [21, 24], with the main difference that when a functional dependency cannot be satisfied, our algorithm does *not* stop.

Instead, our chasing algorithm carries on the computation, returning a database $\Delta^* = (D^*, \mathcal{FD})$ and a set $inc(\mathcal{FD})$ based on which inconsistent and true tuples are shown to be efficiently computed. Before doing so, we illustrate Algorithm 2 in the context of our introductory example.

| $D$ | $Id$ | $K$ | $M$ | $C$ |
|---|---|---|---|---|
| | $i_1$ | $k$ | $m$ | $c$ |
| | $i_1$ | | $m'$ | |
| | $i_1$ | $k$ | | $c$ |
| | $i_2$ | $k'$ | $m'$ | $c$ |
| | $i_2$ | $k'$ | $m''$ | |
| | $i_2$ | $k'$ | | $c'$ |
| | $i_3$ | | $m$ | |
| | $i_3$ | $k'$ | | |

| $D^*$ | $Id$ | $K$ | $M$ | $C$ |
|---|---|---|---|---|
| | $i_1$ | $k$ | $m$ | $c$ |
| | $i_1$ | $k$ | $m'$ | $c$ |
| | $i_2$ | $k'$ | $m'$ | $c$ |
| | $i_2$ | $k'$ | $m''$ | $c$ |
| | $i_2$ | $k'$ | $m'$ | $c'$ |
| | $i_2$ | $k'$ | $m''$ | $c'$ |
| | $i_3$ | $k'$ | $m$ | |

**Fig. 2** The table $D$ and its chased version $D^*$

*Example 8* Running Algorithm 2 with the table $D$ shown in Figure 1, and recalled in Figure 2, produces the table $D^*$ shown in the right of Figure 2 and the set $inc(\mathcal{FD}) = \{inc(Id \to M), inc(Id \to C)\}$ where $inc(Id \to M) = \emptyset$ and $inc(Id \to C) = \{i_2\}$. The main steps of the algorithm work as follows:

• First, $D^*$ is assigned $D$, and $inc(Id \to M)$ and $inc(Id \to C)$ are assigned $\emptyset$.

• Due to the statement on line 9, the first two rows in $D$ (thus in $D^*$) generate the new tuples $(i_1, k, m')$ and $(i_1, m', c)$. Similarly, applying $Id \to K$ to the last two rows in $D$ generates the new tuple $(i_3, k', m)$.
The rows 4 and 5 in $D$ generate $(i_2, k', m'', c)$ and the rows 5 and 6 generate $(i_2, k', m'', c')$. Moreover, due to the statement on line 12, the rows 4 and 6 generate $(i_2, k', m', c')$ and $(i_2, k', c)$ and $i_2$ is inserted in $inc(Id \to C)$.

• With these new tuples at hand, the loop on line 4 proceeds further, generating $(i_1, k, m', c)$ by the statement on line 9. No new tuple is generated at this stage.

• The loop on line 4 is processed once again, producing no new tuple. When running the reduction step against the current state of $D^*$, the following tuples are removed: $(i_1, m')$, $(i_1, k, m')$, $(i_1, m', c)$, $(i_1, k, c)$, $(i_2, k', m'')$, $(i_2, k', c')$, $(i_2, k', c)$, $(i_3, m)$ and $(i_3, k')$.

Thus, the output of Algorithm 2 is indeed as expected. It is important to notice that, although tuples have been added in $D^*$ during the processing, the final number of tuples in $D^*$ is less than that in $D$. Although this particular result cannot be proven in general, it will be shown that in the worst case, the size of $D^*$ remains polynomial in the size of $D$.

We emphasize that some nulls present in $D$ have been replaced by actual values in $D^*$, thanks to the functional dependencies in $\mathcal{FD}$. For example the second tuple in $D$ with two nulls has been 'completed' into a total tuple in $D^*$. However, such a completion has not been possible for *every* tuple in $D^*$. Namely, the $C$-value in the last tuple of $D^*$ is left as null.

Keeping in line with our statement that 'a missing value exists only if it is inferred from the functional dependencies', this indicates that the $C$-value of this tuple could *not* be determined based on the content of $D$ and $\mathcal{FD}$, and no other conclusion can be drawn regarding this null.

14

To see why the two insertions mentioned in the statement on line 12 are needed, we first recall from [24] that, in the traditional case, the chased table characterizes the semantics of the input table, in case no inconsistency has been detected[1]. In this work our goal is similar, but has to be adapted to our context. Namely, we expect that the chased table $D^*$ can provide a syntactical characterization of all possibly true tuples in $\Delta$, that is of all tuples $t$ such that $\Delta \vdash t$ holds.

In the context of our example, if we assume that $(i_2, k', m', c')$ is not inserted during the processing then Algorithm 2 would *not* fit our semantics. Indeed, for every $\mathcal{T}$-mapping $\mu$ such that $\mu \models \Delta$, $\mu(i_2) \subseteq \mu(c')$ holds because of $Id \to C$ applied to the seventh row in $D$. Thus, $\mu(i_2, k', m', c') = \mu(k', m', c')$, and since $\mu(k', m', c')$ is nonempty (due to the fifth row in $D$), $\Delta \vdash (i_2, k', m', c')$. Hence $(i_2, k', m', c')$ must appear in $D^*$ to fulfill our expectation.

Adding such 'new' tuples when chasing a table is one of the main features of our approach, as compared with traditional chase. This step should be seen as a 'by-product' of carrying on the computation even after encountering a violation of a functional dependency. □

The following lemma shows that Algorithm 2 provides an operational means to characterize the tuples $t$ such that $\mu^*(t) \neq \emptyset$.

**Lemma 3** *Algorithm 2 applied to $\Delta = (D, \mathcal{FD})$ always terminates. Moreover, for every tuple $t$, $\mu^*(t) \neq \emptyset$ holds if and only if $t$ is in $\mathsf{LoCl}(D^*)$.*

*Proof* See Appendix C. □

Recalling that $\mathsf{LoCl}(D^*)$ denotes the Lower Closure of $D^*$, that is the set of all sub-tuples of tuples in $D^*$, Lemma 3 shows that $D^*$ is a 'tabular' version of the set of all tuples $t$ such that $\mu^*(t) \neq \emptyset$, that is, by Lemma 1, a 'tabular' version of the set of all tuples $t$ such that $\Delta \vdash t$. Therefore, $D^*$ provides a syntactical characterization of the set of all tuples $t$ such that $\Delta \vdash t$, as expected in the previous example.

4.2 Computing True Tuples and Inconsistent Tuples

As mentioned just above, Lemma 1 and Lemma 3 show that, given $\Delta = (D, \mathcal{FD})$, a tuple $t$ is in $\mathsf{LoCl}(D^*)$ if and only if $\Delta \vdash t$ holds, that is, if and only if $t$ is potentially true in $\Delta$, that is if and only if $t$ is either true or inconsistent in $\Delta$.

To see how to compute the set of all inconsistent tuples, we first recall the notion of *closure of a relation scheme* as defined in relational database theory [24].

Given a set $\mathcal{FD}$ of functional dependencies and a relation scheme $X$, the *closure of $X$ with respect to $\mathcal{FD}$*, or more simply the *closure of $X$*, denoted by $X^+$, is the set of all attributes $A$ in $U$ such that every table $D$ satisfying $\mathcal{FD}$ in the sense of relational tables, also satisfies $X \to A$.

It is well-known that $X^+$ is computed through the following two steps that are quite similar to the steps of Algorithm 1:

---

[1] In traditional chase, the semantics of a table $D$ containing nulls is the set of all tuples true in every instance of $D$, *i.e.,* in every relation $R$ over $U$ with no nulls, that satisfies the functional dependencies and such that for every $t$ in $D$, there exists $r$ in $R$ such that $r.T = t$.

$X^+ := X$
**while** $X^+$ changes **do**
    **for all** $Y \to B$ in $\mathcal{FD}$ such that $Y \subseteq X^+$ **do**
        $X^+ := X^+ \cup \{B\}$
**return** $X^+$

The following proposition shows a strong relationship between the closure of a relation scheme as recalled above and the closure of a tuple as stated in Definition 3.

**Proposition 3** *Let $\Delta = (D, \mathcal{FD})$ and $t$ be such that $\Delta \vdash t$. For every tuple $q$ and every $a$ in $dom(A)$ such that $q \sqsubseteq t$ and $a \sqsubseteq t$, we have: $a$ belongs to $q^+$ if and only if $A$ belongs to $Q^+$.*

*Proof* See Appendix D. ☐

---

**Algorithm 3** Inconsistent tuples in $\Delta = (D, \mathcal{FD})$

---

**Input:** The output of Algorithm 2, that is $\Delta^* = (D^*, \mathcal{FD})$ and $inc(\mathcal{FD})$.
**Output:** The set $\mathsf{Inc}(\Delta)$
 1: $\mathsf{Inc}(\Delta) := \emptyset$
 2: **for all** $t$ in $D^*$ **do**
 3:     **for all** $X \to A$ in $\mathcal{FD}$ such that $XA \subseteq T$ **do**
 4:         **if** $x = t.X$ is in $inc(X \to A)$ **then**
 5:             **for all** $q$ such that $q \sqsubseteq t$ **do**
 6:                 **if** $X \subseteq Q^+$ **then**
 7:                     $\mathsf{Inc}(\Delta) := \mathsf{Inc}(\Delta) \cup \{t.Q\}$
 8: **return** $\mathsf{Inc}(\Delta)$

---

Using the notion of relation scheme closure, we introduce Algorithm 3 which computes the set of inconsistent tuples in $\Delta$. The correctness of this algorithm is shown in Lemma 4.

**Lemma 4** *Given $\Delta = (D, \mathcal{FD})$, a tuple $t$ is inconsistent in $\Delta$ if and only if $t \in \mathsf{Inc}(\Delta)$.*

*Proof* See Appendix E. ☐

The following proposition characterizes inconsistent and true tuples in $\Delta$ based on Algorithm 2 and Algorithm 3.

**Proposition 4** *Given $\Delta = (D, \mathcal{FD})$ and a tuple $t$:*
*1. $t$ is inconsistent in $\Delta$ if and only if $t \in \mathsf{Inc}(\Delta)$.*
*2. $t$ is true in $\Delta$ if and only if $t \in \mathsf{LoCl}(D^*) \setminus \mathsf{Inc}(\Delta)$.*

*Proof* Immediate consequence of Definition 5, Lemma 3 and Lemma 4. ☐

The following examples illustrate Algorithm 3 and Proposition 4.

*Example 9* As in Example 7, let $\Delta = (D, \mathcal{FD})$ over $U = \{A, B, C\}$ where $D = \{abc, ac'\}$ and $\mathcal{FD} = \{A \to B, B \to C\}$. The tabular version of $D$ is shown on the left below, whereas $D^*$ is shown on the right.

**Algorithm 4** Tuple truth value in $\Delta = (D, \mathcal{FD})$

**Input:** A tuple $t$, $\Delta^* = (D^*, \mathcal{FD})$ and $\mathsf{Inc}(\Delta)$
**Output:** The truth value of $t$ as one of the truth values `true`, `false`, `inc` or `unkn`
1: $v := \mathtt{unkn}$
2: **if** $t \in \mathsf{LoCl}(\Delta^*)$ **then**
3:      **if** $t \in \mathsf{Inc}(\Delta)$ **then**
4:          $v := \mathtt{inc}$
5:      **else**
6:          $v := \mathtt{true}$
7: **else**
8:      Compute $D_t^*$ using Algorithm 2 applied to $D^* \cup \{t\}$ and $\mathcal{FD}$
9:      Compute $\mathsf{Inc}(\Delta_t)$ using Algorithm 3 applied to $\Delta_t = (D^* \cup \{t\}, \mathcal{FD})$
10:     **if** $t \in \mathsf{Inc}(\Delta_t)$ **then**
11:        $v := \mathtt{false}$
12: **return** $v$

| $D$ | $A$ | $B$ | $C$ |
|---|---|---|---|
| | $a$ | $b$ | $c$ |
| | $a$ | | $c'$ |

| $D^*$ | $A$ | $B$ | $C$ |
|---|---|---|---|
| | $a$ | $b$ | $c$ |
| | $a$ | $b$ | $c'$ |

Running Algorithm 2, $D^*$ is first set to $D$ and $abc'$ is inserted in $D^*$ by the statement line 9 due to $A \to B$. Then, $b$ is inserted in $inc(B \to C)$ by the statement line 13, due to the tuples $abc$ and $abc'$. Thus, the table $D^*$ output by Algorithm 2 is as shown above and $inc(\mathcal{FD}) = \{inc(A \to B), inc(B \to C)\}$ where $inc(A \to B) = \emptyset$ and $inc(B \to C) = \{b\}$.

When running Algorithm 3 for $abc$ in $D^*$, since $b$ is in $inc(B \to C)$, $b$, $ab$, $bc$ and $abc$ are inserted into $\mathsf{Inc}(\Delta)$, due to the statement on line 7. This is so because the schema $Q$ of each of these tuples contains $B$, and so, satisfies $B \subseteq Q^+$.

Moreover, for $q = a$, due to $A \to B$, we have $A^+ = ABC$ and thus, $B \subseteq A^+$ holds, showing that $a$ is inserted in $\mathsf{Inc}(\Delta)$ on line 7. A similar reasoning holds for $q = ac$ because $B \in (AC)^+$. Thus, $ac$ is also inserted in $\mathsf{Inc}(\Delta)$ on line 7. The only remaining possibility is $q = c$, and does not modify $\mathsf{Inc}(\Delta)$ because $B \not\subseteq C^+$. A similar computation is performed with $abc'$ in $D^*$, adding $bc'$, $abc'$ and $ac'$ in $\mathsf{Inc}(\Delta)$. As no other tuple can be inserted in $\mathsf{Inc}(\Delta)$, Algorithm 3 returns

$$\mathsf{Inc}(\Delta) = \{abc, abc', ab, ac, ac', bc, bc', a, b\},$$

which, by Proposition 4(1), is the set of all inconsistent tuples in $\Delta$. As a consequence, by Proposition 4(2), $c$ and $c'$ are the only true tuples in $\Delta$.

Now, as in Example 7, referring to $\Delta_1 = (D_1, \mathcal{FD})$ with $D_1 = \{ac, ac'\}$, it is easy to see that $D_1^* = D_1$. This implies that $\Delta_1$ is consistent, and that $ac$, $ac'$, $a$, $c$ and $c'$ are true in $\Delta_1$. $\qquad\square$

### 4.3 The Case of False Tuples

As already noticed, computing all tuples false in a given $\Delta = (D, \mathcal{FD})$ is not feasible in case of infinite attribute domains. However, given a tuple $t$ and assuming that $\Delta^*$ and $\mathsf{Inc}(\Delta)$ have been computed, Algorithm 4 allows to compute $v_\Delta(t)$. In this way, instead of being systematically identified, false tuples are identified on demand.

**Proposition 5** *Given $\Delta = (D, \mathcal{FD})$, and a tuple $t$ and assuming that $\Delta^* = (D^*, \mathcal{FD})$ and $\mathsf{Inc}(\Delta)$ have been computed, the truth value returned by Algorithm 4 is equal to $v_\Delta(t)$.*

*Proof* Immediate consequence of Definition 5 and Proposition 2. □

The following example illustrates the algorithm.

*Example 10* It has been seen in Example 8 that in the context of our introductory example, Algorithm 2 returns the table $D^*$ as shown in Figure 2, and $inc(\Delta) = \{inc(Id \rightarrow K), inc(Id \rightarrow C)\}$ where $inc(Id \rightarrow K) = \emptyset$ and $inc(Id \rightarrow C) = \{i_2\}$. Thus, by Algorithm 3, the set $\mathsf{Inc}(\Delta)$ is defined by:

$$\mathsf{Inc}(\Delta) = \quad \{t \mid i_2 \sqsubseteq t \sqsubseteq (i_2, k', m', c)\} \cup \{t \mid i_2 \sqsubseteq t \sqsubseteq (i_2, k', m', c')\} \cup$$
$$\{t \mid i_2 \sqsubseteq t \sqsubseteq (i_2, k', m'', c)\} \cup \{t \mid i_2 \sqsubseteq t \sqsubseteq (i_2, k', m'', c')\}$$

Applying now Algorithm 4, we have the following:

- $v_\Delta(i_1, a, m, c) = v_\Delta(i_1, k, m', c) = \mathtt{true}$, because line 6 changes the value of $v$, since these tuples are in $D^*$ but not in $\mathsf{Inc}(\Delta)$.
- $v_\Delta(i_2) = v_\Delta(i_2, c) = v_\Delta(i_2, c') = \mathtt{inc}$, because line 4 changes the value of $v$, since these tuples are in $\mathsf{Inc}(\Delta)$.
- $v_\Delta(i_1, k') = v_\Delta(i_1, c') = \mathtt{false}$. Indeed, none of these tuples is in $\mathsf{LoCl}(D^*)$, thus implying that neither line 4 nor line 6 applies. Moreover, when running Algorithm 4 with $(i_1, k')$ as input, $(i_1, k, m, c)$ and $(i_1, k')$ are in $D^* \cup \{(i_1, k')\}$. Hence $(i_1, k')$ is in $\mathsf{Inc}(\Delta_t)$, because of $Id \rightarrow K$, and by line 11, $v_\Delta(k'm)$ is set to $\mathtt{false}$. A similar reasoning holds for $(i_1, c')$, but using $Id \rightarrow C$.
- $v_\Delta(k', m) = \mathtt{unkn}$. Indeed, as above, when running Algorithm 4 with $(k', m)$ as input, lines 4 and 6 do not change the value of $v$ (as $(k', m)$ does not occur in $\mathsf{LoCl}(D^*)$). Moreover, as $D^* \cup \{(k', m)\}$ is consistent, the value of $v$ is not changed by the statement line 11. Consequently, $\mathtt{unkn}$ is returned. □

4.4 Complexity Issues

We argue that the computation of inconsistent and true tuples in $\Delta = (D, \mathcal{FD})$ is polynomial in the size of the table $D$ and in the order of the 'number of conflicts with respect to functional dependencies' (to be defined shortly). To see this, denoting by $|E|$ the cardinality of a set $E$, we investigate the complexities of Algorithm 2 and of Algorithm 3.

Regarding Algorithm 2, we first notice that, contrary to the standard chase algorithm [24], rows are added in the table during the computation, and some others are then removed by the reduction statement of line 14. To assess the size of the table $D^*$ during the processing, we point out the following:

- If no inconsistency is found during the processing of the while-loop on line 4, at most one tuple is added in $D^*$ as the 'join' of two tuples in $D$ by statement line 9. Therefore, the cardinality of $D^*$ remains in the same order as that of $D$. Notice in this respect that, upon reduction, *one* 'join' tuple replaces *two* tuples in $D$, which reduces the size of the table $D^*$ output by the algorithm.
- However, when inconsistent tuples occur, the statement line 9 adds more than one tuple and statement line 12 inserts tuples resulting from a cross-product.

To find an upper bound of the size of $D^*$, for every $X \to A$ in $\mathcal{FD}$, let $N(x)$ be the number of different $A$-values $a$ such that $\Delta \vdash xa$ and $x$ belongs to $inc(X \to A)$. We denote by $\delta$ the maximal value of all $N(x)$ for all $x$ in $inc(\mathcal{FD})$; in other words $\delta = \max(\{N(x) \mid x \in inc(\mathcal{FD})\})$. $\delta$ is precisely what was earlier referred to as the 'number of conflicts with respect to functional dependencies'.

Given a tuple in $D$ and a functional dependency $X \to A$ in $\mathcal{FD}$, each of the statements line 9 and line 12 generates at most $\delta$ tuples. Since several functional dependencies may apply to $t$, at most $\delta^{|\mathcal{FD}|}$ tuples are generated for the given tuple $t$. Hence, the number of tuples generated by the statements lines 9 and 12 is in $\mathcal{O}\left(|D|.\delta^{|\mathcal{FD}|}\right)$. We therefore obtain that the size of the table $D^*$ when running Algorithm 2 is in $\mathcal{O}\left(|D|.\left(1 + \delta^{|\mathcal{FD}|}\right)\right)$, that is in $\mathcal{O}\left(|D|.\delta^{|\mathcal{FD}|}\right)$.

Since the number of runs of the while-loop on line 4 is at most equal to the number of tuples added into $D^*$, this number is in $\mathcal{O}\left(|D|.\delta^{|\mathcal{FD}|}\right)$. Since moreover one run of the while-loop is quadratic in the size of $D^*$, the computational complexity of this while-loop is in $\mathcal{O}\left(|D|^3.\delta^{3.|\mathcal{FD}|}\right)$.

The last point to be mentioned here is that the reduction processing on line 14 is performed through a scan $D^*$ whereby for every $t$ in $D^*$ every sub-tuple of $t$ is removed. Such a processing being quadratic in the size of $D^*$, the overall computational complexity of Algorithm 2 is in $\mathcal{O}\left(|D|^3.\delta^{3.|\mathcal{FD}|}\right)$.

As the computational complexity of Algorithm 3 is clearly linear in the size of $D^*$, the global complexity of the computation of inconsistent and true tuples in $\Delta$ is as stated just above, and therefore *polynomial in the size of $D$*.

Regarding Algorithm 4, we notice that its complexity is in $\mathcal{O}\left(|D|^3.\delta^{3.|\mathcal{FD}|}\right)$ as well, because it requires a scan of $D^*$ and then, in case the test line 2 fails, Algorithm 2 is applied to a table whose cardinality is that of $D^*$ plus 1. It should however be kept in mind that, in this case, the algorithm has to be run *once for each tuple*, which shows that computing false tuples is not feasible even if all attribute domains are finite. Indeed, in this case, denoting by $DOM$ the maximal cardinality of attribute domains, the cardinality of $\mathcal{T}$ is in $\mathcal{O}\left(|U|^{DOM}\right)$, thus yielding a computation in $\mathcal{O}\left(|U|^{DOM}.|D|^3.\delta^{3.|\mathcal{FD}|}\right)$.

We draw attention on the following important points regarding these complexity results:

1. Regarding the computation of false tuples, the above result has to be further investigated in the following two directions: first the computation of $D_t^*$ processed in Algorithm 4 is likely to be optimized using an *incremental* algorithm instead of Algorithm 2, and second, it is expected that there exist interesting and relevant cases whereby the computation of $D_t^*$ is *not necessary*. We indeed suspect that this holds in the case of a star schema. This is an important issue that lies out of the scope of the present paper, but that will be investigated in the next future.

2. When the database is consistent, $\delta$ is equal to 1, thus yielding a complexity in $\mathcal{O}(|D|^3)$. This result can be shown independently from the above computations as follows: In the case of traditional chase the maximum of nulls in $D$ being bounded by $|U|.|D|$, the number or iterations when running the algorithm is

19

| $\varphi$ | $\neg\varphi$ |   | $\vee$ | t | b | n | f |   | $\wedge$ | t | b | n | f |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| t | f |   | t | t | t | t | t |   | t | t | b | n | f |
| b | b |   | b | t | b | t | b |   | b | b | b | f | f |
| n | n |   | n | t | t | n | n |   | n | n | f | n | f |
| f | t |   | f | t | b | n | f |   | f | f | f | f | f |

| $\oplus$ | t | b | n | f |   | $\otimes$ | t | b | n | f |
|---|---|---|---|---|---|---|---|---|---|---|
| t | t | b | t | b |   | t | t | t | n | n |
| b | b | b | b | b |   | b | t | b | n | f |
| n | t | b | n | f |   | n | n | n | n | n |
| f | b | b | f | f |   | f | n | f | n | f |

**Fig. 3** Truth tables of basic connectors

also bounded by $|U|.|D|$. Since the run of one iteration is in $|D|^2$, the overall complexity is in $\mathcal{O}(|U|.|D|^3)$, or in $\mathcal{O}(|D|^3)$, as $|U|$ is independent from $|D|$.

3. The above complexity study should be further investigated in order to provide more accurate results regarding the estimation of the number of actual tests necessary to the computation of $D^*$. The results in [11] are likely to be useful for such a more thorough study of this complexity.

## 5 Four-Valued Logic and Table Merging

In this section, we first give a brief overview of Belnap's Four-valued logic and then we show that our approach has a strong relationship with this formalism in the context of merging two or more tables.

### 5.1 Basics of Four-Valued Logic

Four-valued logic was introduced by Belnap in [5], who argued that his formalism is of interest when integrating data from various data sources. To this end, he introduced four truth values denoted by t, b, n and f and read as *true*, *both true and false*, *neither true nor false* and *false*, respectively. An important feature of this Four-valued logic is that its truth values can be compared according to two partial orderings, known as *truth ordering* and *knowledge ordering*, respectively denoted by $\preceq_t$ and $\preceq_k$ and defined as follows:

$$\mathtt{n} \preceq_k \mathtt{t} \preceq_k \mathtt{b} \; ; \mathtt{n} \preceq_k \mathtt{f} \preceq_k \mathtt{b} \qquad \text{and} \qquad \mathtt{f} \preceq_t \mathtt{n} \preceq_t \mathtt{t} \; ; \mathtt{f} \preceq_t \mathtt{b} \preceq_t \mathtt{t}.$$

As a consequence, two new connectors were introduced, denoted by $\oplus$ and $\otimes$, in addition to the standard connectors $\vee$ (disjunction) and $\wedge$ (conjunction). The corresponding truth tables, along with that for negation, are displayed in Figure 3 and show that $\vee$ and $\oplus$ correspond to the least upper bound (lub) with respect to $\preceq_t$ and $\preceq_k$, respectively; whereas $\wedge$ and $\otimes$, correspond to the geatest lower bound (glb) with respect to $\preceq_t$ and $\preceq_k$, respectively . It is also shown in [5,13] that the set $\{\mathtt{t},\mathtt{b},\mathtt{n},\mathtt{f}\}$ equipped with the two orderings $\preceq_t$ and $\preceq_k$ has a distributive bi-lattice structure.

Not surprisingly, some basic properties holding in standard logic do not hold in this setting. For example, Figure 3 shows that formulas of the form $\Phi \vee \neg\Phi$ are not

always true, independently of the truth value of $\Phi$. The reader is referred to the literature [3, 5, 13, 14, 23] for more details on the properties of Four-valued logic.

Based on the truth tables shown in Figure 3, it turns out that the connector $\oplus$ plays a key role in the context of data integration. Indeed, considering $n$ data sources $S_1, \ldots, S_n$ and a fact $\varphi$, for every $i = 1, \ldots, n$, $\varphi$ is assigned one truth value $v_i$, among $\mathtt{t}$, $\mathtt{b}$, $\mathtt{n}$, or $\mathtt{f}$ in each $S_i$. The 'integrated' truth value of $\varphi$, denoted by $v$ is then obtained as the expression $v = v_1 \oplus \ldots \oplus v_n$, due to the following intuition:

- The third row (or third column) of the truth table of $\oplus$ shows that every $v_i$ such that $v_i = \mathtt{n}$ plays no role in the resulting truth value $v$, provided that one of them be distinct from $\mathtt{n}$ (otherwise the 'integrated' truth value of $\varphi$ is obviously $\mathtt{n}$). This fits our intuition that a source in which the truth value of $\varphi$ is unknown does not provide any piece of information regarding the 'integrated' truth value of $\varphi$. We thus assume hereafter that for every $i = 1, \ldots, n$, $v_i \neq \mathtt{n}$.
- For every $\mathtt{v}$ among $\mathtt{t}$, $\mathtt{b}$, $\mathtt{n}$ or $\mathtt{f}$, if $v_1 = \ldots = v_n = \mathtt{v}$, then $v = \mathtt{v}$. The intuition here is that, since all sources agree on truth value $\mathtt{v}$, it is obvious to expect $v$ to be this common value $\mathtt{v}$. For example, if for every $i = 1, \ldots, n$, $v_i = \mathtt{t}$, then it should be obvious that $v$ must be $\mathtt{t}$ as well!
- Now, if there exists $i_0$ such that $v_{i_0} = \mathtt{b}$, then $v = \mathtt{b}$. This fits the intuition that if $\varphi$ is inconsistent in at least one data source, then $\varphi$ remains inconsistent in the integrated source.
- The last case is when there exist distinct $i$ and $j$ in $\{1, \ldots, n\}$ such that $v_i \neq v_j$, and no $v_i$ is equal to $\mathtt{b}$. In this case we have $v_i = \mathtt{t}$ and $v_j = \mathtt{f}$ (or equivalently $v_i = \mathtt{f}$ and $v_j = \mathtt{t}$), which is the standard case of conflicting data sources in practice. In this case, it holds that $v = \mathtt{b}$ (since $\mathtt{t} \oplus \mathtt{f} = \mathtt{b}$). This result again fits our intuition that in case of conflicting data sources, the 'integrated' truth value in inconsistent.

In the next sub-section, we show that, in our approach, the four truth values as defined in Definition 5 also follow this intuition when it comes to merging two or more tables over the same universe $U$.

5.2 Merging two or more Tables

Data merging consists in collecting data from multiple, possibly heterogeneous sources and putting them in a single destination. The data from each source usually comes in the form of a CSV file, along with some hints on the data, referred to as metadata [16, 19]. During this process, different data sources are put together, or merged, into a single data store. Data merging is also related to data consolidation and to data integration.

When data comes from a broad range of sources, consolidation allows organizations to more easily present data, while also facilitating effective data analysis. Data consolidation techniques reduce inefficiencies, like data duplication, costs related to reliance on multiple databases and multiple data management points.

In this section, we consider a simplified, relational scenario of $n$ sources $\Delta_1 = (D_1, \mathcal{FD}_1), \ldots, \Delta_n = (D_n, \mathcal{FD}_n)$, where each source $\Delta_i = (D_i, \mathcal{FD}_i)$ consists of a table $D_i$ over a fixed universe $U$, possibly with nulls, and functional dependencies

$\mathcal{FD}_i$. We then explain how to merge these sources in our approach under the following assumptions:

1. All source tables are over the *same* universe $U$.
2. Merging is done in the simplest possible way, namely $(a)$ the merged table is the *union* (in the set theoretic sense) of the source tables and $(b)$ the set of functional dependencies of the merged table is the union of the sets of functional dependencies of the source tables. That is, the sources are merged through the pair: $\Delta = (D, \mathcal{FD})$, where $D = \bigcup_{i=1}^{i=n} D_i$ and $\mathcal{FD} = \bigcup_{i=1}^{i=n} \mathcal{FD}_i$.

Relying on Belnap's Four-valued logic, we investigate the relationship between the truth values a tuple $t$ has in the source tables and the truth value the tuple $t$ has in the merged table.

First, notice that a 'natural' one-to-one mapping $h$ from our set $\mathsf{Four} = \{\mathtt{true}, \mathtt{inc}, \mathtt{unkn}, \mathtt{false}\}$ to Belnap's set $\mathcal{FOUR} = \{\mathtt{t}, \mathtt{b}, \mathtt{n}, \mathtt{f}\}$, can be defined by: $h(\mathtt{true}) = \mathtt{t}$, $h(\mathtt{inc}) = \mathtt{b}$, $h(\mathtt{unkn}) = \mathtt{n}$ and $h(\mathtt{false}) = \mathtt{f}$. Then, the connector $\oplus$ defined on $\mathcal{FOUR}$ induces a connector $\overline{\oplus}$ over $\mathsf{Four}$ defined by: $\mathtt{v}_1 \overline{\oplus} \mathtt{v}_2 = h^{-1}(h(\mathtt{v}_1) \oplus h(\mathtt{v}_2))$ for all $\mathtt{v}_1$ and $\mathtt{v}_2$ in $\mathsf{Four}$.

Moreover, we can define a partial ordering on $\mathsf{Four}$ isomorphic to the knowledge ordering of $\mathcal{FOUR}$ that allows us to compare truth values in $\mathsf{Four}$. Denoting this partial ordering by $\lhd$, we have:

$$\mathtt{unkn} \lhd \mathtt{false} \lhd \mathtt{inc} \quad \text{and} \quad \mathtt{unkn} \lhd \mathtt{true} \lhd \mathtt{inc}$$

The following proposition shows that the truth value of a tuple $t$ in the merged table is always greater (with respect to $\lhd$) than any of the truth values that $t$ has in the source tables in which it appears. In other words, when merging tables, the knowledge about tuples always increases, compared to the knowledge we have about tuples in the source tables.

**Proposition 6** *Let $\Delta_i = (D_i, \mathcal{FD}_i)$ $(i = 1, \ldots, n)$ be $n$ data sources over the same universe, and let $\Delta = (D, \mathcal{FD})$ be defined by $D = \bigcup_{i=1}^{i=n} D_i$ and $\mathcal{FD} = \bigcup_{i=1}^{i=n} \mathcal{FD}_i$. For every tuple $t$ the following holds:*

$$\overline{\bigoplus}_{i=1}^{i=n} v_{\Delta_i}(t) \ \lhd \ v_{\Delta}(t).$$

*Proof* For every $i = 1, \ldots, n$, let $\Delta_i' = (D_i, \mathcal{FD})$. We first prove that for every tuple $t$, $v_{\Delta_i}(t) \lhd v_{\Delta_i'}(t)$ holds. Indeed, for every $i = 1, \ldots, n$, let $D_i^*$, respectively $(D_i')^*$, the chased table of $D_i$ with respect to $\mathcal{FD}_i$, respectively $\mathcal{FD}$. Since $\mathcal{FD}_i \subseteq \mathcal{FD}$ holds, it is easy to see that for every $q_i$ in $(D_i')^*$ there exists $q$ in $D_i^*$ such that $q_i \sqsubseteq q$. Hence, for every $q$ in $\mathcal{T}$, $[q^+]_i \subseteq [q^+]_i'$, where $[q^+]_i$, respectively $[q^+]_i'$, denotes the closure of $q$ in $\Delta_i$, respectively $\Delta_i'$. Therefore, if $\Delta_i \vdash t$, respectively $\Delta_i \mathrel{\vhook\sim} t$, then $\Delta_i' \vdash t$, respectively $\Delta_i' \mathrel{\vhook\sim} t$, and so, for every $i = 1, \ldots, n$, $v_{\Delta_i}(t) \lhd v_{\Delta_i'}(t)$.

Considering $\Delta_i'$ $(i = 1, \ldots, n)$ and $\Delta$, it can be seen that for every $i = 1, \ldots, n$ and every $q_i$ in $D_i'^*$ there exists $q$ in $D^*$ such that $q_i \sqsubseteq q$. Consequently, for every $i = 1, \ldots, n$, and every $q$ in $\mathcal{T}$, $[q^+]_i' \subseteq q^+$, where $q^+$ denotes the closure of $q$ in $\Delta$. Therefore, if for some $i$, $\Delta_i' \vdash t$, respectively $\Delta_i' \mathrel{\vhook\sim} t$, then $\Delta \vdash t$, respectively $\Delta \mathrel{\vhook\sim} t$, and so, for every $i = 1, \ldots, n$, $v_{\Delta_i'}(t) \lhd v_{\Delta}(t)$. The proposition follows from the transitivity of $\lhd$ and from the fact that $\overline{\oplus}$ defines the least upper bound (lub) with respect to $\lhd$, in the same way as $\oplus$ defines the lub with respect to $\preceq_k$. $\qquad\square$

In what follows, we identify cases where the equality $\overline{\bigoplus}_{i=1}^{i=n} v_{\Delta_i}(t) = v_{\Delta}(t)$ holds and cases where it does not. To simplify, we assume that $n = 2$.

First, if for $i = 1$ or $i = 2$, $v_{\Delta_i}(t) = \texttt{inc}$, then the proposition implies that $v_\Delta(t) = \texttt{inc}$, because $\texttt{inc}$ is maximal with respect to $\lhd$. In this case, the equality always holds. Another case where the equality holds is if $v_{\Delta_1}(t) = \texttt{true}$ and $v_{\Delta_2}(t) = \texttt{false}$. Indeed, in this case we have $\Delta \vdash t$ and $\Delta \not\hspace{-0.3em}\sim t$, showing that $v_\Delta(t) = \texttt{inc}$. Therefore, $v_\Delta(t) = v_{\Delta_1}(t) \overline{\oplus} v_{\Delta_2}(t)$.

To see cases where the equality $v_{\Delta_1}(t) \overline{\oplus} v_{\Delta_2}(t) = v_\Delta(t)$ does not hold, let $U = \{A, B, C\}$, $\Delta_1 = (\{abc\}, \emptyset)$ and $\Delta_2 = (\{bc'\}, \{B \to C\})$.

In this case, $\Delta = (D, \mathcal{FD})$ where $D = \{abc, bc'\}$ and $\mathcal{FD} = \{B \to C\}$. Hence, $D^* = \{abc, abc'\}$ and $\mathsf{Inc}(\Delta) = \{b, bc, bc', abc, abc'\}$, and so:

- $v_{\Delta_1}(b) = v_{\Delta_2}(b) = \texttt{true}$, whereas $v_\Delta(b) = \texttt{inc}$.
- $v_{\Delta_1}(bc') = \texttt{unkn}$, $v_{\Delta_2}(bc') = \texttt{true}$, thus implying that $v_1 \oplus v_2 = \texttt{true}$, whereas $v_\Delta(bc') = \texttt{inc}$.

We further illustrate Proposition 6 in the the context of our introductory example.

*Example 11* We recall that in our introductory example, we have two data sources $\Delta_1 = (D_1, \mathcal{FD})$ and $\Delta_2 = (D_2, \mathcal{FD})$, where $\mathcal{FD} = \{ID \to K, ID \to C\}$.

Based on $D_1$ and $D_2$ as shown in Figure 1 and displayed in Figure 4, applying Algorithm 2 produces $D_1^*$ and $D_2^*$ also shown in Figure 4, and returns $\mathsf{Inc}(\Delta_1) = \mathsf{Inc}(\Delta_2) = \emptyset$.

| $D_1$ | $Id$ | $K$ | $M$ | $C$ |
|---|---|---|---|---|
| | $i_1$ | $k$ | $m$ | $c$ |
| | $i_1$ | | $m'$ | |
| | $i_2$ | $k'$ | $m'$ | $c$ |
| | $i_2$ | $k'$ | $m''$ | |
| | $i_3$ | | $m$ | |

| $D_2$ | $Id$ | $K$ | $M$ | $C$ |
|---|---|---|---|---|
| | $i_1$ | $k$ | | $c$ |
| | $i_2$ | $k'$ | | $c'$ |
| | $i_2$ | $k'$ | $m''$ | |
| | $i_3$ | $k'$ | | |

| $D_1^*$ | $Id$ | $K$ | $M$ | $C$ |
|---|---|---|---|---|
| | $i_1$ | $k$ | $m$ | $c$ |
| | $i_1$ | $k$ | $m'$ | $c$ |
| | $i_2$ | $k'$ | $m'$ | $c$ |
| | $i_2$ | $k'$ | $m''$ | $c$ |
| | $i_3$ | | $m$ | |

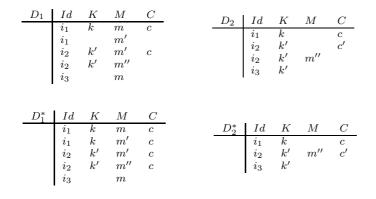| $D_2^*$ | $Id$ | $K$ | $M$ | $C$ |
|---|---|---|---|---|
| | $i_1$ | $k$ | | $c$ |
| | $i_2$ | $k'$ | $m''$ | $c'$ |
| | $i_3$ | $k'$ | | |

**Fig. 4** The source tables of our introductory example and their chased versions

Hence, as already mentioned, $\Delta_1$ and $\Delta_2$ are consistent. Referring to Example 10 and Figure 2, applying Proposition 6 entails the following:

- $v_{\Delta_1}(i_1, k, m, c) = \texttt{true}$, $v_{\Delta_2}(i_1, k, m, c) = \texttt{unkn}$ and $v_\Delta(i_1, k, m, c) = \texttt{true}$.
  $v_{\Delta_1}(i_1, k, m', c) = \texttt{true}$, $v_{\Delta_2}(i_1, k, m', c) = \texttt{unkn}$ and $v_\Delta(i_1, k, m', c) = \texttt{true}$.
  These are cases of equality because $\texttt{true} \overline{\oplus} \texttt{unkn} = \texttt{true}$.
- $v_{\Delta_1}(i_2, c) = \texttt{true}$, $v_{\Delta_2}(i_2, c) = \texttt{false}$ and $v_\Delta(i_2, c) = \texttt{inc}$.
  This is another case of equality because $\texttt{true} \overline{\oplus} \texttt{false} = \texttt{inc}$.
- $v_{\Delta_1}(i_2) = \texttt{true}$, $v_{\Delta_2}(i_2) = \texttt{true}$ and $v_\Delta(i_2) = \texttt{inc}$.
  This is a case where equality does not hold because $\texttt{true} \overline{\oplus} \texttt{true} \neq \texttt{inc}$. Notice however that $\texttt{true} \lhd \texttt{inc}$ holds. □

## 6 Consistent Query Answering

In this section, considering true tuples and false tuples only (*i.e.,* forgetting about false tuples), we address the important problem of *consistent query answering*. We first provide a brief review of the abundant related literature, and then, we show that our approach provides new insights in the problem of consistent query answering. Moreover, we also argue that in our approach, the 'quality' of such consistent answers can be assessed, based on the notion of tuple truth value. However, this issue lies out of the scope of the present paper, and will be the subject of further research in the next future.

### 6.1 Related Work

The problem of query answering in presence of inconsistencies has motivated important research efforts during the past two decades and is still the subject of current research. As mentioned in the introductory section, the most popular approaches in the literature are based on the notion of 'repair', a repair of $\mathcal{D}$ being intuitively *a consistent database $\mathcal{R}$ 'as close as possible' to $\mathcal{D}$*; and an answer to a query $Q$ is consistent if it is present in *every repair $\mathcal{R}$ of $\mathcal{D}$*.

However, it has been recognized that generating *all* repairs is difficult to implement - if not unfeasible. This is a well known problem in practice which explains, for instance, why data cleansing is a very important but tedious task in the management of databases and data warehouses [18]. This issue has been thoroughly investigated in [15], where it has been shown that computing repairs of a given relational table in the presence of functional dependencies is either polynomial or APX-complete[2], depending on the form of the functional dependencies. The reader is referred to [1] for theoretical results on the complexity of testing whether $\mathcal{R}$ is a repair of $\mathcal{D}$, when considering a more generic context than we do in this work (more than one table and constraints other than functional dependencies). A Prolog based approach for the generation of repairs can be found in [4].

Dealing with repairs without generating them is thus an important issue, also known as *Consistent Query Answering in Inconsistent Databases*. One of the first works in this area is [8] and the problem has since been addressed in the context of various database models (mainly the relational model or deductive database models) and under various types of constraints (first order constraints, key constraints, key foreign-key constraints). Seminal papers in this area are [2] and [27], while an overview of works in this area can be found in [6].

The problem considered in all these works can be stated as follows: Given a database $\mathcal{D}$ with integrity constraints $\mathcal{IC}$, assume that $\mathcal{D}$ is inconsistent with respect to $\mathcal{IC}$. Under this assumption, given a query $Q$ against $\mathcal{D}$, what is the *consistent answer* to $Q$? The usual approach to alleviate the impact of inconsistent data on the answers to a query is to consider that an answer to $Q$ is consistent if it is present in *every repair $\mathcal{R}$ of $\mathcal{D}$*.

Complexity results regarding the computation of the consistent answer have been widely studied in [9]. For example one important case is when $\mathcal{IC}$ consists

---

[2] Roughly, APX is the set of NP optimization problems that allow polynomial-time approximation algorithms (source: Wikipedia).

in having one key constraint per database relation and $Q$ is a conjunctive query containing no self-join (*i.e.,* no join of a relation with itself). In this case computing the consistent answer is polynomial whereas if self-joins occur then the problem is co-NP-complete.

Another important problem in considering repairs is that there are many ways of defining the notion of repair. This is so because there are many ways of defining a distance between two database instances, and there is no consensus as to the 'best' definition of distance. Although the distance based on symmetric difference seems to be the most popular, other distances exist as well based for example on sub-sets, on cardinality, on updates or on homomorphism [26]. Notice in this respect that the results in [15] are set for two distances: one based on sub-sets and one based on updates.

6.2 Consistent Query Answering in our Approach

In our work we do *not* use any notion of repair, thus we avoid the above problem of choosing among all possible ways of defining repairs. Instead, we use set theoretic semantics for tuples and functional dependencies that allow us to associate each tuple with one truth value among true, false, inconsistent or unknown.

In what follows, we outline the process of consistent query answering in our approach, and then compare it to the approaches based on repairs. In doing so we follow the intuition of the repairs-approach where an answer to a query is consistent if it is present in every repair; and we transpose it in our approach by considering that a tuple is in the consistent answer to the query if its truth value is `true` in the sense of our model.

As usual when dealing with a single table with nulls, a query $Q$ is an SQL-like expression of one of the following two forms:

$$Q : \mathsf{SELECT}\ X \qquad \text{or} \qquad Q : \mathsf{SELECT}\ X\ \mathsf{WHERE}\ \Gamma$$

In either of these forms, $X$ is an attribute list seen as a relation schema, and in the second form, the WHERE clause specifies a selection condition $\Gamma$. It should thus be clear that, as in SQL, the where clause in a query is optional. The generic form of a query $Q$ is denoted by $Q : \mathsf{SELECT}\ X\ [\mathsf{WHERE}\ \Gamma]$.

A selection condition $\Gamma$ is a well formed formula involving the usual connectors $\neg$, $\vee$ and $\wedge$ and built up from atomic boolean comparisons of one of the forms $A\,\theta\,a$ or $A\,\theta\,A'$, where $\theta$ is a comparison predicate, $A$ and $A'$ are attributes in $U$ whose domain elements are comparable through $\theta$, and $a$ is in $dom(A)$.

Moreover, a tuple $t$ satisfies $A\,\theta\,a$ if $A$ is in $sch(t)$ and if $t.A\,\theta\,a$ holds, and $t$ satisfies $A\,\theta\,A'$ if $A$ and $A'$ are in $sch(t)$ and if $t.A\,\theta\,t.A'$ holds. Based on this, determining whether $t$ satisfies $\Gamma$ follows the rules usual in First Order Logic regarding connectors. For instance, referring to our introductory example, the tuple $t = (k, m)$ such that $sch(t) = KM$ satisfies the conditions $(K = k)$ and $(M = m \vee C = c')$ but does not satisfy the condition $(M = K)$, assuming that $m$ and $k$ are comparable but distinct constants.

Given $\Delta = (D, \mathcal{FD})$, the *answer to $Q$ in $\Delta$* is the set of the restrictions to $X$ of all tuples $t$ in $D^*$ such that $X \subseteq sch(t)$ and such that $t$ satisfies $\Gamma$, when present in $Q$. It follows that answers to queries contain only tuples without nulls.

Now, roughly speaking, the *consistent answer* to $Q$ is the set of all *true* tuples defined over $X$ that satisfy the condition in $Q$. However, as the following example shows, this rough definition should be carefully stated in particular with regard to the functional dependencies to be taken into account for tuple truth value.

*Example 12* In the context of our introductory example, let $\Delta = (D, \mathcal{FD})$ where $\mathcal{FD} = \{Id \rightarrow K, Id \rightarrow C\}$ and where $D$ is displayed in Figure 1. As seen in Example 8, Algorithm 2 returns $D^*$ as shown below and $inc(\mathcal{FD}) = \{inc(Id \rightarrow K), inc(Id \rightarrow C)\}$ where $inc(Id \rightarrow K) = \emptyset$ and $inc(Id \rightarrow C) = \{i_2\}$.

| $D^*$ | $Id$ | $K$ | $M$ | $C$ |
|---|---|---|---|---|
| | $i_1$ | $k$ | $m$ | $c$ |
| | $i_1$ | $k$ | $m'$ | $c$ |
| | $i_2$ | $k'$ | $m'$ | $c$ |
| | $i_2$ | $k'$ | $m''$ | $c$ |
| | $i_2$ | $k'$ | $m'$ | $c'$ |
| | $i_2$ | $k'$ | $m''$ | $c'$ |
| | $i_3$ | $k'$ | $m$ | |

We also recall from Example 10 that $\mathsf{Inc}(\Delta)$ is defined by:

$$\mathsf{Inc}(\Delta) = \quad \{t \mid i_2 \sqsubseteq t \sqsubseteq (i_2, k', m', c)\} \cup \{t \mid i_2 \sqsubseteq t \sqsubseteq (i_2, k', m'', c)\} \cup$$
$$\{t \mid i_2 \sqsubseteq t \sqsubseteq (i_2, k', m', c')\} \cup \{t \mid i_2 \sqsubseteq t \sqsubseteq (i_2, k', m'', c')\}$$

Let $Q_1$ and $Q_2$ be two queries (without conditions) as defined below:

$$Q_1 : \mathsf{SELECT}\ Id, K, C \quad \text{and} \quad Q_2 : \mathsf{SELECT}\ Id, K, M$$

Projecting the tuples in $D^*$ over the attributes $Id$, $K$, $C$ for $Q_1$ and over $Id$, $K$, $M$ for $Q_2$ produces the tables $\Pi_1$ and $\Pi_2$ shown below.

| $\Pi_1$ | $Id$ | $K$ | $C$ |
|---|---|---|---|
| | $i_1$ | $k$ | $c$ |
| | $i_2$ | $k'$ | $c$ |
| | $i_2$ | $k'$ | $c'$ |

| $\Pi_2$ | $Id$ | $K$ | $M$ |
|---|---|---|---|
| | $i_1$ | $k$ | $m$ |
| | $i_1$ | $k$ | $m'$ |
| | $i_2$ | $k'$ | $m'$ |
| | $i_2$ | $k'$ | $m''$ |
| | $i_3$ | $k'$ | $m$ |

Since in these two tables, the tuples whose $Id$-value is $i_2$, are inconsistent in $\Delta$, it seems justified to exclude them from any consistent answer. In other words, according to this intuition, the expected consistent answers to $Q_1$ and $Q_2$ are respectively $\{(i_1, k, c)\}$ and $\{(i_1, k, m), (i_1, k, m'), (i_3, k', m)\}$.

We explain below why it makes sense to exclude the two tuples in the case of $Q_1$, whereas the removal in the case of $Q_2$ is debatable.

1. Regarding $Q_1$, the tuples $(i_2, k', c)$ and $(i_2, k', c')$ in $\Pi_1$ clearly violate $Id \rightarrow C$ from $\mathcal{FD}$, and thus can not occur in the *consistent* answer to $Q_1$.
2. Regarding $Q_2$ however, no functional dependency is violated by the tuples in $\Pi_2$, and thus, there is no reason for removing any of them when producing the *consistent* answer to $Q_2$.

Another way of explaining this situation is to notice that, in $D^*$, the only non satisfied functional dependency is $Id \rightarrow C$ and that

1. attributes $Id$ and $C$ occur in the $\mathsf{SELECT}$ clause of $Q_1$, making it necessary to check functional dependency satisfaction;
2. attribute $C$ does not occur in the $\mathsf{SELECT}$ clause of $Q_2$, implying that checking functional dependency satisfaction makes no sense.

Another important point to take into account is the impact of selection conditions on tuple truth value in the answer to a query. To illustrate this point, first notice that, when considering the query $Q_1$ the only functional dependency to be checked is $Id \rightarrow C$, with respect to which the table $\Pi_1$ shows inconsistencies regarding $i_2$. However, let now $Q_1'$ be the query defined by:

$$Q_1' : \mathsf{SELECT}\ Id, K, C\ \mathsf{WHERE}\ C = c'$$

Only the fifth and sixth tuples in $D^*$ satisfy the selection condition and thus, the only possible tuple in the consistent answer to $Q_1'$ is $(i_2, k', c')$, which alone, trivially satisfies the functional dependency $Id \rightarrow C$.

However, the consistency of the answer to $Q_1'$ may seem counter-intuitive, since the tuple $(i_2, k', c')$ is seen as *inconsistent* in the answer to $Q_1$, where the same attributes are involved. To cope with this counter-intuitive situation, we rather consider that the consistent answer of $Q_1'$ is *empty*, *i.e.,* that consistency has to be checked independently from selection conditions, based only on the functional dependencies involving only attributes from the $\mathsf{SELECT}$ clause in the query.

In what follows, we provide the formalism and the definitions to account for these remarks. □

Given a table $D$ over $U$, a subset $X$ of $U$ and a selection condition $\Gamma$, we denote by $\sigma_\Gamma(D)$, $\pi_X(D)$ and $\pi_X(\mathcal{FD})$ the following sets:

- $\sigma_\Gamma(D)$ is the set of all tuples $t$ in $D$ such that $t$ satisfies $\Gamma$.
- $\pi_X(D)$ is the set of the restrictions to $X$ of all tuples in $D$ whose schema contains $X$; that is $\pi_X(D) = \{t \mid (\exists q \in D)(X \subseteq sch(q),\ t = q.X)\}$.
- $\pi_X(\mathcal{FD})$ is the set of all functional dependencies that involve attributes in $X$ only; that is $\pi_X(\mathcal{FD}) = \{(Y \rightarrow B) \in \mathcal{FD} \mid YB \subseteq X\}$.

These notation are used in the following definition where the notion of *consistent answer to a query* is introduced.

**Definition 6** Given $\Delta = (D, \mathcal{FD})$ and $Q : \mathsf{SELECT}\ X$ [$\mathsf{WHERE}\ \Gamma$], let $\Delta_X$ be defined by $\Delta_X = (\pi_X(D^*), \pi_X(\mathcal{FD}))$.

The *answer to $Q$ in $\Delta$*, denoted by $ans_\Delta(Q)$, is the set $\pi_X(\sigma_\Gamma(D^*))$. Moreover, for every tuple $x$ in $ans_\Delta(Q)$, the *truth value of $x$ in $ans_\Delta(Q)$* is defined by $v_{\Delta_X}(x)$.

The *consistent answer* to $Q$ in $\Delta$, denoted by $ans_\Delta^+(Q)$, is the set of all tuples $x$ in $ans_\Delta(Q)$ such that $v_{\Delta_X}(x) = \mathtt{true}$.

It is important to notice that, according to Definition 6, given $\Delta$ and $Q$, *two distinct* truth values may be given to a tuple $t$, namely, its truth value in $\Delta$, *i.e.,* $v_\Delta(t)$, and its truth value in $\Delta_X$, *i.e.,* $v_{\Delta_X}(t)$. Since these truth values are not determined using the *same* set of functional dependencies, they might be distinct.

Referring to Example 12, based on the notation introduced in Definition 6, for $X_1 = Id\,K\,C$, $\pi_{X_1}(\mathcal{FD}) = \mathcal{FD}$, and so, $\Delta_{X_1} = (\Pi_1, \mathcal{FD})$. In this case, for every tuple $x$ over $X_1$, $v_{\Delta_{X_1}}(x) = v_\Delta(x)$. On the other hand, for $X_2 = Id\,K\,M$, $\pi_{X_2}(\mathcal{FD}) = \{Id \rightarrow K\}$, and so, $\Delta_{X_2} = (\Pi_2, \{Id \rightarrow K\})$. Since $\Pi_2$ satisfies $Id \rightarrow K$, for $x = (i_2, k', m')$, $v_{\Delta_{X_2}}(x) = \mathtt{true}$. However, as $x$ is a super-tuple of $i_2$, we have $v_\Delta(x) = \mathtt{inc}$, showing that $v_{\Delta_{X_2}}(x) \neq v_\Delta(x)$.

The following proposition shows that $ans_\Delta^+(Q)$ is computed from $D^*$ and $inc(\mathcal{FD})$, using Algorithm 5.

---

**Algorithm 5** Consistent answer $ans_\Delta^+(Q)$

---

**Input:** A query $Q$ : SELECT $X$ [WHERE $\Gamma$], $\Delta^* = (D^*, \mathcal{FD})$ and $inc(\mathcal{FD})$
**Output:** The set $ans_\Delta^+(Q)$
1: $ans_\Delta^+(Q) := \emptyset$
2: **for all** $t$ in $D^*$ **do**
3:  **if** $sch(t)$ contains all attributes in $X$ **then**
4:    **if** for every $Y \to B$ in $\pi_X(\mathcal{FD})$, $t.Y$ is not in $inc(Y \to B)$ **then**
5:      **if** $t$ satisfies $\Gamma$ **then**
6:        // This test always succeeds if $Q$ involves no selection condition
7:        $ans_\Delta^+(Q) := ans_\Delta^+(Q) \cup \{t.X\}$
8: **return** $ans_\Delta^+(Q)$

---

**Proposition 7** *Given $\Delta = (D, \mathcal{FD})$ and $Q$ : SELECT $X$ [WHERE $\Gamma$], Algorithm 5 correctly computes $ans_\Delta^+(Q)$.*

*Proof* In this proof, denoting by *ans* the output of Algorithm 5, we prove that $ans = ans_\Delta^+(Q)$. To prove that $ans \subseteq ans_\Delta^+(Q)$, we notice that, by Algorithm 5, every tuple $x$ in *ans* $x$ is the projection over $X$ of a tuple $t$ in $D^*$ satisfying $\Gamma$. Thus, $x$ belongs to $\pi_X(\sigma_\Gamma(D^*))$, that is to $ans_\Delta(Q)$. Moreover, since for every $t$ in $D^*$ such that $t.X$ is in *ans* and every $Y \to B$ in $\pi_X(\mathcal{FD})$, $t.Y$ is not in $inc(Y \to B)$, it holds that $v_{\Delta_X}(x) = \texttt{true}$. It thus follows that $x$ is in $ans_\Delta^+(Q)$.

Conversely, assuming that $x$ is in $ans_\Delta^+(Q)$ implies that $x$ is in $ans_\Delta(Q)$. Hence, $D^*$ contains a tuple $t$ that satisfies $\Gamma$ and $t.X = x$, meaning that $sch(t)$ contains $X$ and that $t$ satisfies the if-condition on line 5 in Algorithm 5. Moreover, since we also have $v_{\Delta_X}(x) = \texttt{true}$, for every $Y \to B$ in $\pi_X(\mathcal{FD})$, $t.Y$ cannot be in $inc(Y \to B)$. This shows that the if-condition on line 4 in Algorithm 5 is satisfied, and thus that $x$ belongs to *ans*, which completes the proof. □

Regarding complexity, Proposition 7 shows that, assuming that $D^*$ has been computed, the computation of the consistent answer is *linear* in the size of $D^*$.

If we assume moreover that $\mathsf{Inc}(\Delta)$ has also been computed, labelling each tuple in $ans_\Delta^+(Q)$ by its truth value in $\Delta$ is an option to investigate, because it has been seen from Definition 6 that the truth value of a tuple $t$ in $\Delta$, *i.e.*, $v_\Delta(t)$, may be different than the truth value of $t$ in $ans_\Delta(Q)$, *i.e.*, $v_{\Delta_X}(t)$.

Knowing that a tuple in the consistent answer, thus having truth value `true` in this answer, has truth value `inc` in the database it comes from, may indeed be relevant in case the user is interested in data quality, as is the case when dealing with data lakes [16]. Investigating further issues related to query answering in our approach, including issues related to data quality is the subject of future work.

*Example 13* Running Algorithm 5 with the queries $Q_1$, $Q_1'$ and $Q_2$ as in Example 12, returns $ans_\Delta^+(Q_1) = \{(i_1, k, c)\}$, $ans_\Delta^+(Q_1') = \emptyset$ and $ans_\Delta^+(Q_2) = \{(i_1, k, m), (i_1, k, m'), (i_2, k', m'), (i_2, k', m''), (i_3, k', m)\}$, as expected.

As earlier noticed regarding $ans_\Delta^+(Q_2)$, for $x = (i_2, k', m')$ or $x = (i_2, k', m'')$, we have $v_\Delta(x) \neq v_{\Delta_{X_2}}(x)$. In this case, smart users could find it relevant to be informed of this situation, which can be done by labelling the two tuples $(i_2, k', m')$ and $(i_2, k', m'')$ by `inc`, that is, their truth value in $\Delta$. We notice that this piece of information cannot be provided by any of the existing approaches.

Considering now the query $Q_3$ : SELECT $M, C$ WHERE $K = k'$, Algorithm 5 discards the first two tuples of $D^*$ (because their $K$-value is not equal to $k'$), and

also the last tuple of $D^*$ (as this tuple has no $C$-value). When processing the remaining four tuples in $D^*$, no functional dependency has to be taken care of, and so, we obtain $ans_\Delta^+(Q_3) = \{(m', c), (m', c'), (m'', c), (m'', c')\}$. □


6.3 Comparison with Repair-Based Approaches

Comparing our approach with approaches to consistent query answering from the literature, we point out that when constraints are functional dependencies only, as in our approach, repairs are defined using set-theoretic inclusion as follows.

**Definition 7** Given $\Delta = (D, \mathcal{FD})$, denoting by $D^*$ the chased table associated with $D$, a *repair* of $\Delta$ is a table $R$ over $U$ such that: (1) $R \subseteq D^*$, (2) $R$ satisfies $\mathcal{FD}$, and (3) $R$ is maximal among the sets satisfying (1) and (2).

We notice that in the above definition, inclusion is understood in its strict set-theoretic meaning, disregarding the presence of nulls in the tuples. For example $\{ab, a'bc\} \subseteq \{abc, a'bc\}$ does not hold whereas $\{ab, a'bc\} \subseteq \{ab, abc, a'bc\}$ does.

Repairs of $\Delta$ can be generated based on the tuples stored in $inc(\mathcal{FD})$ according to the following algorithm:

> $R := D^*$
> **for all** $X \to A$ in $\mathcal{FD}$ **do**
>     **for all** $x$ in $inc(X \to A)$ **do**
>         choose an $A$-value $a$ among all $\alpha$ such that $x\alpha$ occurs in $D^*$
>         $R := R \setminus \{q \mid XA \subseteq sch(q), q.X = x, q.A \neq a\}$
> **return** $R$

Indeed, based on Definition 7, $R$ as computed above is a repair because: (1) $R \subseteq D^*$ clearly holds, (2) $R$ satisfies $\mathcal{FD}$ holds since for every $X \to A$ in $\mathcal{FD}$, there exist $q$ and $q'$ in $R$ such that $q.X = q'.X$ and $q.A \neq q'.A$, and (3) $R$ is maximal because inserting any of the removed tuples leads to violation of a functional dependency.

Given a query $Q$ : SELECT $X$ [WHERE $\Gamma$], denoting by $Rep(\Delta)$ the set of all repairs of $\Delta$, the *consistent answer to $Q$ based on repairs* can be formally defined in the following two ways:

1. $ans_\Delta^\downarrow(Q) = \pi_X \left( \bigcap_{R \in Rep(\Delta)} \sigma_\Gamma(R) \right)$.
2. $ans_\Delta^\uparrow(Q) = \bigcap_{R \in Rep(\Delta)} \pi_X(\sigma_\Gamma(R))$.

Intuitively, $ans_\Delta^\downarrow(Q)$ is obtained by evaluating the query against the intersection of all repairs, whereas $ans_\Delta^\uparrow(Q)$ is obtained by evaluating the query against each repair and by taking the intersection of all these answers.

*Example 14* Computing the repairs of $D^*$ as shown in Example 12 produces the tables $R_1$ and $R_2$ shown below.

| $R_1$ | $Id$ | $K$ | $M$ | $C$ |
|---|---|---|---|---|
| | $i_1$ | $k$ | $m$ | $c$ |
| | $i_1$ | $k$ | $m'$ | $c$ |
| | $i_2$ | $k'$ | $m'$ | $c$ |
| | $i_2$ | $k'$ | $m''$ | $c$ |
| | $i_3$ | $k'$ | $m$ | |

| $R_2$ | $Id$ | $K$ | $M$ | $C$ |
|---|---|---|---|---|
| | $i_1$ | $k$ | $m$ | $c$ |
| | $i_1$ | $k$ | $m'$ | $c$ |
| | $i_2$ | $k'$ | $m'$ | $c'$ |
| | $i_2$ | $k'$ | $m''$ | $c'$ |
| | $i_3$ | $k'$ | $m$ | |

Thus, regarding the queries $Q_1$, $Q_1'$, $Q_2$ and $Q_3$ of Example 12, we have:

---

**Algorithm 6** Repair-based consistent answers $ans^{\downarrow}_{\Delta}(Q)$, $ans^{\uparrow}_{\Delta}(Q)$

---

**Input:** A query $Q$ : SELECT $X$ [WHERE $\Gamma$], $\Delta^* = (D^*, \mathcal{FD})$ and $inc(\mathcal{FD})$
**Output:** The sets $ans^{\downarrow}(Q)$ and $ans^{\uparrow}(Q)$
1: $ans^{\downarrow}(Q) := \emptyset$ ; $ans^{\uparrow}(Q) := \emptyset$
2: **for all** $t$ in $D^*$ **do**
3:     **if** $sch(t)$ contains all attributes in $X$ **then**
4:         **if** $t$ satisfies $\Gamma$ **then**
5:             // This test always succeeds if $Q$ involves no selection condition
6:             **if** for every $Y \to B$ in $\mathcal{FD}$ such that $YB \subseteq sch(t)$, $t.Y$ is not in $inc(Y \to B)$
            **then**
7:                 $ans^{\downarrow}(Q) := ans^{\downarrow}(Q) \cup \{t.X\}$
8:             **if** for every $Y \to B$ in $\mathcal{FD}$ such that $YB \subseteq sch(t)$ and $B \in X$, $t.Y$ is not in
            $inc(Y \to B)$ **then**
9:                 $ans^{\uparrow}(Q) := ans^{\uparrow}(Q) \cup \{t.X\}$
10: **return** $ans^{\downarrow}(Q)$, $ans^{\uparrow}(Q)$

---

- $ans^{\downarrow}_{\Delta}(Q_1) = \{(i_1, k, c)\}$ ; $ans^{\uparrow}_{\Delta}(Q_1) = \{(i_1, k, c)\}$
- $ans^{\downarrow}_{\Delta}(Q'_1) = \emptyset$ ; $ans^{\uparrow}_{\Delta}(Q'_1) = \emptyset$
- $ans^{\downarrow}_{\Delta}(Q_2) = \{(i_1, k, m), (i_1, k, m'), (i_3, k', m)\}$ ;
  $ans^{\uparrow}_{\Delta}(Q_2) = \{(i_1, k, m), (i_1, k, m'), (i_2, k', m'), (i_2, k', m''), (i_3, k', m)\}$
- $ans^{\downarrow}_{\Delta}(Q_3) = \emptyset$ ; $ans^{\uparrow}_{\Delta}(Q_3) = \emptyset$

It should be noticed that computing all repairs before computing the answers is not realistic in practice. In what follows, we provide an efficient algorithm to compute these answers and we prove that they are always 'smaller' with respect to set theoretic inclusion than the answers as defined in Definition 6. $\quad\square$

The following proposition deals with the computation of $ans^{\downarrow}_{\Delta}(Q)$ and of $ans^{\uparrow}_{\Delta}(Q)$, and compares these answers with $ans^{+}_{\Delta}(Q)$.

**Proposition 8** *Given $\Delta = (D, \mathcal{FD})$ and a query $Q$ : SELECT $X$ [WHERE $\Gamma$], Algorithm 6 correctly computes $ans^{\downarrow}_{\Delta}(Q)$ and $ans^{\uparrow}_{\Delta}(Q)$. Moreover, the following holds: $ans^{\downarrow}_{\Delta}(Q) \subseteq ans^{\uparrow}_{\Delta}(Q) \subseteq ans^{+}_{\Delta}(Q)$.*

*Proof* See Appendix F. $\quad\square$

To illustrate the inclusions in Proposition 8, it can be seen from Example 13 and Example 14 that:

- $ans^{\downarrow}_{\Delta}(Q_1) = ans^{\uparrow}_{\Delta}(Q_1) = ans^{+}_{\Delta}(Q_1)$;
- $ans^{\downarrow}_{\Delta}(Q'_1) = ans^{\uparrow}_{\Delta}(Q'_1) = ans^{+}_{\Delta}(Q'_1)$;
- $ans^{\downarrow}_{\Delta}(Q_2) \subset ans^{\uparrow}_{\Delta}(Q_2)$ and $ans^{\uparrow}_{\Delta}(Q_2) = ans^{+}_{\Delta}(Q_2)$;
- $ans^{\downarrow}_{\Delta}(Q_3) = ans^{\uparrow}_{\Delta}(Q_3)$ and $ans^{\uparrow}_{\Delta}(Q_3) \subset ans^{+}_{\Delta}(Q_3)$.

Regarding complexity, is important to note that, if the chased table $D^*$ is available then any of the three ways to compute consistent query answers is *linear* in the size of $D^*$. Moreover, when providing any of these consistent answers, our approach allows for pointing to the user possible problematic tuples, namely those tuples that are *inconsistent in $\Delta$, although not inconsistent in the answer.*

30

## 7 Concluding Remarks

In this paper we have introduced a novel approach to handle inconsistencies in a table with nulls and functional dependencies. Our approach uses set theoretic semantics and relies on an extended version of the well known chase procedure to associate every possible tuple with one of the four truth values true, false, inconsistent and unknown. Moreover, we have seen that true and inconsistent tuples can be computed in time polynomial in the size of the input table. We have also seen that our approach applies to consistent query answering and we have shown that it provides larger answers than the repair-based approaches.

Building upon these results, we currently pursue four lines of research: (1) applying our approach to the particular but important case of key-foreign key constraints in the context of a star schema or a snow-flake schema; (2) designing incremental algorithms to improve performance in case of updates, (3) extending our approach to constraints other than functional dependencies, such as inclusion dependencies as done in [7], (4) investigating the issue of data quality in the framework of our approach, and (5) extending our approach to account for the presence of tuples declared as *false*.

### Declarations

### References

1. Foto N. Afrati and Phokion G. Kolaitis. Repair checking in inconsistent databases: algorithms and complexity. In Ronald Fagin, editor, *Database Theory - ICDT 2009, 12th International Conference,Proceedings*, volume 361 of *ACM International Conference Proceeding Series*, pages 31–41. ACM, 2009.
2. Marcelo Arenas, Leopoldo E. Bertossi, and Jan Chomicki. Consistent query answers in inconsistent databases. In Victor Vianu and Christos H. Papadimitriou, editors, *Proceedings of the Eighteenth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, Pennsylvania, USA*, pages 68–79. ACM Press, 1999.
3. Ofer Arieli and Arnon Avron. The value of the four values. *Artif. Intell.*, 102(1):97–141, 1998.
4. Ofer Arieli, Marc Denecker, Bert Van Nuffelen, and Maurice Bruynooghe. Computational methods for database repair by signed formulae. *Ann. Math. Artif. Intell.*, 46(1-2):4–37, 2006.
5. Nuel D. Belnap. A useful four-valued logic. In J. Michael Dunn and George Epstein, editors, *Modern Uses of Multiple-Valued Logic*, pages 5–37", isbn="978–94–010–1161–7, Dordrecht, 1977. Springer Netherlands.
6. Leopoldo E. Bertossi. *Database Repairing and Consistent Query Answering*. Synthesis Lectures on Data Management. Morgan & Claypool Publishers, 2011.

7. Loreto Bravo and Leopoldo E. Bertossi. Semantically correct query answers in the presence of null values. In Torsten Grust, Hagen Höpfner, Arantza Illarramendi, Stefan Jablonski, Marco Mesiti, Sascha Müller, Paula-Lavinia Patranjan, Kai-Uwe Sattler, Myra Spiliopoulou, and Jef Wijsen, editors, *Current Trends in Database Technology - EDBT 2006, EDBT 2006 Workshops PhD, DataX, IIDB, IIHA, ICSNW, QLQP, PIM, PaRMA, and Reactivity on the Web, Munich, Germany, March 26-31, 2006, Revised Selected Papers*, volume 4254 of *Lecture Notes in Computer Science*, pages 336–357. Springer, 2006.

8. François Bry. Query answering in information systems with integrity constraints. In Sushil Jajodia, William List, Graeme W. McGregor, and Leon Strous, editors, *Integrity and Internal Control in Information Systems*, volume 109 of *IFIP Conference Proceedings*, pages 113–130. Chapman Hall, 1997.

9. Andrea Calì, Domenico Lembo, and Riccardo Rosati. On the decidability and complexity of query answering over inconsistent and incomplete databases. In Frank Neven, Catriel Beeri, and Tova Milo, editors, *Proceedings of the Twenty-Second ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, June 9-12, 2003, San Diego, CA, USA*, pages 260–271. ACM, 2003.

10. S. Ceri, G. Gottlob, and L. Tanca. *Logic Programming and Databases*. Surveys in Computer Science, Springer Verlag, 1990.

11. Stavros S. Cosmadakis, Paris C. Kanellakis, and Nicolas Spyratos. Partition semantics for relations. *J. Comput. Syst. Sci.*, 33(2):203–233, 1986.

12. Ronald Fagin, Alberto O. Mendelzon, and Jeffrey D. Ullman. A simplified universal relation assumption and its properties. *ACM Trans. Database Syst.*, 7(3):343–360, 1982.

13. Melvin Fitting. Bilattices and the semantics of logic programming. *J. Log. Program.*, 11(1&2):91–116, 1991.

14. Dominique Laurent. 4-valued semantics under the OWA: A deductive database approach. In Giorgos Flouris, Dominique Laurent, Dimitris Plexousakis, Nicolas Spyratos, and Yuzuru Tanaka, editors, *Information Search, Integration, and Personalization - 13th International Workshop, ISIP, Revised Selected Papers*, volume 1197 of *Communications in Computer and Information Science*, pages 101–116. Springer, 2019.

15. Ester Livshits, Benny Kimelfeld, and Sudeepa Roy. Computing optimal repairs for functional dependencies. *ACM Trans. Database Syst.*, 45(1):4:1–4:46, 2020.

16. Cedrine Madera and Anne Laurent. The next information architecture evolution: The data lake wave. In *Proceedings of the 8th International Conference on Management of Digital EcoSystems*, MEDES, pages 174–180, New York, NY, USA, 2016. ACM.

17. Francesco Parisi and John Grant. Inconsistency measures for relational databases. *CoRR*, abs/1904.03403, 2019.

18. Erhard Rahm and Hong Hai Do. Data cleaning: Problems and current approaches. *IEEE Data Eng. Bull.*, 23(4):3–13, 2000.

19. Franck Ravat and Yan Zhao. Data lakes: Trends and perspectives. In Sven Hartmann, Josef Küng, Sharma Chakravarthy, Gabriele Anderst-Kotsis, A Min Tjoa, and Ismail Khalil, editors, *Database and Expert Systems Applications - 30th International Conference, DEXA, Proceedings, Part I*, volume 11706 of *Lecture Notes in Computer Science*, pages 304–313. Springer, 2019.

20. Raymond Reiter. On closed world data bases. In Hervé Gallaire and Jack Minker, editors, *Logic and Data Bases, Symposium on Logic and Data Bases, Centre d'études et de recherches de Toulouse, France, 1977*, Advances in Data Base Theory, pages 55–76, New York, 1977. Plenum Press.

21. Nicolas Spyratos. The partition model: A deductive database model. *ACM Trans. Database Syst.*, 12(1):1–37, 1987.

22. Nicolas Spyratos and Christophe Lécluse. Incorporating functional dependencies in deductive query answering. In *Proceedings of the Third International Conference on Data Engineering, February 3-5, 1987, Los Angeles, California, USA*, pages 658–664. IEEE Computer Society, 1987.

23. Alexis Tsoukiàs. A first order, four-valued, weakly paraconsistent logic and its relation with rough sets semantics. *Foundations of Computing and Decision Sciences*, 27(2):77–96, 2002.

24. Jeffrey D. Ullman. *Principles of Databases and Knowledge-Base Systems*, volume 1-2. Computer Science Press, 1988.

25. Moshe Y. Vardi. The universal-relation data model for logical independence. *IEEE Softw.*, 5(2):80–85, 1988.

26. Jef Wijsen. Database repairing using updates. *ACM Trans. Database Syst.*, 30(3):722–768, 2005.
27. Jef Wijsen. On the consistent rewriting of conjunctive queries under primary key constraints. *Inf. Syst.*, 34(7):578–601, 2009.

# A Proof of Lemma 1

**Lemma 1.** *For every $\Delta = (D, \mathcal{FD})$, the sequence $(\mu_i)_{i \geq 0}$ has a unique limit $\mu^*$ that satisfies that $\mu^* \models \Delta$. Moreover:*

1. *For all $a_1$ and $a_2$ in the same attribute domain $dom(A)$, if $\mu^*(a_1) \cap \mu^*(a_2) \neq \emptyset$ then there exist $X \to A$ in $\mathcal{FD}$ and $x$ over $X$ such that $\mu^*(x) \neq \emptyset$ and $\mu^*(x) \subseteq \mu^*(a_1) \cap \mu^*(a_2)$.*
2. *For all $\alpha$ and $\beta$, $\Delta \vdash (\alpha \sqcap \beta)$ holds if and only if $\mu^*(\alpha) \cap \mu^*(\beta) \neq \emptyset$ holds.*

*Proof* We recall that the sequence $(\mu_i)_{i \geq 0}$ is defined by the following steps:

1. For every $t$ in $D$, assign a 'fresh' integer $id(t)$ to $t$;
2. Let $\mu_0$ be the mapping defined for every domain constant $a$ by:
   $\mu_0(a) = \{id(t) \mid t \in D \text{ and } a \sqsubseteq t\}$;
3. While there exists $X \to A$ in $\mathcal{FD}$, $x$ over $X$ and $a$ in $dom(A)$ such that $\mu(xa) \neq \emptyset$ and $\mu(x) \not\subseteq \mu(a)$, define $\mu_{i+1}$ by: $\mu_{i+1}(a) = \mu_i(a) \cup \mu_i(x)$ and $\mu_{i+1}(\alpha) = \mu_i(\alpha)$ for any other constant $\alpha$.

The sequence $(\mu_i)_{i \geq 0}$ is increasing in the sense that for every $\alpha$, $\mu_i(\alpha) \subseteq \mu_{i+1}(\alpha)$, and bounded in the sense that for every $\alpha$, $\mu_i(\alpha) \subseteq \{id(t) \mid t \in \Delta\}$. Hence the sequence has a unique limit. Moreover, for every $t$ in $D$, $\mu^*(t) \neq \emptyset$ holds because $id(t)$ always belongs to $\mu^*(t)$, and $\mu^* \models \mathcal{FD}$, because otherwise $\mu^*$ would not be the limit of the sequence. Therefore $\mu^* \models \Delta$, which shows the first part of the lemma.

(1) Regarding the first item in the second part of the lemma, we first notice that by definition of $\mu_0$, we have $\mu_0(a_1) \cap \mu_0(a_2) = \emptyset$, because it is not possible that a tuple in $D$ has two distinct values over an attribute.

Since we assume that $\mu^*(a_1) \cap \mu^*(a_2) \neq \emptyset$, there exists $i_0 \geq 0$ such that $\mu_{i_0}(a_1) \cap \mu_{i_0}(a_2) = \emptyset$ and $\mu_{i_0+1}(a_1) \cap \mu_{i_0+1}(a_2) \neq \emptyset$. By definition of the sequence $(\mu_i)_{i \geq 0}$, for $j = 1, 2$, $\mu_{i_0+1}(a_j) = \mu_{i_0}(a_j) \cup M(a_j)$ where $M(a_j)$ is the union of all $\mu_{i_0}(x_j)$ such that $X_j \to A$ is in $\mathcal{FD}$, $\mu_{i_0}(x_j) \cap \mu_{i_0}(a_j) \neq \emptyset$ and $\mu_{i_0}(x_j) \not\subseteq \mu_{i_0}(a_j)$. Hence,

$$\begin{aligned}
\mu_{i_0+1}(a_1) \cap \mu_{i_0+1}(a_2) &= (\mu_{i_0}(a_1) \cup M(a_1)) \cap (\mu_{i_0}(a_2) \cup M(a_2)) \\
&= (\mu_{i_0}(a_1) \cap \mu_{i_0}(a_2)) \cup (\mu_{i_0}(a_1) \cap M(a_2)) \cup \\
&\quad (M(a_1) \cap \mu_{i_0}(a_2)) \cup (M(a_1) \cap M(a_2))
\end{aligned}$$

Since $\mu_{i_0+1}(a_1) \cap \mu_{i_0+1}(a_2) \neq \emptyset$, at least one of the four terms of the above union is not empty. But since $\mu_{i_0}(a_1) \cap \mu_{i_0}(a_2) = \emptyset$, only the last three cases are investigated below.

($i$) If $\mu_{i_0}(a_1) \cap M(a_2) \neq \emptyset$, $M(a_2)$ contains $x_2$ such that $\mu_{i_0}(a_1) \cap \mu_{i_0}(x_2) \neq \emptyset$. Thus, there exists $X_2 \to A$ is in $\mathcal{FD}$ such that $X_2 = sch(x_2)$, $\mu_{i_0}(a_1) \cap \mu_{i_0}(x_2) \neq \emptyset$ and $\mu_{i_0}(a_2) \cap \mu_{i_0}(x_2) \neq \emptyset$. Since both $a_1$ and $a_2$ are in $dom(A)$, we have $\mu_{i_0+1}(x_2) \subseteq \mu_{i_0+1}(a_1)$ and $\mu_{i_0+1}(x_2) \subseteq \mu_{i_0+1}(a_2)$. Thus $\mu^*(x_2) \subseteq \mu^*(a_1) \cap \mu^*(a_2)$.

($ii$) If $\mu_{i_0}(a_2) \cap M(a_1) \neq \emptyset$, it can be shown in a similar way that there exist $X_1 \to A$ is in $\mathcal{FD}$ and $x_1$ over $X_1$ such that $\mu^*(x_1) \subseteq \mu^*(a_1) \cap \mu^*(a_2)$. The proof is omitted.

($iii$) If $M(a_1) \cap M(a_2) \neq \emptyset$, for $j = 1, 2$, $M(a_j)$ contains $x_j$ such that $\mu_{i_0}(x_1) \cap \mu_{i_0}(x_2) \neq \emptyset$. Thus, for $j = 1, 2$, there exist $X_j \to A$ in $\mathcal{FD}$ such that $X_j = sch(x_j)$, $\mu_{i_0}(x_j) \cap \mu_{i_0}(a_j) \neq \emptyset$ and $\mu_{i_0}(x_1) \cap \mu_{i_0}(x_2) \neq \emptyset$. Hence, $\mu_{i_0+1}(x_j) \subseteq \mu_{i_0+1}(a_j)$, for $j = 1, 2$ and $\mu_{i_0+1}(x_1) \cap \mu_{i_0+1}(x_2) \neq \emptyset$. It follows that, when computing $\mu_{i_0+2}$, we obtain the additional inclusions $\mu_{i_0+2}(x_1) \subseteq \mu_{i_0+2}(a_2)$ and $\mu_{i_0+2}(x_2) \subseteq \mu_{i_0+2}(a_1)$, which implies that for $j = 1, 2$, $\mu^*(x_j) \subseteq \mu^*(a_1) \cap \mu^*(a_2)$ holds. This part of the proof is thus complete.

(2) Regarding the second item in the second part of the lemma, assume first that $\Delta \vdash (\alpha \sqcap \beta)$. Since $\mu^* \models \Delta$, we obviously have that $\mu^*(\alpha) \cap \mu^*(\beta) \neq \emptyset$.

Conversely, assuming that $\mu^*(\alpha) \cap \mu^*(\beta) \neq \emptyset$, we show that $\Delta \vdash (\alpha \sqcap \beta)$, that is, for every $\mu$ such that $\mu \models \Delta$, $\mu(\alpha) \cap \mu(\beta) \neq \emptyset$. The proof is by induction on the steps of the construction of $\mu^*$, assuming $\alpha$ in $dom(A)$ and $\beta$ in $dom(B)$.

• The result holds for $i = 0$. Indeed, if $\mu_0(\alpha) \cap \mu_0(\beta) \neq \emptyset$ then there exists $u$ in $D$ such that

$\alpha \sqsubseteq u$ and $\beta \sqsubseteq u$. Hence for every $\mu$ such that $\mu \models \Delta$, we have $\mu(u) \neq \emptyset$ and $\mu(u) \subseteq \mu(\alpha) \cap \mu(\beta)$, implying that $\mu(\alpha) \cap \mu(\beta) \neq \emptyset$ holds.

• For $i_0 > 0$, assuming that $\mu_{i_0}$ satisfies that for all $\zeta$ and $\eta$ such that $\mu_{i_0}(\zeta) \cap \mu_{i_0}(\eta) \neq \emptyset$, we have $\mu(\zeta) \cap \mu(\eta) \neq \emptyset$ for every $\mu$ such that $\mu \models \Delta$, we show that the result holds for $\mu_{i_0+1}$.

Indeed, let $i_0$ such that $\mu_{i_0}(\alpha) \cap \mu_{i_0}(\beta) = \emptyset$ and $\mu_{i_0+1}(\alpha) \cap \mu_{i_0+1}(\beta) \neq \emptyset$. By definition of the sequence $(\mu_i)_{i \geq 0}$, and as in (1) just above, $\mu_{i_0+1}(\alpha) = \mu_{i_0}(\alpha) \cup M(\alpha)$ where $M(\alpha)$ is the union of all $\mu_{i_0}(x)$ such that $X \to A$ is in $\mathcal{FD}$, $\mu_{i_0}(x) \cap \mu_{i_0}(\alpha) \neq \emptyset$ and $\mu_{i_0}(x) \not\subseteq \mu_{i_0}(\alpha)$. Similarly, $\mu_{i_0+1}(\beta) = \mu_{i_0}(\beta) \cup M(\beta)$ where $M(\beta)$ is the union of all $\mu_{i_0}(y)$ such that $Y \to B$ is in $\mathcal{FD}$, $\mu_{i_0}(y) \cap \mu_{i_0}(\beta) \neq \emptyset$ and $\mu_{i_0}(y) \not\subseteq \mu_{i_0}(\beta)$. Thus:

$$
\begin{aligned}
\mu_{i_0+1}(\alpha) \cap \mu_{i_0+1}(\beta) \quad &= (\mu_{i_0}(\alpha) \cup M(\alpha)) \cap (\mu_{i_0}(\beta) \cup M(\beta)) \\
&= (\mu_{i_0}(\alpha) \cap \mu_{i_0}(\beta)) \cup (\mu_{i_0}(\alpha) \cap M(\beta)) \ \cup \\
&\qquad (M(\alpha) \cap \mu_{i_0}(\beta)) \cup (M(\alpha) \cap M(\beta))
\end{aligned}
$$

Since $\mu_{i_0+1}(\alpha) \cap \mu_{i_0+1}(\beta) \neq \emptyset$, at least one of the four terms of the above union is non empty. But since $\mu_{i_0}(\alpha) \cap \mu_{i_0}(\beta) = \emptyset$, only the last three cases are investigated below.

($i$) If $\mu_{i_0}(\alpha) \cap M(\beta) \neq \emptyset$, there exist $Y \to B$ in $\mathcal{FD}$ and $y$ over $Y$ such that $\mu_{i_0}(\alpha) \cap \mu_{i_0}(y) \neq \emptyset$ and $\mu_{i_0}(\beta) \cap \mu_{i_0}(y) \neq \emptyset$. By our induction hypothesis, for every $\mu$ such that $\mu \models \Delta$, we have $\mu(\alpha) \cap \mu(y) \neq \emptyset$ and $\mu(y) \subseteq \mu(\beta)$, which implies that $\mu(\alpha) \cap \mu(\beta) \neq \emptyset$.

($ii$) If $\mu_{i_0}(\beta) \cap M(\alpha) \neq \emptyset$, the case is similar to ($i$) above. The proof is omitted.

($iii$) If $M(\alpha) \cap M(\beta) \neq \emptyset$, there exist $X \to A$ and $Y \to B$ in $\mathcal{FD}$, $x$ over $X$ and $y$ over $Y$, such that $\mu_{i_0}(x) \cap \mu_{i_0}(y) \neq \emptyset$, $\mu_{i_0}(\alpha) \cap \mu_{i_0}(x) \neq \emptyset$ and $\mu_{i_0}(\beta) \cap \mu_{i_0}(y) \neq \emptyset$. By our induction hypothesis, for every $\mu$ such that $\mu \models \Delta$, we have $\mu(x) \cap \mu(y) \neq \emptyset$, $\mu(x) \subseteq \mu(\alpha)$ and $\mu(y) \subseteq \mu(\beta)$. Hence, $\mu(\alpha) \cap \mu(\beta) \neq \emptyset$ also holds in this case, and the proof is complete. □


## B Proof of Lemma 2

**Lemma 2.** *Let $\Delta = (D, \mathcal{FD})$ and $t$ a tuple. Then Algorithm 1 computes correctly the closure $t^+$ of $t$.*

*Proof* In this proof, we denote by $cl(t)$ the output of Algorithm 1, and we show that $cl(t) = t^+$, that is that $cl(t) \subseteq t^+$ and $t^+ \subseteq cl(t)$ both hold. Before proceeding to these proofs, we draw attention on that for every $\mathcal{T}$-mapping $\mu$ such that $\mu(t) \neq \emptyset$, $\mu \models \Delta$ if and only if $\mu \models \Delta_t$, where $\Delta_t$ is defined by the statement line 1 in Algorithm 1. Indeed:

• If $\mu \models \Delta_t$ then for every $q \in D_t$ $\mu(q) \neq \emptyset$ and $\mu \models \mathcal{FD}$. Since $D \subset D_t$, $\mu(q) \neq \emptyset$ for every $q \in D$, implying that $\mu \models \Delta$ holds.

• Conversely, if $\mu \models \Delta$ then as $\mu(t)$ is supposed to be nonempty, $\mu(q) \neq \emptyset$ for every $q$ in $D_t$. Since $\mu \models \mathcal{FD}$ holds, $\mu \models \Delta_t$ also holds.

To first prove that $cl(t) \subseteq t^+$, we consider a $\mathcal{T}$-mapping $\mu$ such that $\mu \models \Delta$, and we prove that $\mu(t) \subseteq \mu(a)$ for every $a$ in $cl(t)$. We first observe that if $\mu(t) = \emptyset$ then $\mu(t) \subseteq \mu(\alpha)$ holds for every constant $\alpha$. Therefore, $\mu(t) \subseteq \mu(a)$ holds.

Now, if $\mu(t) \neq \emptyset$ then $\mu \models \Delta_t$, as shown above. The proof that $\mu(t) \subseteq \mu(a)$ is done by induction on the steps of the execution of Algorithm 1. Denoting by $cl^0$, $cl^1$, ... the sequence of the assignments of $cl(t)$ during execution, the following holds for every $a$ in $cl(t)$.

• If $a$ is in $cl^0$ as computed on line 2, $a$ occurs in $t$. It is thus clear that $\mu(t) \subseteq \mu(a)$.

• We now assume that, for $j \geq 0$, every $\alpha$ in $cl^j$ is such that $\mu(t) \subseteq \mu(\alpha)$ and we show that this holds for $a$ in $cl^{j+1}$ but not in $cl^j$. In this case, according to the condition in line 5 of Algorithm 1, there exist $X \to A$ in $\mathcal{FD}$ and $x$ over $X$ such that $\Delta_t \vdash xa$ and for every $b$ in $x$, $b \in cl^j$. Thus $\mu(x) \cap \mu(a) \neq \emptyset$ (because $\mu \models \Delta_t$) and $\mu(t) \subseteq \mu(b)$ for every $b$ in $x$ (by our induction hypothesis, because $b$ is in $cl^j$). Hence $\mu(t) \subseteq \mu(x)$ and $\mu(x) \subseteq \mu(a)$ hold, thus implying that $\mu(t) \subseteq \mu(a)$.

As a consequence, we have shown that for every $\mu$ such that $\mu \models \Delta$, for every $a$ in $cl(t)$, $\mu(t) \subseteq \mu(a)$. Therefore, by Definition 3, $cl(t) \subseteq t^+$ holds.

Conversely, $t^+ \subseteq cl(t)$ is shown by contraposition: assuming that $a \notin cl(t)$, we prove that $a \notin t^+$. To this end, we exhibit a $\mathcal{T}$-mapping $\mu_t$ such that $\mu_t \models \Delta$ and $\mu_t(t) \not\subseteq \mu_t(a)$.

We denote by $\mu_t^*$ the $\mathcal{T}$-mapping built up as $\mu^*$, but starting from $\Delta_t$ as defined line 1 in Algorithm 1. Thus, $\mu_t^* \models \Delta_t$, and since $\mu_t^*(t) \neq \emptyset$, it has been seen above that $\mu_t^* \models \Delta$.

Thus, if $\mu_t^*(t) \not\subseteq \mu_t^*(a)$ then $\mu_t^*$ is the $\mathcal{T}$-mapping we are looking for, and thus, we set $\mu_t = \mu_t^*$. Assuming that $\mu^*(t) \subseteq \mu^*(a)$, let $k$ be an integer not in $\mu_t^*(\alpha)$ for any $\alpha$ occurring in

$\Delta_t$, and let $\mu_t$ be the $\mathcal{T}$-mapping defined for every constant $\alpha$ by:
- $\mu_t(\alpha) = \mu_t^*(\alpha) \cup \{k\}$, if $\alpha \in cl(t)$
- $\mu_t(\alpha) = \mu_t^*(\alpha)$, otherwise.
We show that $\mu_t$ satisfies that: (1) $\mu_t(t) \not\subseteq \mu_t(a)$ and (2) $\mu_t \models \Delta$.

(1) Since every $\alpha$ in $t$ is in $cl(t)$, $k$ is in $\mu_t(t)$ and since $a$ is not in $cl(t)$, $k$ is not in $cl(a)$. It thus follows that $\mu_t(t) \not\subseteq \mu_t(a)$.

(2) Since for every constant $\alpha$, $\mu_t^*(\alpha) \subseteq \mu_t(\alpha)$ holds, for every $q$ in $D$, it holds that $\mu_t^*(q) \subseteq \mu_t(q)$, which implies $\mu_t(q) \neq \emptyset$, because $\mu_t^*(q) \neq \emptyset$ holds as a consequence of $\mu_t^* \models \Delta$.

To prove that $\mu_t \models Y \to B$ for every $Y \to B$ in $\mathcal{FD}$, let $y$ over $Y$ and $b$ in $dom(B)$ such $\mu_t(y) \cap \mu_t(b) \neq \emptyset$. To show that $\mu_t(y) \subseteq \mu_t(b)$, we consider the two cases according to which $\mu_t^*(y) \cap \mu_t^*(b)$ is or not empty.
- If $\mu_t^*(y) \cap \mu_t^*(b) = \emptyset$, then by definition of $\mu_t$, for $\mu_t(y) \cap \mu_t(b)$ to be nonempty, it must be that $\mu_t(y) = \mu_t^*(y) \cup \{k\}$ and $\mu_t(b) = \mu_t^*(b) \cup \{k\}$. Writing $y$ as $\beta_1 \dots \beta_p$, this implies that every $\beta_i$ $(i = 1, \dots, p)$, and $b$ are in $cl(t)$. Then, as we know that $cl(t) \subseteq t^+$ holds, all these constants are in $t^+$, implying that $\mu_t^*(t) \subseteq \mu_t^*(\beta_i)$ $(i = 1, \dots, p)$ and $\mu_t^*(t) \subseteq \mu_t^*(b)$, because $\mu_t^* \models \Delta$. Since $\mu_t^*(t) \neq \emptyset$, we have $\mu_t^*(y) \cap \mu_t^*(b) \neq \emptyset$, which contradicts our hypothesis that $\mu_t^*(y) \cap \mu_t^*(b) = \emptyset$. This case is thus not possible.
- If $\mu_t^*(y) \cap \mu_t^*(b) \neq \emptyset$, then as $\mu_t^* \models \mathcal{FD}$, $\mu_t^*(y) \subseteq \mu_t^*(b)$ holds, and by Lemma 1 applied to $\Delta_t$, we also have that $\Delta_t \vdash yb$. Since $\mu_t^*(y) \subseteq \mu_t^*(b)$ holds, assuming that $\mu_t(y) \subseteq \mu_t(b)$ does not hold implies that $k$ belongs to $\mu_t(y)$ but not to $\mu_t(b)$. Hence, every $\beta_i$ $(i = 1, \dots, p)$ is in $cl(t)$ whereas $b$ is not. This is a contradiction with line 5 of Algorithm 1, where it is stated that $\beta$ is inserted into $cl(t)$ (because $\Delta_t \vdash yb$ and every $\beta_i$ $(i = 1, \dots, p)$ is in $cl(t)$). Thus, $\mu_t(y) \subseteq \mu_t(b)$ holds showing that $\Delta_t \models Y \to B$. The proof is therefore complete. $\qquad\square$


## C Proof of Lemma 3

**Lemma 3.** *Algorithm 2 applied to $\Delta = (D, \mathcal{FD})$ always terminates. Moreover, for every tuple $t$, $\mu^*(t) \neq \emptyset$ holds if and only if $t$ is in $\mathsf{LoCl}(D^*)$.*

*Proof* The tuples inserted into $D^*$ when running the while-loop line 4 of Algorithm 2 are built up using only constants occurring in $\Delta$. Thus, the number of these tuples is finite, and so, Algorithm 2 terminates.

The proof that for every $t$ in $\mathsf{LoCl}(D^*)$, $\mu^*(t) \neq \emptyset$ holds is conducted by induction on the steps of Algorithm 2. If $(D_k)_{k \geq 0}$ denotes the sequence of the states of $D^*$ during the execution, we first note that since $D_0 = D$, for every $t$ in $\mathsf{LoCl}(D_0)$, $\mu^*(t) \neq \emptyset$ holds.

Assuming now that for $i > 0$, for every $t$ in $\mathsf{LoCl}(D_i)$, $\mu^*(t) \neq \emptyset$, we prove the result for every $t$ in $\mathsf{LoCl}(D_{i+1})$. Indeed, let $t'$ in $D_{i+1}$ such that $t \sqsubseteq t'$. If $t'$ is in $D_i$, the proof is immediate; we thus now assume that $t'$ is not in $D_i$, that is that $t'$ occurs in $D_{i+1}$ when running Algorithm 2, that is, there exist $X \to A$ in $\mathcal{FD}$, $t_1$ and $t_2$ in $D_i$ such that $t_1.X = t_2.X = x$, $t_1.A = a$ and either (i) $t_2.A$ is not defined or (ii) $t_2.A$ is defined but not equal to $t_1.A$. Writing $t_1$ as $t_1'xa$, we have the following:

(i) If $t_2.A$ is not defined, then $t_2$ is written as $t_2'x$ and, according to the statement line 9, $t'$ is of the form $t_2'xa$. By our induction hypothesis, $\mu^*(t_1)$ and $\mu^*(t_2)$ are nonempty, and thus $\mu^*(x) \cap \mu^*(a) \neq \emptyset$. Hence, $\mu^*(x) \subseteq \mu^*(a)$ (because $\mu^* \models X \to A$), and so, $\mu^*(t') = \mu^*(t_2') \cap \mu^*(x) \cap \mu^*(a) = \mu^*(t_2') \cap \mu^*(x)$, showing that $\mu^*(t') = \mu^*(t_2)$. Hence $\mu^*(t') \neq \emptyset$, and so, $\mu^*(t) \neq \emptyset$ also holds, since $\mu^*(t') \subseteq \mu^*(t)$.
(ii) If $t_2.A$ is defined but $t_1.A \neq t_2.A$. for $i = 1, 2$, $t_i$ is written as $t_i'xa_i$ where $a_i = t_i.A$. statement line 12, $t'$ is one of the tuples $t_1'xa_2$ or $t_2'xa_1$, and each of these cases can be treated as in (i) above,

We therefore have shown that if $t$ is in $\mathsf{LoCl}(D^*)$ as computed by the main loop line 4 of Algorithm 2, then $\mu^*(t) \neq \emptyset$. Since the last loop line 14 does not change this set $\mathsf{LoCl}(D^*)$, this part of the proof is complete.

Conversely, we show that for every $t$, if $\mu^*(t) \neq \emptyset$ then $t$ is in $\mathsf{LoCl}(D^*)$. The proof is done by induction on the construction of $\mu^*$. By definition of $\mu_0$, it is clear that if $\mu_0(t) \neq \emptyset$ then $t$ is in $\mathsf{LoCl}(D)$ and thus in $\mathsf{LoCl}(D^*)$. Now, if we assume that for every $i > 0$ and every $t$, if $\mu_i(t) \neq \emptyset$ then $t$ belongs to $\mathsf{LoCl}(D^*)$, we prove that this result holds for $\mu_{i+1}$.

Let $t$ be such that $\mu_i(t) = \emptyset$ and $\mu_{i+1}(t) \neq \emptyset$. For every $\alpha$, writing $\mu_{i+1}(\alpha)$ as $\mu_i(\alpha) \cup M(\alpha)$, where $M(\alpha)$ is the union of all $\mu_i(x)$ such that $x$ is a tuple over $X$, where $X \to A \in \mathcal{FD}$, $\alpha \in dom(A)$, $\mu_i(x) \cap \mu_i(\alpha) \neq \emptyset$, and $\mu_i(x) \not\sqsubseteq \mu_i(\alpha)$, we have the following:

$$
\begin{aligned}
\mu_{i+1}(t) \quad &= \bigcap_{\alpha \sqsubseteq t} \mu_{i+1}(\alpha) \\
&= \bigcap_{\alpha \sqsubseteq t} \left( \mu_i(\alpha) \cup M(\alpha) \right) \qquad\qquad\qquad\qquad (1) \\
&= \mu_i(t) \cup \left( \bigcup_{t = t_1 t_2} \left( \mu_i(t_1) \cap \left( \bigcap_{\beta \sqsubseteq t_2} M(\beta) \right) \right) \right) \cup \left( \bigcap_{\alpha \sqsubseteq t} M(\alpha) \right) \qquad (2)
\end{aligned}
$$

Equality (2) above is obtained from (1) by applying the distributivity of intersection over union with the convention that $t = t_1 t_2$ refers to any split of $t$ into two tuples $t_1$ and $t_2$. Assuming $\mu_i(t) = \emptyset$ and $\mu_{i+1}(t) \neq \emptyset$ implies that in Equality (2) either the second or the last term of the union is nonempty.

• If $\bigcup_{t = t_1 t_2} \left( \mu_i(t_1) \cap \left( \bigcap_{\beta \sqsubseteq t_2} M(\beta) \right) \right) \neq \emptyset$, there exist $t_1$ and $t_2$ such that $t = t_1 t_2$ and $\mu_i(t_1) \cap \left( \bigcap_{\beta \sqsubseteq t_2} M(\beta) \right) \neq \emptyset$. Given such a split of $t$, writing $t_2$ as $\beta_1 \ldots \beta_p$ implies that, for $k = 1, \ldots, p$, $M(\beta_k)$ contains $y_k$ such that $Y_k \to B_k$ is in $\mathcal{FD}$ and $\mu_i(y_k) \cap \mu_i(\beta_k) \neq \emptyset$. Moreover, we have that $\mu_i(t_1) \cap \left( \bigcap_{k=1}^{k=p} \mu_i(y_k) \right) \neq \emptyset$. Thus by our induction hypothesis, $\mathsf{LoCl}(D^*)$ contains a tuple of the form $q_1 t_1 y_1 \ldots y_p$ and $p$ tuples of the form $q'_k y_k \beta_k$ $(k = 1, \ldots, p)$.

Now, given $k = 1, \ldots, p$, if $q_1 t_1 y_1 \ldots y_p$ is not defined over $B_k$, $q_1 t_1 y_1 \ldots y_p \beta_k$ appears in $D^*$ due to the statement line 9 of Algorithm 2. Assume now that $q_1 t_1 y_1 \ldots y_p$ is defined over $B_k$ but with a value different than $\beta_k$, say $\beta'_k$.

By construction of $t_1$ and $t_2$, $B_k$ is not in $sch(t_1)$, and so, $B_k$ is either in $sch(q_1)$ or in $Y_i$ for some $i = 1, \ldots, p$. In any case, denoting $sch(q_1 y_1 \ldots y_p)$ by $Q$, we write $q_1 t_1 y_1 \ldots y_p$ as $r^k t_1 b'_k$ where $r^k = (q_1 y_1 \ldots y_p).(Q \setminus B_k)$. Considering that $r^k t_1 b'_k$ and $q'_k y_k \beta_k$ have the same $Y_k$-value $y_k$, the statement line 12 of Algorithm 2 applies and $r^k t_1 \beta_k$ is inserted in $D^*$. During the subsequent iterations, a similar argument shows that $D^*$ contains a tuple of the form $r t_1 \beta_1 \ldots \beta_p$, that is $r t_1 t_2$ or $rt$. It thus follows that $t$ is in $\mathsf{LoCl}(D^*)$.

• If $\bigcap_{\alpha \sqsubseteq t} M(\alpha) \neq \emptyset$, the same reasoning as above applies considering that $t_1$ is empty and $t_2 = t$. After the iterations, $D^*$ contains a tuple of the form $r \beta_1 \ldots \beta_p$, that is $rt$. Thus, in this case again, $t$ is in $\mathsf{LoCl}(D^*)$, and the proof is complete. □


## D Proof of Proposition 3

**Proposition 3.** *Let $\Delta = (D, \mathcal{FD})$ and $t$ be such that $\Delta \vdash t$. For every tuple $q$ and every $a$ in $dom(A)$ such that $q \sqsubseteq t$ and $a \sqsubseteq t$, we have: $a$ belongs to $q^+$ if and only if $A$ belongs to $Q^+$.*

*Proof* Assuming first $a$ in $q^+$, we show by induction on the steps of Algorithm 1 that $A$ is in $Q^+$. It is important to notice that since $q \sqsubseteq t$ and $\Delta \vdash t$, $\Delta \vdash q$ holds. Hence when running Algorithm 1 with $\Delta$ and $q$ as input, as shown in the proof of Lemma 2, $\mu \models \Delta$ holds if and only if $\mu \models \Delta_q$ holds. Thus, for every tuple $\tau$, $\Delta \vdash \tau$ holds if and only if $\Delta_q \vdash \tau$ holds.

If $a$ is in $q^+$ because of line 2 in Algorithm 1, then $A$ is in $Q$, showing that $A$ is in $Q^+$. If $a$ is inserted in $q^+$ because of line 5, then there exist $X \to A$ in $\mathcal{FD}$ and $x$ over $X$ such that every $b$ in $x$ belongs to $q^+$ and $\Delta_q \vdash xa$, that is $\Delta \vdash xa$. Assuming that the proposition holds for every $b$ in $x$ implies that every $B$ in $X$ is in $Q^+$. Thus, $X \subseteq Q^+$ holds, and so $A$ is in $Q^+$.

Conversely, let $A$ be in $Q^+$. If $A$ is in $Q$, then $q.A = a$, and so, $a$ is in $q^+$. Let us now assume that $A$ is not in $Q$, and let us show by induction on the execution of the loop computing $Q^+$ that $a$ belongs to $q^+$. Indeed, denoting by $Q'$ the current value of $Q^+$ when $A$ is inserted in $Q^+$, there exists $X \to A$ in $\mathcal{FD}$ such that $X \subseteq Q'$. Thus, by our induction hypothesis, every $\alpha$ in $q.X$ is in $q^+$. Moreover, since $\Delta \vdash t$ and $xa = t.XA$, $\Delta \vdash xa$. Hence, $\Delta_q \vdash xa$, and by the statement line 5 of Algorithm 1, $a$ belongs to $q^+$. The proof is therefore complete. □


## E Proof of Lemma 4

**Lemma 4.** *Given $\Delta = (D, \mathcal{FD})$, a tuple $t$ is inconsistent in $\Delta$ if and only if $t \in \mathsf{Inc}(\Delta)$.*

*Proof* We note first that for every $x$ in $inc(X \to A)$ there exist $a_1, \ldots, a_k$ $(k \geq 2)$ in $dom(A)$ such that for every $i = 1, \ldots, k$, $xa_i \in \mathsf{LoCl}(D^*)$, thus such that $\Delta \vdash xa_i$. Therefore, for every $i = 1, \ldots, k$, $a_i$ belongs to $x^+$, and so, $\Delta \vdash (x \preceq a_1 \sqcap \ldots \sqcap a_k)$ holds, showing that $x$ is inconsistent in $\Delta$.

We now prove that if $q$ belongs to $\mathsf{Inc}(\Delta)$ then $q$ is inconsistent in $\Delta$. Indeed, by Algorithm 3, there exist $t$ in $D^*$, $X \to A$ in $\mathcal{FD}$, such that $Q \subseteq T$, $t.Q = q$, $t.X \in inc(X \to A)$, and $X \subseteq Q^+$. Since $\Delta \vdash t$, Proposition 3 applies, showing that for every $\alpha$ in $x$, $\alpha$ belongs to $q^+$, where $q = t.Q$. Hence, every $a_i$ in $x^+$ is also in $q^+$, and thus for every $i = 1, \ldots, k$, $\Delta \vdash (q \preceq a_i)$, implying that $q$ is inconsistent in $\Delta$.

Conversely, if $q$ is inconsistent in $\Delta$, then $\Delta \vdash q$ and $\Delta \not\vdash q$. Thus, there exist $A$ in $U$ and $a$ and $a'$ in $dom(A)$ such that $\Delta \vdash (q \preceq a \sqcap a')$, implying that $\Delta \vdash qa$ and $\Delta \vdash qa'$. By Lemma 3, $D^*$ contains two rows $t$ and $t'$ such that $qa \sqsubseteq t$ and $qa' \sqsubseteq t'$. This implies that $A$ can not be in $Q$ because otherwise, we would for instance have $qa = q$ and thus $qa' = qaa'$, which does not define a tuple. Since, by Definition 3, $\Delta \vdash (q \preceq a \sqcap a')$ implies that $a$ and $a'$ are in $q^+$, by Proposition 3, $A$ is in $Q^+$. Since $A$ is not in $Q$, $\mathcal{FD}$ contains $X \to A$ such that $X \subseteq Q^+$. It follows that $A$ is in $X^+$, $t.XA = xa$ and $t'.XA = xa'$. Therefore $x$ belongs to $inc(X \to A)$.

Summing up, we have found a tuple $t$ in $D^*$ and $X \to A$ in $\mathcal{FD}$ such that $t.X$ belongs to $inc(X \to A)$, $q \sqsubseteq t$ and $X \subseteq Q^+$. It thus follows from line 7 of Algorithm 3 that $q$ belongs to $\mathsf{Inc}(\Delta)$, which completes the proof. □

# F Proof of Proposition 8

**Proposition 8.** *Given $\Delta = (D, \mathcal{FD})$ and a query $Q$ : SELECT $X$ [WHERE $\Gamma$], Algorithm 6 correctly computes $ans^{\downarrow}_{\Delta}(Q)$ and $ans^{\uparrow}_{\Delta}(Q)$. Moreover, the following holds: $ans^{\downarrow}_{\Delta}(Q) \subseteq ans^{\uparrow}_{\Delta}(Q) \subseteq ans^{+}_{\Delta}(Q)$.*

*Proof* In this proof we respectively denote by $ans^{\downarrow}$ and $ans^{\uparrow}$ the two sets returned by Algorithm 6 and we successively show that $ans^{\downarrow} = ans^{\downarrow}(Q)$ and $ans^{\uparrow} = ans^{\uparrow}(Q)$.

First it is clear that all selected tuples are defined over $X$ and that they satisfy $\Gamma$. Moreover, assuming that the previous two conditions are satisfied, a tuple $t$ generates an $X$-value in $ans^{\downarrow}_{\Delta}(Q)$, if and only if $t$ is in every repair $R$ of $\Delta$, that is if and only if $t$ contains no conflicting value with respect to some dependency in $\mathcal{FD}$. This condition being precisely that on line 6 of Algorithm 6, we obtain that $ans^{\downarrow} = ans^{\downarrow}_{\Delta}(Q)$.

Now, given $x$ in $ans^{\uparrow}_{\Delta}(Q)$, assume that the condition on line 8 is not satisfied. In this case there exist $t$ in $D^*$ and $Y \to B$ in $\mathcal{FD}$ such that $t$ satisfies $\Gamma$, $t.Y$ is in $inc(Y \to B)$, $t.X = x$ and $t.B$ occurs in $x$. Thus, by the statement on line 9 in Algorithm 2, $D^*$ contains a tuple $t'$ such that $sch(t) \subseteq sch(t')$, $t'.Y = t.Y$ and $t'.B \neq t.B$. Hence, writing $x$ as $x'b$, $Rep(\Delta)$ contains a repair $R$ where $x'b$ occurs and a repair $R'$ where $x'b$ does not occur, showing that $x$ cannot belong to $ans^{\uparrow}_{\Delta}(Q)$. This is a contradiction showing that $ans^{\uparrow}_{\Delta}(Q) \subseteq ans^{\uparrow}$ holds.

Conversely, we first notice that for every tuple $q$ occurring in a repair $R$ but not in another repair $R'$, there exist $q'$ in $R'$, $B$ in $sch(q)$ and $Y \to B$ in $\mathcal{FD}$ such that $q.Y = q'.Y = y$, $y \in inc(Y \to B)$ and $q.B \neq q'.B$. Now, if $x$ is a tuple over $X$ for which the condition on line 8 is satisfied, then there exists $t$ in $D^*$ such that $t$ satisfies $\Gamma$, $t.X = x$ and for every $Y \to B$ in $\mathcal{FD}$ such that $y$ is in $inc(Y \to B)$, $B$ is not in $X$. Therefore, it turns out that $x = t.X$ occurs in $\pi_X(\sigma_\Gamma(R))$ for every $R$ in $Rep(\Delta)$, which shows that $x$ is in $ans^{\uparrow}_{\Delta}(Q)$.

As for the inclusions, in Algorithm 6, the condition on line 6 implies that on line 8, showing the first inclusion. Moreover, this second condition implies the one on line 4 of Algorithm 5, showing the second inclusion. The proof is therefore complete. □

37