



BIROn - Birkbeck Institutional Research Online

Levene, Mark and Vincent, Millist W. (2000) Justification for inclusion dependency normal form. *IEEE Transactions on Knowledge and Data Engineering* 12 (2), pp. 281-291. ISSN 1041-4347.

Downloaded from: <http://eprints.bbk.ac.uk/id/eprint/196/>

Usage Guidelines:

Please refer to usage guidelines at <https://eprints.bbk.ac.uk/policies.html>
contact lib-eprints@bbk.ac.uk.

or alternatively

Birkbeck ePrints: an open access repository of the research output of Birkbeck College

<http://eprints.bbk.ac.uk>

Levene, Mark and Vincent, Millist W. (2000) Justification for inclusion dependency normal form. *IEEE Transactions on Knowledge and Data Engineering* **12** (2) 281-291.

This is an exact copy of a paper published in *IEEE Transactions on Knowledge and Data Engineering* (ISSN 1041-4347). It is reproduced with permission from the publisher. Personal use of this material is permitted. However, permission to reprint/republish this material for advertising or promotional purposes or for creating new collective works for resale or redistribution to servers or lists, or to reuse any copyrighted component of this work in other works must be obtained from the IEEE. © 2000 IEEE.

Copyright and all rights therein are retained by authors or by other copyright holders. All persons downloading this information are expected to adhere to the terms and constraints invoked by copyright. This document or any part thereof may not be reposted without the explicit permission of the copyright holder.

Citation for this copy:

Levene, Mark and Vincent, Millist W. (2000) Justification for inclusion dependency normal form. *London: Birkbeck ePrints*. Available at: <http://eprints.bbk.ac.uk/archive/00000196>

Citation as published:

Levene, Mark and Vincent, Millist W. (2000) Justification for inclusion dependency normal form. *IEEE Transactions on Knowledge and Data Engineering* **12** (2) 281-291.

Justification for Inclusion Dependency Normal Form

Mark Levene and Millist W. Vincent

Abstract—Functional dependencies (FDs) and inclusion dependencies (INDs) are the most fundamental integrity constraints that arise in practice in relational databases. In this paper, we address the issue of normalization in the presence of FDs and INDs and, in particular, the semantic justification for *Inclusion Dependency Normal Form* (IDNF), a normal form which combines Boyce-Codd normal form with the restriction on the INDs that they be noncircular and key-based. We motivate and formalize three goals of database design in the presence of FDs and INDs: noninteraction between FDs and INDs, elimination of redundancy and update anomalies, and preservation of entity integrity. We show that, as for FDs, in the presence of INDs being free of redundancy is equivalent to being free of update anomalies. Then, for each of these properties, we derive equivalent syntactic conditions on the database design. Individually, each of these syntactic conditions is weaker than IDNF and the restriction that an FD not be embedded in the righthand side of an IND is common to three of the conditions. However, we also show that, for these three goals of database design to be satisfied simultaneously, IDNF is both a necessary and sufficient condition.

Index Terms—Relational database design, normal forms, functional dependency, inclusion dependency.

1 INTRODUCTION

FUNCTIONAL dependencies (FDs) [2], [25], [31], [1] generalize the notions of *entity integrity* and *keys* [13] and inclusion dependencies (INDs) [27], [8] generalize the notions of *referential integrity* and *foreign keys* [13], [17]. In this sense, FDs and INDs are the most fundamental data dependencies that arise in practice.

Relational database design in the presence of FDs is an established area in database theory which has been researched for more than 20 years [12], [3], [25], [31], [1]. The semantic justification of the normal forms in the presence of FDs is well-understood in terms of eliminating the so-called update anomalies and redundancy problems that can arise in a relation satisfying a set of FDs [6], [20], [9], [33], [34]. The advice that is given as a result of this investigation of the semantics of the normal forms is that, in order to eliminate the above-mentioned problems, we should design database schemas which are in *Boyce-Codd Normal Form* (BCNF) [12].

Despite the importance of INDs as integrity constraints, little research has been carried out on how they should be integrated into the normalization process of a relational database. Such an integration is fundamental to the success of a design since the enforcement of referential integrity is no simple matter [14]. Normal forms which include FDs and INDs have been considered in [7], [28], [23], [29], [4], but necessary and sufficient conditions in terms of remov-

ing the update anomalies and redundancy problems were not given. It is our goal in this paper to fill in this gap by providing sufficient and necessary semantics for *Inclusion Dependency Normal Form* (IDNF).

We consider some of the problems that occur in the presence of FDs and INDs through two examples. The first example illustrates the situation when an attribute is redundant due to interaction between FDs and INDs.

Example 1.1. Let HEAD be a relation schema, with attributes H and D, where H stands for head of department and D stands for department, and let LECT be a relation schema, with attributes L and D, where L stands for lecturer and, as before, D stands for department. Furthermore, let $F = \{\text{HEAD} : H \rightarrow D, \text{LECT} : L \rightarrow D\}$ be a set of FDs over a database schema $R = \{\text{HEAD}, \text{LECT}\}$, stating that a head of department manages a unique department and a lecturer works in a unique department, and $I = \{\text{HEAD}[\text{HD}] \subseteq \text{LECT}[\text{LD}]\}$ be a set of INDs over R stating that a head of department also works as a lecturer in the same department. We note that $I \cup (F - \{\text{HEAD} : H \rightarrow D\}) \models \text{HEAD} : H \rightarrow D$, where \models denotes logical implication, by the *pullback* inference rule (see Proposition 2.1) and, thus, the FD $\text{HEAD} : H \rightarrow D$ in F is redundant. Also, note that we have *not* assumed that $\text{HEAD} : D \rightarrow H$ in F and, thus, a department may have more than one head.

Two problems arise with respect to R and $F \cup I$. First, the interaction between F and I may lead to the logical implication of data dependencies that were not envisaged by the database designer and may not be easy to detect; in general, the implication problem for FDs and INDs is intractable (see the discussion in Section 2). In this example, the pullback inference rule implies that an FD in F is redundant.

Second, the IND $\text{HEAD}[\text{HD}] \subseteq \text{LECT}[\text{LD}]$ combined with the FD $\text{LECT} : L \rightarrow D$ imply that the attribute D in

• M. Levene is with the Department of Computer Science, University College London, Gower Street, London WC1E 6BT, UK.
E-mail: mlevene@cs.ucl.ac.uk.

• M.W. Vincent is with the Advanced Computing Research Centre, School of Computer and Information Science, University of South Australia, Adelaide, Australia 5095.
E-mail: millist.vincent@unisa.edu.au.

Manuscript received 5 May 1997; accepted 7 Jan. 1999.

For information on obtaining reprints of this article, please send e-mail to: tkde@computer.org, and reference IEEECS Log Number 104996.

HEAD is redundant since the department of a head can be inferred from the fact that L is a key for LECT. (Formally, this inference can be done with the aid of a relational algebra expression which uses renaming, join and projection; see [25], [31], [1] for details on the relational algebra.) Thus, $HEAD[HD] \subseteq LECT[LD]$ can be replaced by $HEAD[H] \subseteq LECT[L]$ and the attribute D in HEAD can be removed without any loss of information.

The second example illustrates the situation when the propagation of insertions due to INDs may result in the violation of entity integrity.

Example 1.2. Let EMP be a relation schema, with attributes E and P, where E stands for employee name and P stands for project title, and let PROJ be a relation schema, with attributes P and L, where, as before, P stands for project title and L stands for project location. Furthermore, let $F = \{EMP : E \rightarrow P\}$ be a set of FDs over a database schema $\mathbf{R} = \{EMP, PROJ\}$, stating that an employee works on a unique project, and $I = \{EMP[P] \subseteq PROJ[P]\}$ be a set of INDs over \mathbf{R} stating that an employee's project is one of the listed projects. We note that a project may be situated in several locations and, correspondingly, a location may be associated with several projects and, thus, $\{P, L\}$ is the primary key of PROJ.

The problem that arises with respect to \mathbf{R} and $F \cup I$ is that the righthand side, P, of the IND $EMP[P] \subseteq PROJ[P]$ is a proper subset of the primary key of PROJ. Let r_1 and r_2 be relations over EMP and PROJ, respectively. Suppose that an employee is assigned to a new project which has not yet been allocated a location and is thus not yet recorded in r_2 . Now, due to the IND in I, the insertion of the employee tuple into r_1 , having this new project, should be propagated to r_2 by inserting into r_2 a tuple recording the new project. But, since the location of the project is still unknown, then, due to entity integrity, it is not possible to propagate this insertion to r_2 .

We summarize the problems that we would like to avoid when designing relational databases in the presence of FDs and INDs. First, we should avoid redundant attributes, second, we should avoid the violation of entity integrity when propagating insertions, and, last, we should avoid any interaction between FDs and INDs due to the intractability of the joint implication problem for FDs and INDs. The main contributions of this paper are the formalization of these design problems and the result that if the database schema is in IDNF, then all of these problems are eliminated. We also demonstrate the robustness of IDNF by showing that, as for BCNF, removing redundancy from the database schema in the presence of FDs and INDs is also equivalent to eliminating update anomalies from the database schema.

The layout of the rest of the paper is as follows: In Section 2, we formally define FDs, INDs and their satisfaction and introduce the chase procedure as a means of testing and enforcing the satisfaction of a set of FDs and INDs. In Section 3, we formalize the notion of no interaction between a set of FDs and INDs. In Section 4, we characterize

redundancy in the presence of FDs and INDs. In Section 5, we characterize insertion and modification anomalies in the presence of FDs and INDs and show an equivalence between being free of either insertion or modification anomalies and being free of redundancy. In Section 6, we characterize a generalization of entity integrity in the presence of FDs and INDs. In Section 7, we define IDNF and present our main result that establishes the semantics of IDNF in terms of either the update anomalies or redundancy problems and the satisfaction of generalized entity integrity. Finally, in Section 8, we give our concluding remarks and indicate our current research direction.

Definition 1.1 (Notation). We denote the cardinality of a set S by $|S|$. The size of a set S is defined to be the cardinality of a standard encoding of S.

If S is a subset of T, we write $S \subseteq T$ and if S is a proper subset of T, we write $S \subset T$. We often denote the singleton $\{A\}$ simply by A. In addition, we often denote the union of two sets S, T, i.e., $S \cup T$, simply by ST.

2 FUNCTIONAL AND INCLUSION DEPENDENCIES

We formalize the notions of FDs and INDs and their satisfaction and define some useful subclasses of FDs and INDs. We also present the chase procedure for testing and enforcing the satisfaction of a set of FDs and INDs. The chase procedure is instrumental in proving our main results.

Definition 2.1 (Database schema and database). Let U be a finite set of attributes. A relation schema R is a finite sequence of distinct attributes from U. A database schema is a finite set $\mathbf{R} = \{R_1, R_2, \dots, R_n\}$ such that each $R_i \in \mathbf{R}$ is a relation schema and $\bigcup_i R_i = U$.

We assume a countably infinite domain of values, \mathcal{D} ; without loss of generality, we assume that \mathcal{D} is linearly ordered. An R-tuple (or, simply, a tuple whenever R is understood from context) is a member of the Cartesian product $\mathcal{D} \times \dots \times \mathcal{D}$ ($|R|$ times).

A relation r over R is a finite (possibly empty) set of R-tuples. A database d over \mathbf{R} is a family of n relations $\{r_1, r_2, \dots, r_n\}$ such that each $r_i \in d$ is over $R_i \in \mathbf{R}$. Given a tuple t over R and assuming that $r \in d$ is the relation in d over R, we denote the insertion of t into r by $d \cup \{t\}$ and the deletion of t from r by $d - \{t\}$.

From now on, we let \mathbf{R} be a database schema and d be a database over \mathbf{R} . Furthermore, we let $r \in d$ be a relation over the relation schema $R \in \mathbf{R}$.

Definition 2.2 (Projection). The projection of an R-tuple t onto a set of attributes $Y \subseteq R$, denoted by $t[Y]$ (also called the Y-value of t), is the restriction of t to Y. The projection of a relation r onto Y, denoted as $\pi_Y(r)$, is defined by $\pi_Y(r) = \{t[Y] \mid t \in r\}$.

Definition 2.3 (Functional Dependency). A functional dependency (or, simply, an FD) over a database schema \mathbf{R} is a statement of the form $R: X \rightarrow Y$ (alternatively, $X \rightarrow Y$ is an FD over the relation schema R), where $R \in \mathbf{R}$ and $X, Y \subseteq R$ are sets of attributes. An FD of the form $R: X \rightarrow Y$ is said to be trivial if $Y \subseteq X$; it is said to be standard if $X \neq \emptyset$.

An FD $R: X \rightarrow Y$ is satisfied in d, denoted by $d \models R: X \rightarrow Y$, whenever $\forall t_1, t_2 \in r$, if $t_1[X] = t_2[X]$, then $t_1[Y] = t_2[Y]$.

Definition 2.4 (Inclusion Dependency). An inclusion dependency (or, simply, an IND) over a database schema \mathbf{R}

is a statement of the form $R_i[X] \subseteq R_j[Y]$, where $R_i, R_j \in \mathbf{R}$ and $X \subseteq R_i, Y \subseteq R_j$ are sequences of distinct attributes such that $|X| = |Y|$. An IND is said to be trivial if it is of the form $R[X] \subseteq R[X]$. An IND $R[X] \subseteq S[Y]$ is said to be unary if $|X| = 1$. An IND is said to be typed if it is of the form $R[X] \subseteq S[X]$.

An IND $R_i[X] \subseteq R_j[Y]$ over \mathbf{R} is satisfied in d , denoted by $d \models R_i[X] \subseteq R_j[Y]$, whenever $\pi_X(r_i) \subseteq \pi_Y(r_j)$, where $r_i, r_j \in d$ are the relations over R_i and R_j , respectively.

From now on, we let F be a set of FDs over \mathbf{R} and $F_i = \{R_i : X \rightarrow Y \in F\}$, $\{1, 2, \dots, n\}$, be the set of FDs in F over $R_i \in \mathbf{R}$. Furthermore, we let I be a set of INDs over \mathbf{R} and let $\Sigma = F \cup I$.

Definition 2.5 (Logical implication). Σ is satisfied in d , denoted by $d \models \Sigma$, if $\forall \sigma \in \Sigma, d \models \sigma$.

Σ logically implies an FD or an IND σ , written $\Sigma \models \sigma$, if, whenever d is a database over \mathbf{R} , then the following condition is true:

if $d \models \Sigma$ holds then $d \models \sigma$ also holds.

Σ logically implies a set Γ of FDs and INDs over \mathbf{R} , written $\Sigma \models \Gamma$, if $\forall \sigma \in \Gamma, \Sigma \models \sigma$. We let Σ^+ denote the set of all FDs and INDs that are logically implied by Σ .

Definition 2.6 (Keys, BCNF and key-based INDs). A set of attributes $X \subseteq R_i$ is a superkey for R_i with respect to F_i if $F_i \models R_i : X \rightarrow R_i$ holds; X is a key for R_i with respect to F_i if it is a superkey for R_i with respect to F_i and for no proper subset $Y \subset X$ is Y a superkey for R_i with respect to F_i . We let $\text{KEYS}(F)$ be the set of all FDs of the form $X \rightarrow R_i$, where X is a key for R_i with respect to F_i , for $i \in \{1, 2, \dots, n\}$.

A database schema \mathbf{R} is in Boyce-Codd Normal Form (or, simply, BCNF) with respect to F if for all $R_i \in \mathbf{R}$, for all nontrivial FDs $R_i : X \rightarrow Y \in F_i$, X is a superkey for R_i with respect to F_i .

An IND $R_i[X] \subseteq R_j[Y]$ is superkey-based, respectively, key-based, if Y is a superkey, respectively, a key, for R_j with respect to F_j .

Definition 2.7. (Circular and noncircular sets of INDs). A set I of INDs over \mathbf{R} is circular if either

1. I contains a nontrivial IND $R[X] \subseteq R[Y]$, or
2. there exist m distinct relation schemas, $R_1, R_2, R_3, \dots, R_m \in \mathbf{R}$, with $m > 1$, such that I contains the INDs:

$$\begin{aligned} R_1[X_1] \subseteq R_2[Y_2], R_2[X_2] \\ \subseteq R_3[Y_3], \dots, R_m[X_m] \subseteq R_1[Y_1]. \end{aligned}$$

A set of INDs I is noncircular if it is not circular.

The class of proper circular INDs [21] defined below includes the class of noncircular INDs as a special case.

Definition 2.8 (Proper circular sets of INDs). A set I of INDs over \mathbf{R} is proper circular if it is either noncircular or whenever there exist m distinct relation schemas, $R_1, R_2, R_3, \dots, R_m \in \mathbf{R}$, with $m > 1$, such that I contains the INDs:

$$\begin{aligned} R_1[X_1] \subseteq R_2[Y_2], R_2[X_2] \subseteq R_3[Y_3], \dots, R_{m-1}[X_{m-1}] \\ \subseteq R_m[Y_m], R_m[X_m] \subseteq R_1[Y_1], \end{aligned}$$

then, for all $i \in \{1, 2, \dots, m\}$, we have $X_i = Y_i$.

It is well-known that Armstrong's axiom system [2], [25], [31], [1] can be used to compute F^+ and that Casanovas et al.'s axiom system [8] can be used to compute I^+ . However, when we consider FDs and INDs together computing Σ^+ was shown to be undecidable [27], [16]. On the other hand, when I is noncircular, then Mitchell's axiom system [27] can be used to compute Σ^+ [10]. Moreover, in the special case, when I is a set of unary INDs, then Cosmadakis's et al. axiom system [11] can be used to compute Σ^+ .

The implication problem is the problem of deciding whether $\sigma \in \Sigma^+$, where σ is an FD or IND and Σ is a set of FDs and INDs. It is well-known that the implication problem for FDs on their own is decidable in linear time [3]. On the other hand, the implication problem for INDs is, in general, PSPACE-complete [8]. The implication problem for noncircular INDs is NP-complete [26], [10]. Typed INDs have a polynomial time implication problem [15]. Unary INDs have a linear time implication problem [11]. When we consider FDs and INDs together, the implication problem is undecidable, as mentioned above. The implication problem for FDs and noncircular INDs is EXPTIME-complete and if the noncircular INDs are typed, then the implication problem is NP-hard [10]. FDs and unary INDs have a polynomial time implication problem [11].

The next proposition describes the *pullback* inference rule [27], [8], which allows us to infer an FD from an FD and an IND.

Proposition 2.1. If $\Sigma \models \{R[XY] \subseteq S[WZ], S : W \rightarrow Z\}$ and $|X| = |W|$, then $\Sigma \models R : X \rightarrow Y$.

Definition 2.9. (Reduced set of FDs and INDs). The projection of a set of FDs F_i over R_i onto a set of attributes $Y \subseteq R_i$, denoted by $F_i[Y]$, is given by $F_i[Y] = \{R_i : W \rightarrow Z \mid R_i : W \rightarrow Z \in F_i^+ \text{ and } WZ \subseteq Y\}$.

A set of attributes $Y \subseteq R_i$ is said to be reduced with respect to R_i and a set of FDs F_i over R_i (or, simply, reduced with respect to F_i if R_i is understood from context) if $F_i[Y]$ contains only trivial FDs. A set of FDs and INDs $\Sigma = F \cup I$ is said to be reduced if $\forall R_i[X] \subseteq R_j[Y] \in I, Y$ is reduced with respect to F_j .

It can easily be shown that it can be decided in polynomial time in the size of Σ whether Σ is reduced or not.

The chase procedure provides us with an algorithm which forces a database to satisfy a set of FDs and INDs.

Definition 2.10 (The chase procedure for INDs). The chase of d with respect to Σ , denoted by $\text{CHASE}(d, \Sigma)$, is the result of applying the following chase rules, that is, the FD and the IND rules, to the current state of d as long as possible. (The current state of d prior to the first application of the chase rule is its state upon input to the chase procedure.)

FD rule: If $R_j : X \rightarrow Y \in F_j$ and $\exists t_1, t_2 \in r_j$ such that $t_1[X] = t_2[X]$, but $t_1[Y] \neq t_2[Y]$, then, $\forall A \in Y$, change all the occurrences in d of the larger of the values of $t_1[A]$ and $t_2[A]$ to the smaller of the values of $t_1[A]$ and $t_2[A]$.

IND rule: If $R_i[X] \subseteq R_j[Y] \in I$ and $\exists t \in r_i$ such that $t[X] \notin \pi_Y(r_j)$, then add a tuple u over R_j to r_j , where $u[Y] = t[X]$ and $\forall A \in R_j - Y, u[A]$ is assigned a new value greater than any other current value occurring in the tuples of relations in the current state of d .

We observe that if we allow I to be circular, then the chase procedure does not always terminate [22]. (When the chase of d with respect to Σ does not terminate, then $\text{CHASE}(d, \Sigma)$ is said to violate a set of FDs G over \mathbf{R} , i.e., $\text{CHASE}(d, \Sigma) \not\models G$, if after some finite number of applications of the IND rule to the current state of d , resulting in d' , we have that $d' \not\models G$.) In the special case when I is in the class of *proper circular* INDs, then it was shown that the chase procedure always terminates [21].

The following theorem is a consequence of results in [29, Chapter 10].

Theorem 2.2. *Let $\Sigma = F \cup I$ be a set of FDs and proper circular INDs over a database schema \mathbf{R} . Then, the following two statements are true:*

1. $\text{CHASE}(d, \Sigma) \models \Sigma$.
2. $\text{CHASE}(d, \Sigma)$ terminates after a finite number of applications of the IND rule to the current state of d .

3 INTERACTION BETWEEN FDs AND INDs

As demonstrated by Proposition 2.1, FDs and INDs may interact in the sense that there may be FDs and INDs implied by a set of FDs and INDs which are not implied by the FDs or INDs taken separately. From the point of view of database design, interaction is undesirable since a database design may be normalized with respect to the set of FDs, but not with respect to the combined set of FDs and INDs. This is illustrated in the following example.

Example 3.1. Consider the database schema $\mathbf{R} = \{R, S\}$, where $R = S = ABC$ and a set Σ of FDs F and INDs I over \mathbf{R} given by $F = \{S : A \rightarrow BC\}$ and $I = \{R[AB] \subseteq S[AB]\}$. It can easily be verified that \mathbf{R} is in BCNF with respect to F . However, if we augment F with the FD $R : A \rightarrow B$, which is logically implied by Σ on using Proposition 2.1, then \mathbf{R} is not in BCNF with respect to the augmented set of FDs F since A is not a superkey for R with respect to the set of FDs $\{R : A \rightarrow B\}$.

The other major difficulty that occurs in database design when the set of FDs and INDs interact is a result of the fact that, as noted earlier, the implication problem for an arbitrary set of FDs and INDs is undecidable [27], [16]. Because of this, in general, it cannot be determined whether a database design is in BCNF with respect to an arbitrary set of FDs and INDs since the set of all logically implied FDs cannot be effectively computed. As a consequence, a desirable goal of database design is that the set of FDs and INDs do not interact. We now formalize the notion of noninteraction and characterize, for proper circular INDs, a special case when this noninteraction occurs.

Definition 3.1 (No interaction occurring between FDs and INDs). *A set of FDs F over \mathbf{R} is said not to interact with of set of INDs I over \mathbf{R} if*

1. for all FDs α over \mathbf{R} , for all subsets $G \subseteq F$, $G \cup I \models \alpha$ if and only if $G \models \alpha$, and
2. for all INDs β over \mathbf{R} , for all subsets $J \subseteq I$, $F \cup J \models \beta$ if and only if $J \models \beta$.

The following theorem is proven in [24] (see [29, Chapter 10]).

Theorem 3.1. *If \mathbf{R} is in BCNF with respect to a set of FDs F over \mathbf{R} , I is a proper circular set of of INDs over \mathbf{R} and $\Sigma = F \cup I$ is reduced, then F and I do not interact.*

As the next example shows, we cannot, in general, extend Theorem 3.1 to the case when the set of INDs I is not proper circular. In particular, by Proposition 2.1, Σ being reduced is a necessary condition for no interaction to occur between F and I , but it is not a sufficient condition for noninteraction.

Example 3.2. Consider a database schema $\mathbf{R} = \{R, S\}$, where $R = S = AB$, and a set Σ of FDs F and INDs I over \mathbf{R} given by $F = \{R : A \rightarrow B, S : B \rightarrow A\}$ and $I = \{R[A] \subseteq S[A], S[B] \subseteq R[B]\}$. It can easily be verified that Σ is reduced and that I is circular. On using the axiom system of [11], it follows that $\Sigma \models \{R : B \rightarrow A, S : A \rightarrow B\}$ and, thus, F and I interact.

As another example let $F = \{R : A \rightarrow B\}$ and $I = \{R[A] \subseteq R[B]\}$. It can easily be verified that Σ is reduced and I is circular. Again, F and I interact since, on using the axiom system of [11], $\Sigma \models \{R : B \rightarrow A, R[B] \subseteq R[A]\}$.

4 ATTRIBUTE REDUNDANCY

In this section, we investigate the conditions on database design which ensure the elimination of redundancy, a goal which has been long cited as one of the principal motivations for the use of normalization in database design [25], [31], [29], [1]. However, it has proven to be somewhat difficult to formalize the intuitive notion of redundancy and it was only relatively recently that this notion was formalized and its relationship to the classical normal forms was established [33], [34]. This definition of redundancy, and the associated normal form that guarantees redundancy elimination, is as follows.

Definition 4.1 (Value redundancy). *Let d be a database over \mathbf{R} that satisfies F and let $t \in r$ be a tuple, where $r \in d$ is a relation over a relation schema $\mathbf{R} \in \mathbf{R}$. The occurrence of a value $t[A]$, where $A \in \mathbf{R}$ is an attribute, is redundant in d with respect to F if, for every replacement of $t[A]$ by a distinct value $v \in \mathcal{D}$ such that $v \neq t[A]$, resulting in the database d' , we have that $d' \not\models F$.*

A database schema \mathbf{R} is said to be in Value Redundancy Free Normal Form (or, simply, VRFNF) with respect to a set of FDs F over \mathbf{R} if there does not exist a database d over \mathbf{R} and an occurrence of a value $t[A]$ that is redundant in d with respect to F .

We now illustrate the definition by a simple example.

Example 4.1. Consider the single relation scheme $\mathbf{R} = \{R\}$, where $R = ABC$, and a set F of FDs given by $F = \{R : A \rightarrow B\}$. Then, \mathbf{R} is not in VRFNF since, if we consider the relation r over \mathbf{R} , shown in Table 1, then the B -value, 2, present in both tuples, is redundant since $r \models F$ and replacing the value 2 in either tuple by another value results in F being violated.

The following result, which was established in [33], [34], shows that, given a set of FDs F , VRFNF is equivalent to

TABLE 1
The Relation r

A	B	C
1	2	3
1	2	4

BCNF. For the sake of completeness, we provide a sketch of the proof.

Theorem 4.1. *A database schema \mathbf{R} is in BCNF with respect to F if and only if \mathbf{R} is in VRFNF with respect to F .*

Proof. The if part follows by showing the contrapositive that if \mathbf{R} is not in BCNF, then it is not in VRFNF. This result follows because if \mathbf{R} is not in BCNF, then there exists a nontrivial implied FD $X \rightarrow A$, where X is not a superkey and, using a well-known construction (Theorem 7.1 in [31]), there exists a two tuple relation r in which the tuples are identical on X . The relation r is not in VRFNF since changing either of the A -values results in $X \rightarrow A$ being violated. The only if part follows from the observation that if an occurrence $t[A]$ is redundant in a relation r , then there must exist $X \rightarrow Y \in F$, with $A \in Y$, and another tuple $t_1 \in r$ such that $t[X] = t_1[X]$. This implies that X cannot be a superkey and, hence, that \mathbf{R} is not in BCNF and, so, establishes the only if part. \square

In [34], it was shown that, in the presence of multivalued dependencies, 4NF (fourth normal form [18]) is also equivalent to VRFNF. Somewhat surprisingly, the syntactic equivalent for VRFNF in the most general case, where join dependencies are present, is a new normal form that is weaker than PJ/NF (project-join normal form [19]) and 5NF (fifth normal form [25]) [34].

The concept of value redundancy is of little use though in evaluating database designs in the presence of INDs because of the following result. It demonstrates that no design where the constraints involve nontrivial INDs can be in VRFNF with respect to the set of INDs.

Lemma 4.2. *Let $\Sigma = F \cup I$ be a set of FDs and noncircular INDs over a database schema \mathbf{R} . Then, \mathbf{R} is not in VRFNF with respect to Σ if I contains at least one nontrivial IND.*

Proof. Let $R_i[X] \subseteq R_j[Y]$ be a nontrivial IND in I . We construct a database d such that the relation r_i in d over R_i has in it a single tuple t_i containing zeros and every other relation r_k in d is empty. Now, let $d' = \text{CHASE}(d, \Sigma)$ and, thus, by Theorem 2.2, $d' \models \Sigma$. Due to the non-circularity of I , we have in d' that $r_i = r'_i$, where r'_i is the current state of r_i in d' . Let r'_j be the current state of the relation r_j over R_j in d' . Then, the Y -values of the tuple in r'_j must contain zeros since $d' \models \Sigma$. Thus, all of the zeros in the single tuple in r'_i are redundant since changing any of them results in $R_i[X] \subseteq R_j[Y]$ being violated. \square

Due to the above lemma, we require only that the design be in VRFNF with respect to the FDs, but not with respect to the INDs. However, as noted by others [30], [23], [29], an even stronger form of redundancy can occur in a database in the presence of INDs. We refer to this as attribute

redundancy, which was illustrated in Example 1.1 given in the introduction.

Definition 4.2 (Attribute redundancy). *An attribute A in a relation schema $\mathbf{R} \in \mathbf{R}$ is redundant with respect to Σ if, whenever d is a database over \mathbf{R} which satisfies Σ and $r \in d$ is a nonempty relation over \mathbf{R} , then, for every tuple $t \in r$, if $t[A]$ is replaced by a distinct value $v \in \mathcal{D}$ such that $v \neq t[A]$, resulting in the database d' , then $d' \not\models \Sigma$.*

A database schema \mathbf{R} is said to be in Attribute Redundancy Free Normal Form (or, simply, ARFNF) with respect to a set of FDs and INDs Σ over \mathbf{R} if there does not exist an attribute A in a relation schema $\mathbf{R} \in \mathbf{R}$ which is redundant with respect to Σ .

The next example shows that ARFNF is too weak when Σ contains only FDs, highlighting the difference between VRFNF and ARFNF.

Example 4.2. Consider the relation r of Example 4.1, shown in Table 1, and let r' be the result of adding the tuple $t = \langle 5, 2, 4 \rangle$ to r . Then, it can easily be verified that, although \mathbf{R} is not in VRFNF with respect to F , \mathbf{R} is in ARFNF with respect to F since if we replace any value in t by a distinct value, then the resulting database still satisfies F .

Combining Definitions 4.1 and 4.2, we can define redundancy free normal form.

Definition 4.3 (Redundancy free normal form). *A database schema \mathbf{R} is said to be in Redundancy Free Normal Form (or, simply, RFNF) with respect to a set of FDs and INDs Σ over \mathbf{R} if it is in VRFNF with respect to F and in ARFNF with respect to Σ .*

The next theorem shows that, when the set of INDs is noncircular, then RFNF is equivalent to the set of FDs and INDs being reduced and to the database schema being in BCNF.

Theorem 4.3. *Let $\Sigma = F \cup I$ be a set of FDs and noncircular INDs over a database schema \mathbf{R} . Then, \mathbf{R} is in RFNF with respect to Σ if and only if Σ is reduced and \mathbf{R} is in BCNF with respect to F .*

Proof. If. By Theorem 4.1, \mathbf{R} is in VRFNF with respect to F . So, it remains to show that \mathbf{R} is in ARFNF.

In the proof, we utilize a directed graph representation, $G_I = (N, E)$, of the set of INDs I , which is constructed as follows (see [30]): Each relation schema \mathbf{R} in \mathbf{R} has a separate node in N labeled R ; we do not distinguish between nodes and their labels. There is an arc $(R, S) \in E$ if and only if there is a nontrivial IND $R[X] \subseteq S[Y] \in I$. It can easily be verified that there is a path in G_I from R to S if and only if, for some IND $R[X] \subseteq S[Y]$, we have $I \models R[X] \subseteq S[Y]$. Moreover, since I is noncircular, we have that G_I is acyclic.

Let $A \in R_i$ be an attribute, where $R_i \in \mathbf{R}$ is a relation schema. We construct a database d having a nonempty relation $r_i \in d$, which exhibits the fact that A is nonredundant with respect to Σ .

We first initialize the database d to be a database d_0 as follows: Let $r_i = r_i^0$ have a single tuple t such that, for all $B \in R_i - A$, $t[B] = 0$ and $t[A] = 1$. All other relations r_k^0

over relation schemas R_k are initialized to be empty in d_0 . Therefore, by Theorem 2.2, we have that $d_1 = \text{CHASE}(d_0, \Sigma) \models \Sigma$. Let r_i^1 in d_1 be the current state of r_i . Then, by Theorem 3.1, we have $r_i^1 = r_i$ since $d_0 \models F$ and the current state r_k^1 in d_1 of a relation r_k^0 is empty if there does not exist a path in G_I from R_i to R_k .

Let $d_2 = (d_1 - \{r_i^1\}) \cup \{r_i^2\}$, where r_i^2 has a single tuple t' such that, for all $B \in R_i$, including A , $t[B] = 0$. Therefore, by Theorem 2.2, we have that $d_3 = \text{CHASE}(d_2, \Sigma) \models \Sigma$. Moreover, as above, $r_i^3 = r_i^2$, where r_i^3 in d_3 is the current state of r_i since $d_2 \models F$. Now, let $d = (d_3 - \{r_i^3\}) \cup \{r_i\}$ be the final state of the initialization of d . Then, $d \models F$ since $d_3 \models F$. We claim that it is also the case that $d \models I$.

Let us call a nontrivial IND $R_i[X] \subseteq R_j[Y] \in I$ a *source* IND if $A \in X$. By the projection and permutation inference rule for INDs [8], we assume, without loss of generality, that a source IND has the form $R_i[VA] \subseteq R_j[WB]$.

Due to the noncircularity of I , any current state r_k in d of a relation r_k^3 is empty if there does not exist a path in G_I from R_i to R_k ; therefore, for such r_k , we have $r_k = r_k^3 = r_k^2 = r_k^1 = r_k^0 = \emptyset$. Now, if there is an arc from R_i to R_j in G_I , then there is some IND $R_i[X] \subseteq R_j[Y]$ in I . There are two cases to consider.

First, if $A \notin X$, then $t[X] = t'[X]$ contains only zeros and, thus, $d \models R_i[X] \subseteq R_j[Y]$. Second, if $A \in X$, then $R_i[X] \subseteq R_j[Y]$ is a source IND $R_i[VA] \subseteq R_j[WB]$. Let r_j in d be the current state of r_j^3 , i.e., $r_j = r_j^3$. Then, $t[VA] \in \pi_{VA}(r_j)$ since $r_j^0 \subseteq r_j^1 \subseteq r_j^2 \subseteq r_j^3 = r_j$ and $t[VA] \in \pi_{VA}(r_j^0)$. Therefore, $d \models R_i[VA] \subseteq R_j[WB]$. It follows that, for any IND, $R_i[X] \subseteq R_j[Y]$ such that there is a path from R_i to R_j and such that $I \models R_i[X] \subseteq R_j[Y]$, we have $d \models R_i[X] \subseteq R_j[Y]$ since $d_3 \models R_i[X] \subseteq R_j[Y]$. Thus, $d \models I$ as required. The if part is now concluded since d_3 and d differ only by the replacement of $t[A] = 0$ by $t[A] = 1$.

Only if. By Theorem 4.1, if \mathbf{R} is not in BCNF with respect to F , then it is not in VRFNF. So, assuming that Σ is not reduced, it remains to show that \mathbf{R} is not in ARFNF. By this assumption, there exists an IND $R_i[X] \subseteq R_j[Y] \in I$ such that $W \rightarrow B \in F_j[Y]$ is a nontrivial FD. By the pullback inference rule, there is a nontrivial FD $V \rightarrow A \in F_i[X]$ since $R_i[VA] \subseteq R_j[WB]$ by the projection and permutation inference rule for INDs [8]. Now, let $t \in r_i$ be a tuple, where $r_i \in d$ is a nonempty relation over R_i and assume that $d \models \Sigma$. It follows that A is redundant with respect to Σ since, whenever we replace $t[A]$ by a distinct value resulting in a database d' , it can be seen that $d' \not\models R_i[VA] \subseteq R_j[WB]$, otherwise $d' \not\models R_j : W \rightarrow B$ contrary to assumption. \square

We next construct two examples which demonstrate that if the conditions of Theorem 4.3 are violated, then \mathbf{R} is not in ARFNF.

Example 4.3. Let \mathbf{R} be the database schema from Example 1.1 and consider the database d over \mathbf{R} , shown in Tables 2 and 3, respectively. It can be verified that the set of FDs and INDs for this example is not reduced, but that \mathbf{R} is in BCNF with respect to the set of FDs. The attribute D of the relation schema HEAD can be seen to

TABLE 2
The Relation in d over HEAD

H	D
h_1	e_1

be redundant since changing e_1 in Table 2 causes I to be violated. A similar situation occurs for every other database defined over \mathbf{R} and, so, \mathbf{R} is not in ARFNF.

Example 4.4. Let $\mathbf{R} = \{\text{STUDENT}, \text{ENROL}\}$, be a database schema, with $\text{STUDENT} = \{\text{Stud_ID}, \text{Name}\}$ and $\text{ENROL} = \{\text{Stud_ID}, \text{Course}, \text{Address}\}$. Furthermore, let $F = \{\text{ENROL} : \text{Stud_ID} \rightarrow \text{Address}\}$ be a set of FDs over \mathbf{R} and $I = \{\text{ENROL}[\text{Stud_ID}] \subseteq \text{STUDENT}[\text{Stud_ID}]\}$ be a set of INDs over \mathbf{R} . It can be verified that the set of FDs and INDs for this example is reduced, but that \mathbf{R} is not in BCNF with respect to the set of FDs. Then, \mathbf{R} is not in RFNF because it is not in VRFNF since both occurrences of a_1 in ENROL are redundant in the database d shown in Tables 4 and 5, respectively.

As the next example shows, we cannot extend Theorem 4.3 to the case when the set of INDs I is circular.

Example 4.5. Consider a database schema $\mathbf{R} = \{R, S\}$, where $R = AB$ and $S = A$, and a set Σ of FDs F and INDs I over \mathbf{R} given by $F = \{R : A \rightarrow B, R : B \rightarrow A\}$ and $I = \{R[A] \subseteq S[A], S[A] \subseteq R[A]\}$. It can be verified that Σ is reduced, \mathbf{R} is in BCNF with respect to F , and that I is proper circular, unary, typed, and also key-based.

Let d be a database over \mathbf{R} such that $d \models \Sigma$, let $r \in d$ be a nonempty relation over R , and let $t \in r$ be a tuple. Assume, without loss of generality, that d' is the database resulting from replacing $t[A] = 0$ by a distinct value 1, resulting in a tuple t' , with $t'[A] = 1$. In order to conclude the example, we show that $d' \not\models \Sigma$. Assume to the contrary that $d' \models \Sigma$. Thus, there must be a tuple $u \in r$ which is distinct from t and such that $u[A] = 1$. If this is not the case, then $d' \not\models R[A] \subseteq S[A]$ since $\pi_A(r) = \pi_A(s)$, where $s \in d$ is a relation over S . Moreover, $u[B] = t[B] = t'[B]$, otherwise $d' \not\models R : A \rightarrow B$. It follows that $u[B] = t[B]$, but $u[A] \neq t[A]$ and, thus, $d' \not\models R : B \rightarrow A$, contrary to assumption. Therefore, the attribute $A \in R$ must be redundant with respect to Σ . The reader can easily verify that the attribute $A \in S$ is also redundant with respect to Σ even if F was empty. It appears that, in this example, the relation schema S can be removed from \mathbf{R} without any loss of semantics.

5 INSERTION AND MODIFICATION ANOMALIES

In this section, we investigate the conditions under which a database design ensures the elimination of key-based update anomalies (as distinct from other types of update anomalies as investigated in [6], [33]). This concept was originally introduced in [20] to deal with the insertion and deletion of tuples and was later extended in [32], [33] to include the modifications of tuples. A key-based update anomaly is defined to occur when an update to a relation, which can either be an insertion or a deletion or a modification, results in the new relation satisfying key

TABLE 3
The Relation in d over LECT

L	D
h_1	e_1
l_1	e_1

uniqueness, but violating some other constraint on the relation. The reason for this being considered undesirable is that the enforcement of key uniqueness can be implemented via relational database software in a much more efficient manner than the enforcement of more general constraints such as FDs [25], [31], [1]. So, if the satisfaction of all the constraints on a relation is a result of key uniqueness then the integrity of the relation after an update can be easily enforced, whereas the existence of a key-based update anomaly implies the converse. Herein, we formalize these concepts based on this approach with the only difference being that, because of the presence of INDs, we allow an update to propagate to other relations by the chase procedure. We show that being free of insertion anomalies is equivalent to being free of modification anomalies. In addition, we show that, when the INDs are noncircular, then being free of either insertion or modification anomalies is equivalent to the set of FDs and INDs being reduced and the database schema being in BCNF with respect to the set of FDs. We do not consider deletion anomalies since, in the presence of FDs, removing a tuple from a relation that satisfies a set of FDs does not cause any violation of an FD in the set.

Definition 5.1 (Compatible tuple). A tuple t over \mathbf{R} is compatible with d with respect to a set of FDs and INDs $\Sigma = F \cup I$ over \mathbf{R} (or, simply, compatible with d whenever Σ is understood from context) if $d \cup \{t\} \models \text{KEYS}(F)$.

Definition 5.2 (Free of insertion anomalies). A database d over \mathbf{R} has an insertion violation with respect to a set of FDs and INDs $\Sigma = F \cup I$ over \mathbf{R} (or, simply, d has an insertion violation whenever Σ is understood from context) if

1. $d \models \Sigma$, and
2. there exists a tuple t over \mathbf{R} which is compatible with d but $\text{CHASE}(d \cup \{t\}, I) \not\models \Sigma$.

A database schema \mathbf{R} is free of insertion anomalies with respect to Σ (or, simply, \mathbf{R} is free of insertion anomalies if Σ is understood from context) if there does not exist a database d over \mathbf{R} which has an insertion violation.

We note that, in Definition 5.2, we have utilized the chase procedure to enforce the propagation of insertions of tuples due to the INDs in I . As an example of an insertion violation, consider the database schema \mathbf{R} in Example 1.1 and let d be the database, where r_1 , the relation over HEAD,

TABLE 4
The Relation in d over STUDENT

Stud_ID	Name
s_1	n_1

TABLE 5
The Relation in d over ENROL

Stud_ID	Course	Address
s_1	c_1	a_1
s_1	c_2	a_1

is empty and r_2 , the relation over LECT, contains the single tuple $\langle 0, 0 \rangle$. Then, d has an insertion violation when the tuple $\langle 0, 1 \rangle$ is inserted into r_1 since applying the chase procedure results in $\langle 0, 1 \rangle$ being added to the relation r_2 and, thus, violating the FD $L \rightarrow D$.

The next theorem shows that, assuming that the set of INDs is noncircular, being in BCNF and the set of FDs and INDs being reduced is equivalent to being free of insertion anomalies.

Theorem 5.1. Let $\Sigma = F \cup I$ be a set of FDs and noncircular INDs over a database schema \mathbf{R} . Then, \mathbf{R} is free of insertion anomalies if and only if Σ is reduced and \mathbf{R} is in BCNF with respect to F .

Proof. *If.* Let d be a database over \mathbf{R} such that $d \models \Sigma$ and let t over R_i , where $r_i \in \mathbf{R}$, be a tuple which is compatible with d . It remains to show that $\text{CHASE}(d \cup \{t\}, I) \models \Sigma$. We first claim that $\text{CHASE}(d \cup \{t\}, I) = \text{CHASE}(d \cup \{t\}, \Sigma)$. This holds due to Theorem 3.1, implying that the FD rule need never be invoked during the computation of $\text{CHASE}(d \cup \{t\}, \Sigma)$. Moreover, $d \cup \{t\} \models F$ due to the fact that \mathbf{R} is in BCNF with respect to F and t is compatible with d . So, by Theorem 2.2, we have $\text{CHASE}(d \cup \{t\}, \Sigma) \models \Sigma$ and, thus, $\text{CHASE}(d \cup \{t\}, I) \models \Sigma$ as required.

Only if. There are two cases to consider.

Case 1. If \mathbf{R} is not in BCNF with respect to F , then some $R_i \in \mathbf{R}$ is not in BCNF with respect to F_i . Thus, there is a nontrivial FD $X \rightarrow Y \in F_i$ such that X is not a superkey for R_i with respect to F_i . Assume that X is reduced with respect to R_i and F_i ; otherwise, replace X by a reduced subset W of X such that $W \rightarrow X \in F_i^+$. It follows that X is a proper subset of a superkey of R_i with respect to F_i . Let r_i over R_i contain a single tuple containing zeros and let all other relations in d be empty. We can assume without loss of generality that $d \models \Sigma$; otherwise, we let d be $\text{CHASE}(d, \Sigma)$. Due to I being noncircular, the state of r_i remains unchanged in $\text{CHASE}(d, \Sigma)$. Now, let t be a tuple whose X -values are zeros and such that all its other values are ones. Then, t is compatible with d , but $\text{CHASE}(d \cup \{t\}, I) \not\models \Sigma$ since the FD $X \rightarrow Y$ will be violated in the current state of r_i . Therefore, \mathbf{R} is not free of insertion anomalies.

Case 2. If Σ is not reduced, but \mathbf{R} is in BCNF with respect to F , then we have an IND $R_i[X] \subseteq R_j[Y]$, where Y is a proper superset of a key, say W , for R_j with respect to F_j . We let r_j over R_j contain a single tuple, say t_j , containing zeros and let all other relations in d , including r_i over R_i , be empty. We can assume without loss of generality that $d \models \Sigma$; otherwise, we let d be $\text{CHASE}(d, \Sigma)$. Due to I being noncircular, the state of r_i remains unchanged in $\text{CHASE}(d, \Sigma)$. Now, let t over R_i

TABLE 6
The Relation in d over STUDENT

Stud_ID	Name
s_1	n_1

be a tuple which agrees with t_j on its W -value but disagrees with t_j on the rest of its values. Then, t is compatible with d , but $\text{CHASE}(d \cup \{t\}, I) \not\models \Sigma$ since the FD $W \rightarrow R_j$ will be violated in the resulting current state of r_j . \square

To illustrate this theorem, we note first that the example given before Theorem 5.1 demonstrates the case where a database schema has an insertion anomaly when the set of dependencies is not reduced. Alternatively, the following example demonstrates the case of a database schema not being in BCNF and having an insertion anomaly.

Example 5.1. Let $\mathbf{R}, \Sigma = F \cup I$, be as in Example 4.4. We start with the database d shown in Tables 6 and 7, respectively. If we then insert the tuple $\langle s_1, c_2, a_2 \rangle$, which is compatible with ENROL, into the ENROL relation, applying the chase procedure results in the database d' shown in Tables 8 and 9, respectively, where n_2 is a new value. It can be seen that d' violates Σ and, so, \mathbf{R} is not free of insertion anomalies.

In the next example, we show that the only if part of Theorem 5.1 is, in general, false, even when I is a proper circular set of INDs.

Example 5.2. Consider a database schema $\mathbf{R} = \{R, S\}$, where $R = S = AB$, and a set Σ of FDs F and INDs I over \mathbf{R} given by, $F = \{R : A \rightarrow B, S : A \rightarrow B\}$ and $I = \{R[AB] \subseteq S[AB], S[AB] \subseteq R[AB]\}$. It can easily be verified that I is proper circular, \mathbf{R} is in BCNF, but that Σ is not reduced. In addition, \mathbf{R} is free of insertion anomalies since, for any database $d = r, s$ such that $d \models \Sigma$, where r and s are the relations in d over R and S , respectively, we have that $r = s$ due to I . If we drop $S[AB] \subseteq R[AB]$ from I , then, as in the proof of the only if part of Theorem 5.1, \mathbf{R} has an insertion violation.

The next example illustrates that we cannot, in general, extend Theorem 5.1 to the case when the set of INDs I is circular, even when Σ is reduced, due to possible interaction between the FDs and INDs.

Example 5.3. Consider a database schema $\mathbf{R} = \{R\}$, where $R = AB$, and a set Σ of FDs F and INDs I over \mathbf{R} given by, $F = \{R : A \rightarrow B\}$ and $I = \{R[A] \subseteq R[B]\}$. It can be verified that Σ is reduced, but I is circular. As was shown in Example 3.2, although Σ is reduced $\Sigma \models \{R : B \rightarrow A, R[B] \subseteq R[A]\}$ and, thus, F and I interact.

TABLE 7
The Relation in d over ENROL

Stud_ID	Course	Address
s_1	c_1	a_1

TABLE 8
The Relation in d' over STUDENT

Stud_ID	Name
s_1	n_1
s_1	n_2

Let d be a database over \mathbf{R} such that the relation r over \mathbf{R} contains the single tuple $\langle 0, 0 \rangle$ and let t be the tuple $\langle 1, 0 \rangle$. Then, d has an insertion violation since $d \models \Sigma$, t is compatible with d , but $\text{CHASE}(d \cup \{t\}, I) \not\models R : B \rightarrow A$. (In fact, in this case, the chase procedure does not terminate, but, since t is inserted into r and the chase procedure does not modify any of the tuples in its input database, then it does not satisfy Σ ; see the comment after Definition 2.10.)

We now formally define the second type of key-based update anomaly, a modification anomaly, following the approach in [32], [33] with the only difference again being that the chase procedure is used to propagate the effects of the change into other relations.

Definition 5.3 (Free of modification anomalies). A database d over \mathbf{R} has a modification violation with respect to a set of FDs and INDs $\Sigma = F \cup I$ over \mathbf{R} (or, simply, d has a modification violation whenever Σ is understood from context) if

1. $d \models \Sigma$ and
2. there exists a tuple $u \in r$, where $r \in d$ is the relation over \mathbf{R} , and a tuple t over \mathbf{R} which is compatible with $d - \{u\}$ but $\text{CHASE}((d - \{u\}) \cup \{t\}, I) \not\models \Sigma$.

A database schema \mathbf{R} is free of modification anomalies with respect to Σ (or, simply, \mathbf{R} is free of modification anomalies if Σ is understood from context) if there does not exist a database d over \mathbf{R} which has a modification violation.

Theorem 5.2. Let $\Sigma = F \cup I$ be a set of FDs and noncircular INDs over a database schema \mathbf{R} . Then, \mathbf{R} is free of modification anomalies if and only if \mathbf{R} is free of insertion anomalies.

Proof. If. Let d be a database over \mathbf{R} such that $d \models \Sigma$, t be a tuple that is compatible with d , and $u \in r$ be a tuple, where $r \in d$ is the relation over \mathbf{R} . It follows that t is compatible with $d - \{u\}$. We need to show that if $\text{CHASE}(d \cup \{t\}, I) \models \Sigma$, then $\text{CHASE}((d - \{u\}) \cup \{t\}, I) \models \Sigma$. By Theorem 2.2 of the chase procedure, $\text{CHASE}((d - \{u\}) \cup \{t\}, I) \models I$, so it remains to show that $\text{CHASE}((d - \{u\}) \cup \{t\}, I) \models F$. Let $\text{CHASE}(d \cup \{t\}, I) = \{r_1, r_2, \dots, r_n\}$ and let $\text{CHASE}((d - \{u\}) \cup \{t\}, I) = \{s_1, s_2, \dots, s_n\}$. Then, by Definition 2.10 of the chase procedure, we have that for all $i \in \{1, 2, \dots, n\}$, $s_i \subseteq r_i$. It

TABLE 9
The Relation in d' over ENROL

Stud_ID	Course	Address
s_1	c_1	a_1
s_1	c_2	a_2

follows that $\text{CHASE}((d - \{u\}) \cup \{t\}, I) \models \Sigma$ since, by Definition 5.2, $\text{CHASE}(d \cup \{t\}, I) \models F$.

Only if. If \mathbf{R} is free of modification anomalies, then, by a similar argument to that made in the *only if* part of the proof of Theorem 5.1, it follows that Σ is reduced and \mathbf{R} is in BCNF. In the first case, we add an additional tuple u over R_i , which contains ones, to the original state of r_i and, in the second case, we add an additional tuple u over R_i , which contains zeros, to the original state of r_i . The result now follows by the *if* part of Theorem 5.1. \square

Combining Theorems 4.3, 5.1, and 5.2, we obtain the next result.

Corollary 5.3. *Let $\Sigma = F \cup I$ be a set of FDs and noncircular INDs over a database schema \mathbf{R} . Then, the following statements are equivalent:*

1. \mathbf{R} is free of insertion anomalies.
2. \mathbf{R} is free of modification anomalies.
3. \mathbf{R} is in RFNF.

6 GENERALIZED ENTITY INTEGRITY

In this section, we justify superkey-based INDs on the basis that they do not cause the propagation of the insertion of tuples that represent undefined entities, thus causing the violation of entity integrity. This problem was illustrated in Example 1.2 given in the introduction.

In the next definition, we view the chase procedure as a mechanism which enforces the propagation of insertions of tuples due to the INDs in I .

Definition 6.1 (Generalized entity integrity). *Let t be a tuple that is added to a relation r_i over R_i , in the current state of a database d , during the computation of $\text{CHASE}(d, \Sigma)$. Then, t is entity-based if there exists at least one key X for R_i with respect to F_i such that for all $A \in X$, $t[A]$ is not a new value that is assigned to t as a result of invoking the IND rule.*

A database schema \mathbf{R} satisfies generalized entity integrity with respect to a set $\Sigma = F \cup I$ of FDs and INDs over \mathbf{R} if, for all databases d over \mathbf{R} , all the tuples that are added to relations in the current state of d during the computation of $\text{CHASE}(d, \Sigma)$ are entity-based.

The next theorem shows that satisfaction of generalized entity integrity is equivalent to the set of INDs being superkey-based.

Theorem 6.1. *A database schema \mathbf{R} satisfies generalized entity integrity with respect to a set of FDs and INDs $\Sigma = F \cup I$ if and only if I is superkey-based.*

Proof. If I is superkey-based, then the result immediately follows by the definition of the IND rule (see Definition 2.10). On the other hand, if I is not superkey-based, then there is some IND $R_i[X] \subseteq R_j[Y] \in I$ such that Y is not a superkey for R_j with respect to F_j . Let d be a database over \mathbf{R} such that all its relations apart for r_i over R_i are empty. The relation r_i has a single tuple. By the definition of the FD rule (see Definition 2.10), we have that, for every key, say K , of R_j with respect to F_j , there is at least one attribute, say $A \in K$, such that the tuple t_j added to r_j over R_j is assigned a new A -value by the

IND rule; otherwise, contrary to assumption, we can deduce that Y is a superkey for R_j with respect to F_j . It follows that \mathbf{R} does not satisfy generalized entity integrity, concluding the proof. \square

7 INCLUSION DEPENDENCY NORMAL FORM

A database schema is in IDNF with respect to a set of FDs and INDs if it is in BCNF with respect to the set of FDs and the set of INDs is noncircular and key-based. We show that a database schema is in IDNF if and only if it satisfies generalized entity integrity and is either free of insertion anomalies or free of modification anomalies or in redundancy free normal form.

We next formally define IDNF (cf. [28], [29]).

Definition 7.1 (Inclusion dependency normal form). *A database schema \mathbf{R} is in Inclusion Dependency Normal Form (IDNF) with respect to a set of Σ of FDs F and INDs I over \mathbf{R} (or, simply, in IDNF if Σ is understood from context) if*

1. \mathbf{R} is in BCNF with respect to F and
2. I is a noncircular and key-based set of INDs.

We note that if the set of INDs I is empty, then \mathbf{R} being in IDNF is equivalent to \mathbf{R} being in BCNF. We further note that we have not restricted the FDs in F to be standard.

The next result follows from Corollary 5.3, Theorem 6.1, and Definition 7.1.

Theorem 7.1. *Let $\Sigma = F \cup I$ be a set of FDs and noncircular INDs over a database schema \mathbf{R} . Then, the following statements are equivalent:*

1. \mathbf{R} is in IDNF
2. \mathbf{R} is free of insertion anomalies and satisfies generalized entity integrity.
3. \mathbf{R} is free of modification anomalies and satisfies generalized entity integrity.
4. \mathbf{R} is in RFNF and satisfies generalized entity integrity.

8 CONCLUDING REMARKS

We have identified three problems that may arise when designing databases in the presence of FDs and INDs, apart from the update anomalies and redundancy problems that may arise in each relation due to the FDs considered on their own. The first problem is that of attribute redundancy, the second problem is the potential violation of entity integrity when propagating insertions, and the third problem concerns avoiding the complex interaction which may occur between FDs and INDs and the intractability of determining such interaction. The first problem was formalized through RFNF and it was shown in Corollary 5.3 that a database schema is in RFNF with respect to a set of FDs and INDs if and only if it is free of insertion anomalies or equivalently free of modification anomalies. This result can be viewed as an extension of a similar result when considering FDs on their own. The second problem was formalized through generalized entity integrity and it was shown in Theorem 6.1 that a database schema satisfies generalized entity integrity with respect to

a set of FDs and INDs if and only the set of INDs is superkey-based. The third problem was formalized through the noninteraction of the implication problem for FDs and INDs and it was shown in Theorem 3.1 that a set of FDs and INDs do not interact when the set of INDs is proper circular, the set of FDs and INDs are reduced, and the database schema is in BCNF. Combining all these result together, we obtained, in Theorem 7.1, three equivalent semantic characterizations of IDNF. Theorem 7.1 justifies IDNF as a robust normal form that eliminates both redundancy and update anomalies from the database schema.

If the goal of normalization is to reduce redundancy, then it seems that, apart from \mathbf{R} being in BCNF, in general, we must restrict the set of INDs to be noncircular (see Example 4.5). Nonetheless, circular sets of INDs arise in practice, for example, when we want to express pairwise consistency. (Two relation schemas R and S are *consistent* if the set of INDs I includes the two INDs: $R[R \cap S] \subseteq S[R \cap S]$ and $S[R \cap S] \subseteq R[R \cap S]$. A database schema \mathbf{R} is *pairwise consistent* if every pair of its relation schemas are consistent [5]; we note that pairwise consistency can be expressed by a set of proper circular INDs.) In this case, we need alternative semantics to express the goal of normalization. A minimal requirement is that the FDs and INDs have no interaction. By Theorem 3.1, as long as the set of FDs and INDs $\Sigma = F \cup I$ is reduced and \mathbf{R} is in BCNF with respect to F , then, when the set of INDs I expresses pairwise consistency, F does not interact with I since I is proper circular. Consider a BCNF database schema \mathbf{R} with relation schemas $EMP = \{ENAME, DNAME\}$ and $DEPT = \{DNAME, LOCATION, MGR\}$, with $F = \{ENAME \rightarrow DNAME, DNAME \rightarrow \{LOCATION, MGR\}, MGR \rightarrow DNAME\}$ and $I = \{EMP[DNAME] \subseteq DEPT[DNAME], DEPT[DNAME] \subseteq EMP[DNAME]\}$. The set of INDs I expresses the fact that all employees work in departments that exist and all departments have at least one employee. The first IND is key-based, but the second is not. Despite this fact, it can be verified that F and I do not interact. Moreover, if managers are also employees, then we could add the IND $DEPT[MGR] \subseteq EMP[ENAME]$ to I and it can be verified by exhibiting the appropriate counterexamples that it is still true that F and I do not interact although now that I is not even proper circular. The database schema \mathbf{R} seems to be a reasonable design, but it is not in IDNF. Further research needs to be carried out to determine the semantics of normal forms for such FDs and INDs. We conclude the paper by proposing such a normal form. A database schema \mathbf{R} is in *Interaction Free Inclusion Dependency Normal Form* with respect to a set of Σ of FDs F and INDs I over R if:

1. \mathbf{R} is in BCNF with respect to F ,
2. All the INDs in I are either key-based or express pairwise consistency, and
3. F and I do not interact.

REFERENCES

- [1] P. Atzeni and V. De Antonellis, *Relational Database Theory*. Redwood City, Calif.: Benjamin/Cummings, 1993.
- [2] W.W. Armstrong, "Dependency Structures of Data Base Relationships," *Proc. IFIP Congress*, pp. 580-583, 1974.
- [3] C. Beeri and P.A. Bernstein, "Computational Problems Related to the Design of Normal Form Relational Schemas," *ACM Trans. Database Systems*, vol. 4, pp. 30-59, 1979.
- [4] J. Biskup and P. Dublish, "Objects in Relational Database Schemes with Functional, Inclusion and Exclusion Dependencies," *Theoretical Informatics and Applications*, vol. 27, pp. 183-219, 1993.
- [5] C. Beeri, R. Fagin, D. Maier, and M. Yannakakis, "On the Desirability of Acyclic Database Schemes," *J. ACM*, vol. 30, pp. 479-513, 1983.
- [6] P.A. Bernstein and N. Goodman, "What Does Boyce-Codd Normal Form Do?" *Proc. Int'l Conf. Very Large Data Bases*, pp. 245-259, 1980.
- [7] M.A. Casanova and J.E. Amaral de Sa, "Mapping Uninterpreted Schemes into Entity-Relationship Diagrams: Two Applications to Conceptual Schema Design," *IBM J. Research and Development*, vol. 28, pp. 82-94, 1984.
- [8] M.A. Casanova, R. Fagin, and C.H. Papadimitriou, "Inclusion Dependencies and Their Interaction with Functional Dependencies," *J. Computer and System Sciences*, vol. 28, pp. 29-59, 1984.
- [9] E.P.F. Chan, "A Design Theory for Solving the Anomalies Problem," *SIAM J. Computing*, vol. 18, pp. 429-448, 1989.
- [10] S.S. Cosmadakis and P.C. Kanellakis, "Functional and Inclusion Dependencies: A Graph Theoretic Approach," *Advances in Computing Research*, P.C. Kanellakis and F. Preparata, eds., vol. 3, pp. 163-184. Greenwich: JAI Press, 1986.
- [11] S.S. Cosmadakis, P.C. Kanellakis, and M.Y. Vardi, "Polynomial-Time Implication Problems for Unary Inclusion Dependencies," *J. ACM*, vol. 37, pp. 15-46, 1990.
- [12] E.F. Codd, "Recent Investigations in Relational Data Base Systems," *Proc. IFIP Congress*, pp. 1,017-1,021, 1974.
- [13] E.F. Codd, "Extending the Database Relational Model to Capture More Meaning," *ACM Trans. Database Systems*, vol. 4, pp. 397-434, 1979.
- [14] M.A. Casanova, L. Tucheran, and A.L. Furtado, "Enforcing Inclusion Dependencies and Referential Integrity," *Proc. Int'l Conf. Very Large Data Bases*, pp. 38-49, 1988.
- [15] M.A. Casanova and V.M.P. Vidal, "Towards a Sound View Integration Methodology," *Proc. ACM Symp. Principles of Database Systems*, pp. 36-47, 1983.
- [16] A.K. Chandra and M.Y. Vardi, "The Implication Problem for Functional and Inclusion Dependencies Is Undecidable," *SIAM J. Computing*, vol. 14, pp. 671-677, 1985.
- [17] C.J. Date, "Referential Integrity," *Relational Database: Selected Writings*, pp. 41-63. Reading, Mass.: Addison-Wesley, 1986.
- [18] R. Fagin, "Multivalued Dependencies and a New Normal Form for Relational Databases," *ACM Trans. Database Systems*, vol. 2, pp. 262-278, 1977.
- [19] R. Fagin, "Normal Forms and Relational Database Operators," *Proc. ACM SIGMOD Conf. Management of Data*, pp. 153-160, 1979.
- [20] R. Fagin, "A Normal Form for Relational Databases that Is Based on Domains and Keys," *ACM Trans. Database Systems*, vol. 6, pp. 387-415, 1981.
- [21] T. Imielinski, "Abstraction in Query Processing," *J. ACM*, vol. 38, pp. 534-558, 1991.
- [22] D.S. Johnson and A. Klug, "Testing Containment of Conjunctive Queries under Functional and Inclusion Dependencies," *J. Computer and System Sciences*, vol. 28, pp. 167-189, 1984.
- [23] T.-W. Ling and C.H. Goh, "Logical Database Design with Inclusion Dependencies," *Proc. Int'l Conf. Data Eng.*, pp. 642-649, 1992.
- [24] M. Levene and G. Loizou, "How to Prevent Interaction of Functional and Inclusion Dependencies," *Information Processing Letters*, vol. 71, pp. 115-125, 1995.
- [25] D. Maier, *The Theory of Relational Databases*. Rockville, Md.: Computer Science Press, 1983.
- [26] H. Mannila, "On the Complexity of the Inference Problem for Subclasses of Inclusion Dependencies," *Proc. Winter School on Theoretical Computer Science*, pp. 182-193, 1984.
- [27] J.C. Mitchell, "The Implication Problem for Functional and Inclusion Dependencies," *Information and Control*, vol. 56, pp. 154-173, 1983.
- [28] H. Mannila and K.-J. Rähkä, "Inclusion Dependencies in Database Design," *Proc. Int'l Conf. Data Eng.*, pp. 713-718, 1986.
- [29] H. Mannila and K.-J. Rähkä, *The Design of Relational Databases*. Reading, Mass.: Addison-Wesley, 1992.

- [30] E. Sciore, "Comparing the Universal Instance and Relational Data Models," *Advances in Computing Research*, P.C. Kanellakis and F. Preparata, eds., vol. 3, pp. 139-162. Greenwich: JAI Press, 1986.
- [31] J.D. Ullman, *Principles of Database and Knowledge-Base Systems*, vol. 1. Rockville, Md.: Computer Science Press, 1988.
- [32] M.W. Vincent, "Modification Anomalies and Boyce-Codd Normal Form," *Research and Practical Issues in Data Management*, B. Srinivasan and J. Zeleznikow, eds., pp. 251-264. Singapore: World Scientific, 1992.
- [33] M.W. Vincent, "The Semantic Justification for Normal Forms in Relational Database Design," PhD thesis, Dept. of Computer Science, Monash Univ., Melbourne, Australia, 1994.
- [34] M.W. Vincent, "Redundancy Elimination and a New Normal Form for Relational Databases," *Semantics in Databases*, B. Thalheim and L. Libkin, eds., pp. 247-264, Berlin: Springer Verlag, 1998.



interests are database theory and hypertext.



Mark Levene received his PhD degree in computer science in 1990 from Birkbeck College, which is part of the University of London. Dr. Levene is currently a senior lecturer in the Department of Computer Science at University College London, also part of the University of London. He has published extensively in the area of database theory and has recently coauthored a comprehensive book on relational databases and its extensions. His main research

Millist W. Vincent received his PhD from Monash University in 1994 for a dissertation in the area of database design. He is a senior lecturer in the School of Computer and Information Science at the University of South Australia. He has published widely in the areas of database design, database theory, and view maintenance, and his current interests are in the areas of database theory and view maintenance.