

# Relational Database Schema Design for Uncertain Data

Sebastian Link  
Department of Computer Science  
The University of Auckland  
New Zealand  
s.link@auckland.ac.nz

Henri Prade  
IRIT, CNRS and  
Université de Toulouse III  
France  
henri.prade@irit.fr

## ABSTRACT

We investigate the impact of uncertainty on relational database schema design. **Uncertainty is modeled qualitatively by assigning to tuples a degree of possibility with which they occur, and assigning to functional dependencies a degree of certainty which says to which tuples they apply.** A design theory is developed for possibilistic functional dependencies, including efficient axiomatic and algorithmic characterizations of their implication problem. Naturally, the possibility degrees of tuples result in a scale of different degrees of data redundancy. Scaled versions of the classical syntactic Boyce-Codd and Third Normal Forms are established and semantically justified in terms of avoiding data redundancy of different degrees. Classical decomposition and synthesis techniques are scaled as well. Therefore, **possibilistic functional dependencies do not just enable designers to control the levels of data integrity and losslessness targeted but also to balance the classical trade-off between query and update efficiency.** Extensive experiments confirm the efficiency of our framework and provide original insight into relational schema design.

## CCS Concepts

•Information systems → Database design and models; Uncertainty; •Theory of computation → Database constraints theory;

## Keywords

Axioms; Boyce-Codd normal form; Data redundancy; Implication problem; Third normal form; Possibility theory

## 1. INTRODUCTION

Relational databases were developed for applications with certain data, such as accounting, inventory, and payroll. Modern applications such as information extraction, data integration and cleaning, require techniques for uncertain data. Research on uncertain data has been prolific, yet

two trends can be observed: Queries are the dominant focus point, and uncertainty is mostly modeled quantitatively as probabilistic data. Here, the impact of uncertainty on database schema design is studied qualitatively, targeting the efficient processing of frequent queries and updates.

**In classical schema design, update inefficiencies are avoided by removing data redundancy caused by functional dependencies (FDs). Intuitively, if data is uncertain, then so is any redundancy that results from this data.** The more possible it is for data redundancy to occur, the more FDs can cause this redundancy, and the harder the normalization effort will be to remove that redundancy. In other words, the removal of redundancy from less possible data requires normalization for a smaller number of FDs. This is great news as data that are less possible are intuitively subject to more updates. Therefore, most of the frequent updates can be supported efficiently with less normalization effort. In turn, less normalization results in better query efficiency. Our findings will establish a practical and precise framework that enables database designers to take full advantage of this intuitive impact of uncertainty.

We model uncertainty by assigning to each tuple a degree of possibility (p-degree) by which it is perceived to occur. The possibility degree comes from any given finite scale of linearly ordered values. The framework meets well the intuition and ability of people to reason qualitatively, and not quantitatively. After all, humans like to classify items into a few categories and are uncomfortable with assigning exact values such as probabilities. **The framework results in a nested chain of possible worlds, each of which is a classical relation.** The smallest world contains only tuples that are fully possible while the largest world contains all tuples, excluding those which are impossible to occur. Moreover, only the smallest world is regarded as the part of the database which is certain. **The certainty by which an FD holds is derived from the p-degree of the smallest world in which the FD is violated.** As extremes, FDs that are fully certain even hold in the largest world, and FDs that are not certain at all do not even hold in the smallest world. Our contributions are as follows:

1. We formalize uncertainty by assigning p-degrees to tuples in a database. The approach is well-founded as it results in a possibility distribution over possible worlds that form a linear chain of relations.

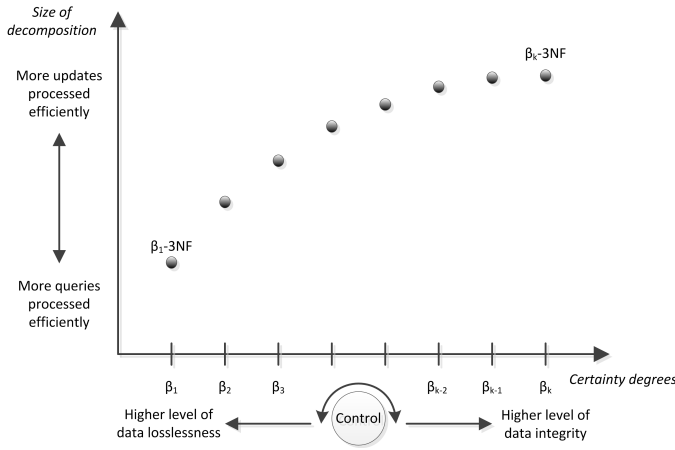
2. We define a possibilistic FD (pFD) as classical FDs with a degree of certainty (c-degree), derived from the p-degree of the smallest possible world in which it is violated.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from Permissions@acm.org.

CIKM'16, October 24–28, 2016, Indianapolis, IN, USA

© 2016 Copyright held by the owner/author(s). Publication rights licensed to ACM. ISBN 978-1-4503-4073-1/16/10...\$15.00

DOI: <http://dx.doi.org/10.1145/2983323.2983801>



**Figure 1: Exploring c-degrees as a mechanism to control levels of data integrity, data losslessness, query efficiency, and update efficiency**

3. We establish a full design theory for pFDs, including axioms and algorithms for their implication problem.

4. We apply the design theory to establish a new normalization framework, including scaled versions of Boyce-Codd (BCNF) and Third normal forms (3NF), their semantic justification in terms of dependency-preservation and removal of data redundancy, as well as normalization algorithms. Our techniques provide database designers with a full scale of normalized schemata, from which better-informed selections of the final schema can be made. We envision two main use-case scenarios: i) An organization primarily focused on data integrity and losslessness may choose the FDs that apply to the target c-degree. In this case, our algorithms compute the normal forms that exactly meet that target. ii) For an organization focused on efficiently processing its workload of queries and updates, a database designer can use our techniques to derive the normal forms that best match the workload. In that case, it is clear what levels of data integrity and losslessness are achieved. Hence, c-degrees provide an efficient mechanism to balance the trade-offs between data integrity and losslessness, and between query and update efficiency. This is illustrated in Figure 1, showing how the selection of different c-degrees  $\beta_i$  results in different decompositions, for example,  $\beta_i$ -3NF for  $i = 1, \dots, k$ . Our findings can justify classically de-normalized schema in terms of permitting data redundancy caused by pFDs whose c-degree is smaller than that targeted.

5. Extensive experiments with our algorithms confirm their efficiency in practice, and provide new insight into classical normalization trade-offs.

**Organization.** A motivating example with applications is introduced in Section 2. Related work is discussed in Section 3. Our data model of uncertainty is introduced in Section 4. PFDs are defined in Section 5, and their design theory is established in Section 6. Scaled versions of BCNF and 3NF are defined in Section 7, and their semantic justifications are derived. Normalization algorithms are established in Section 8. Experimental results are discussed in Section 9. Section 10 concludes and comments on future work. Proofs, more examples, algorithms and experiments can be found in the supplementary technical report [21].

<i>Project</i>	<i>Time</i>	<i>Manager</i>	<i>Room</i>	<i>Poss. degree</i>
Eagle	Mon, 9am	Ann	Aqua	$\alpha_1$
Hippo	Mon, 1pm	Ann	Aqua	$\alpha_1$
Kiwi	Mon, 1pm	Pete	Buff	$\alpha_1$
Kiwi	Tue, 2pm	Pete	Buff	$\alpha_1$
Lion	Tue, 4pm	Gill	Buff	$\alpha_1$
Lion	Wed, 9am	Gill	Cyan	$\alpha_1$
Lion	Wed, 11am	Bob	Cyan	$\alpha_2$
Lion	Wed, 11am	Jack	Cyan	$\alpha_3$
Lion	Wed, 11am	Pam	Lava	$\alpha_3$
Tiger	Wed, 11am	Pam	Lava	$\alpha_4$

**Table 1: Possibilistic relation**

## 2. APPLICATION AND MOTIVATION

As an example application, sufficiently simple to motivate our research and explain our findings, we consider an employee who extracts information from web-sites about weekly project meetings in her company. Attributes involve *Project*, storing projects with unique names, *Time*, for the weekday and start time, *Manager*, for the managers of the project that attend, and *Room*, for the unique name of a room. The employee classifies the possibility with which tuples occur in the relation according to the levels of trust associated with the data source. Tuples from the official web-site are assigned p-degree  $\alpha_1$ , indicating they are fully possible, tuples from a project manager's web-site get degree  $\alpha_2$ , tuples from a project member's web-site get degree  $\alpha_3$ , and tuples that originate from rumors are assigned p-degree  $\alpha_4$ . Implicitly, any other tuple has p-degree  $\alpha_5$ , indicating that it is impossible to occur. A different interpretation may result from already held meetings, confirmed meetings, requested meetings, planned meetings, and all other meetings. The p-degrees may have numerical interpretations, e.g.  $1 > 0.75 > 0.5 > 0.25 > 0$ . Either way, the employee has chosen 5 p-degrees to assign qualitative levels of uncertainty to tuples, with top degree  $\alpha_1$  and bottom degree  $\alpha_5$ . Table 1 shows a possibilistic relation (p-relation), that is, a relation where a p-degree is assigned to each tuple.

Naturally, the assignment of p-degrees results in a linearly ordered chain of possible worlds: For  $i = 1, \dots, 4$ , the relation  $r_i$  consists of tuples with p-degree  $\alpha_i$  or higher, i.e.  $\alpha_j$  with  $j \leq i$ . The p-degree of world  $r_i$  is  $\alpha_i$ . In particular, fully possible tuples occur in every possible world, and are therefore also fully certain to occur. The possible worlds of the p-relation in Table 1 are illustrated in Figure 2.

A given FD is either satisfied by the largest possible world or there is a smallest possible world in which it is violated. For example, the FD  $\sigma_1 : \text{Manager}, \text{Time} \rightarrow \text{Room}$  is satisfied by the world  $r_4$ , and thus holds in every possible world. Consequently, it is assigned the top c-degree, denoted by  $\beta_1$ . The smallest relation that violates  $\sigma_2 : \text{Room}, \text{Time} \rightarrow \text{Project}$  is  $r_4$ , that is, the FD is assigned the second highest c-degree,  $\beta_2$ . The smallest relation that violates  $\sigma_3 : \text{Project}, \text{Time} \rightarrow \text{Manager}$  is  $r_3$ , that is, the FD is assigned the third highest c-degree,  $\beta_3$ . The smallest relation that violates  $\sigma_4 : \text{Project} \rightarrow \text{Manager}$  is  $r_2$ , and the FD thus holds with c-degree  $\beta_4$ . The FD  $\text{Manager}, \text{Room} \rightarrow \text{Time}$  is violated even by the smallest possible world  $r_1$ , and is thus assigned the bottom c-degree  $\beta_5 = \beta_{k+1}$ . Hence, the p-degree  $\alpha_i$  of the smallest possible world  $r_i$  in which the FD is violated, determines the c-degree

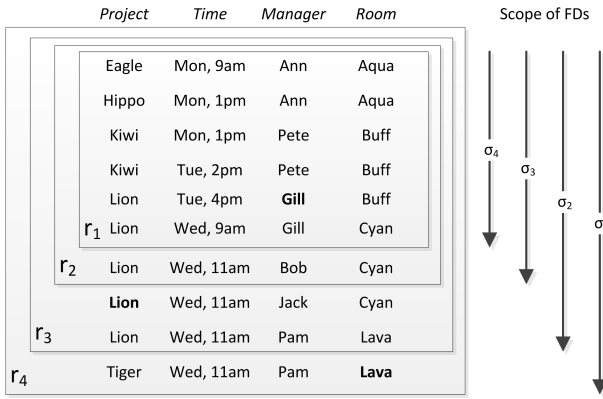


Figure 2: Worlds of p-relation from Table 1 with some redundant data values in bold and FD scopes

$\beta_{k+2-i}$  with which the FD holds. The possible worlds on which the FDs  $\sigma_1, \dots, \sigma_4$  hold, i.e. their scope, are illustrated in Figure 2.

Clearly, the higher the possibility of a possible world the less tuples it contains and the more FDs may hold on it, and therefore, the more data occurrences may be redundant. Therefore, whenever we want to eliminate redundant data occurrences in more possible data, then we must expect the normalization effort to grow. Redundancy in the smallest possible world  $r_1$ , for example, could be caused by FDs that hold with any c-degree, except the bottom degree, while redundancy in the largest possible world  $r_4$  can only be caused by FDs with the top c-degree, see Figure 2. In our running example, the bold occurrence of *Gill* in  $r_1$  is redundant for the FD  $\sigma_4 : Project \rightarrow Manager$  that holds with degree  $\beta_4$ . Indeed, any replacement of this data occurrence by a different value will result in a c-degree of  $\sigma_4$  that is smaller than  $\beta_4$ , with which it must hold. This data value occurrence is therefore  $\alpha_1$ -redundant, and could be caused by any pFD in  $\Sigma$  - all of which apply to tuples with p-degree  $\alpha_1$ . The bold occurrence of the data value *Lion* in  $r_3$  is  $\alpha_3$ -redundant. It could be caused by any FD in  $\Sigma$  with c-degree  $\beta_1$  or  $\beta_2$ , in this case by the FD  $\sigma_2 : Room, Time \rightarrow Project$  with c-degree  $\beta_2$ . Finally, the bold occurrence of the data value *Lava* in  $r_4$  is  $\alpha_4$ -redundant. Its redundancy can only be caused by FDs that hold with c-degree  $\beta_1$ , in this case by  $\sigma_1 : Manager, Time \rightarrow Room$ .

Therefore, the impact we expect uncertainty to have on schema design is that the more frequent updates are likely to be processed efficiently with less normalization effort. This means that more queries can be processed efficiently, as less joins are required. Uncertainty results in a greater variety of normalized schema designs from which the best match for the target workload of updates and queries can be chosen. In other words, uncertainty can be understood as an effective mechanism to better control the classical trade-off between update and query efficiency, see Figure 1. Figure 3 illustrates how different Boyce-Codd normal form decompositions apply to the p-relation of our running example.

### 3. RELATED WORK

Probabilistic databases have received much interest [25] due to the need to deal with uncertain data. Constraints present a key challenge here: “When the data is uncertain,

constraints can be used to increase the quality of the data, and hence they are an important tool in managing data with uncertainties” [8]. Suciu et al. emphasize that “the main use of probabilities is to record the degree of uncertainty in the data and to rank the outputs to a query; in some applications, the exact output probabilities matter less to the user than the ranking of the outputs” [25]. This suggests that a qualitative approach to uncertainty, such as possibility theory [9], can avoid the high computational complexity in obtaining and processing probabilities, while guaranteeing the same qualitative outcome.

Research on probabilistic databases has naturally focused on query processing, e.g. [3, 13]. Common is the desire to extend trusted relational technology to handle uncertainty. This desire is also inherent in our approach: We show how to exploit relational normalization to design database schemata for applications with uncertain data.

No previous research has taken advantage of uncertainty information to introduce different degrees of data redundancy. These different degrees of data redundancy enable us to introduce a whole variety of normal forms each targeting different workloads and different levels of data integrity and losslessness. The linear order of our p-degrees is the counterpart of the linear order of probabilities. Extensions to partial orders are possible, but not our focus.

A companion paper [22] discusses the origins of our research in possibility theory, contains an up-to-date review of FDs for uncertain data, and characterizes pFD implication by that of Horn clauses in possibilistic logic. Design theory, normal forms, normalization, redundancy, and experiments are new contributions of the current article. Keys and cardinality constraints have also been investigated in the context of our possibilistic model [12, 14].

Few papers address schema design for uncertain data [5, 24]. Das Sarma, Ullman and Widom develop an “FD theory for data models whose basic construct for uncertainty is *alternatives*” [24]. Their work is thus fundamentally different from our approach. This is also true for Chaudhry, Moyne and Rundensteiner who model fuzziness in an extended Entity-Relationship model [5].

Finally, classical findings on relational schema design enabled us to develop our approach towards schema design for uncertain data. FDs are already mentioned in the seminal paper by Codd [6]. Armstrong axiomatized FDs [1], and linear-time algorithms to decide their implication problem are known [2]. These enable relational schema design, including 3NF [4], BCNF [7] and their semantic justification in terms of dependency-preservation and elimination of data redundancy [19, 27]. Deciding whether a given schema is in 3NF is NP-complete, and in P-time for BCNF, but deciding whether a projection of a given schema is in BCNF is coNP-complete [2]. Tsou and Fischer gave an algorithm to compute a lossless BCNF decomposition in P-time [26]. Our experiments reveal original insight into the trade-off between classical 3NF synthesis and BCNF decomposition.

### 4. UNCERTAIN DATABASES

We summarize the possibilistic data model from [22]. The main motivation for this data model is the schema design framework we develop in the current article.

A relation schema is a finite non-empty set  $R$  of attributes. Each attribute  $A \in R$  has a domain  $dom(A)$  of values. A tuple  $t$  over  $R$  is an element of the Cartesian product

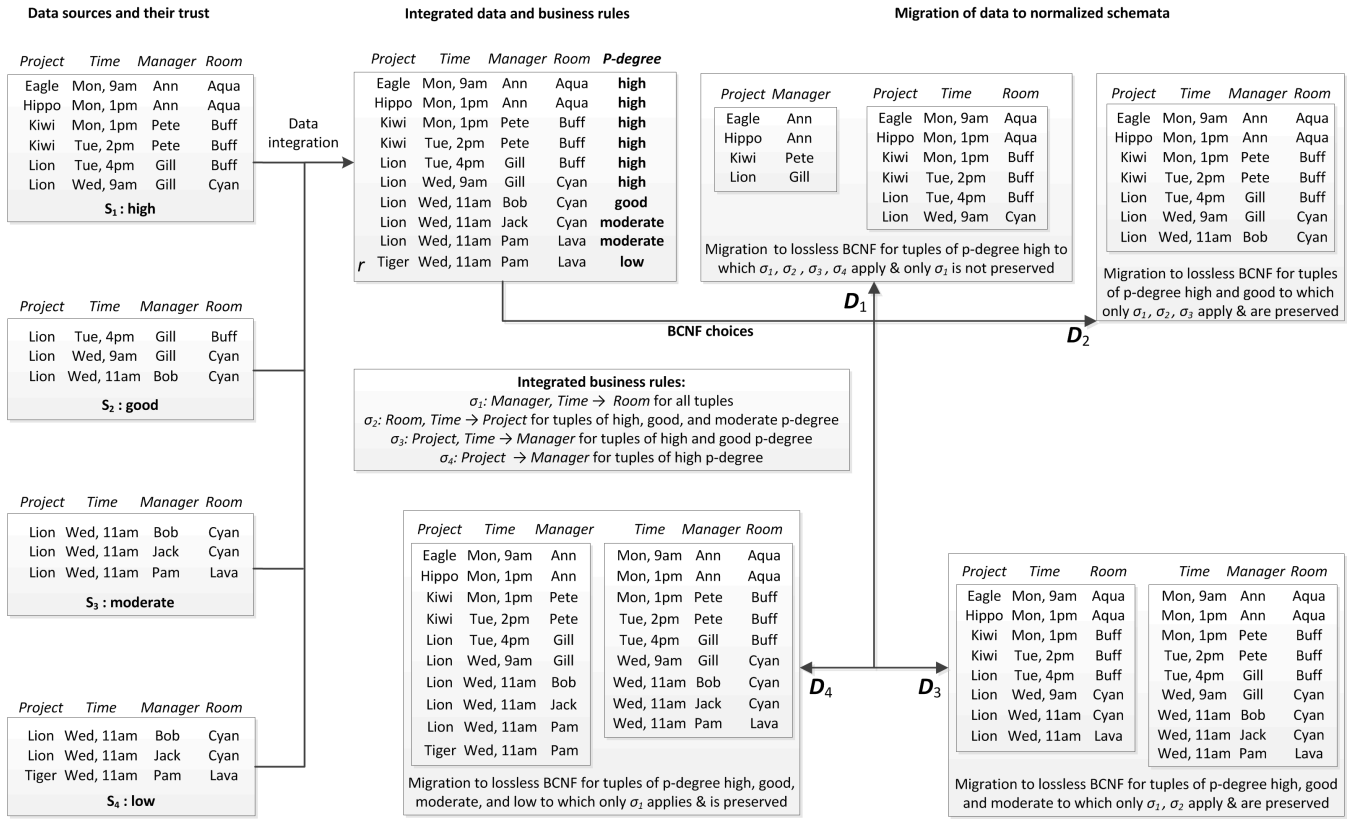


Figure 3: Using degrees of uncertainty for data integration and migration into normalized schemata

$\prod_{A \in R} \text{dom}(A)$ . For  $X \subseteq R$  we denote by  $t(X)$  the projection of  $t$  on  $X$ . A relation over  $R$  is a finite set  $r$  of tuples over  $R$ . As example we use the relation schema  $\text{MEETING} = \{\text{Project}, \text{Time}, \text{Manager}, \text{Room}\}$  from Section 2.

Tuples either belong to a relation or they do not, there is no room for uncertainty. For example, we cannot express less confidence for Bob attending a meeting of project Lion on Wednesday at 11am in room Cyan than for Gill attending this meeting. Database support for new applications, such as data integration, requires us to make the most out of available uncertainty information.

We define uncertain relations as relations where each tuple is associated with some confidence. The confidence expresses with which degree of possibility a tuple occurs in a relation. Formally, we model the confidence as a *scale of possibility*, i.e., a finite strict linear order  $\mathcal{S} = (S, <)$  with  $k+1$  elements where  $k$  is some positive integer, which we denote by  $\alpha_1 > \dots > \alpha_k > \alpha_{k+1}$ . The elements  $\alpha_i$  are *possibility degrees* (p-degrees). The top p-degree  $\alpha_1$  is reserved for tuples that are ‘fully possible’ while the bottom p-degree  $\alpha_{k+1}$  is for tuples that are ‘not possible at all’, that is ‘impossible’, to occur in a relation. Humans like to use simple scales in everyday life, for instance to communicate, compare, or rank. Simple usually means to classify items qualitatively, rather than quantitatively by putting a precise value on it. Classical relations use a scale with two elements, i.e. where  $k = 1$ . As a running example, we will use a scale of five p-degrees, ranging from *fully possible* ( $\alpha_1$ ), to *rather possible* ( $\alpha_2$ ), to *medium possible* ( $\alpha_3$ ), to *somewhat possible* ( $\alpha_4$ ), and finally to *not possible at all* ( $\alpha_5$ ).

Formally, a *possibilistic relation schema* (p-schema)  $(R, \mathcal{S})$  consists of a relation schema  $R$  and a possibility scale  $\mathcal{S}$ . A *possibilistic relation* (p-relation) over  $(R, \mathcal{S})$  consists of a relation  $r$  over  $R$ , together with a function  $\text{Poss}_r$  that maps each tuple  $t \in r$  to a p-degree  $\text{Poss}_r(t)$  in the possibility scale  $\mathcal{S}$ . For example, Table 1 shows a p-relation over  $(\text{MEETING}, \mathcal{S} = \{\alpha_1, \dots, \alpha_5\})$ .

P-relations enjoy a well-founded semantics in terms of possible worlds. In fact, a p-relation gives rise to a possibility distribution over possible worlds of relations. For  $i = 1, \dots, k$  let  $r_i$  denote the relation that consists of all tuples in  $r$  that have a p-degree of at least  $\alpha_i$ , i.e.,  $r_i = \{t \in r \mid \text{Poss}_r(t) \geq \alpha_i\}$ . The linear order of the p-degrees results in a (reversed) linear order of possible worlds of relations. Indeed, we have  $r_1 \subseteq r_2 \subseteq \dots \subseteq r_k$ . The possibility distribution  $\pi_r$  for this linear chain of possible worlds is defined by  $\pi_r(r_i) = \alpha_i$ . Note that  $r_{k+1}$  is not considered to be a possible world, since its possibility  $\pi(r_{k+1}) = \alpha_{k+1}$  means ‘not possible at all’. Vice versa, the possibility  $\text{Poss}_r(t)$  of a tuple  $t \in r$  is the possibility of the smallest possible world in which  $t$  occurs, i.e., the maximum possibility  $\max\{\alpha_i \mid t \in r_i\}$  of a world to which  $t$  belongs. If  $t \notin r_k$ , then  $\text{Poss}_r(t) = \alpha_{k+1}$ . The top p-degree  $\alpha_1$  takes on a distinguished role: every tuple that is ‘fully possible’ occurs in every possible world - and is thus - ‘fully certain’. This formally confirms our intuition that uncertain relations subsume relations (of fully certain tuples) as a special case. Figure 2 shows the possible worlds  $r_1 \subsetneq r_2 \subsetneq r_3 \subsetneq r_4$  of the p-relation in Table 1.

## 5. POSSIBILISTIC FDS

We introduce possibilistic FDs (pFDs) as classical FDs with a degree of certainty. In the same way FDs are fundamental to relational database design, pFDs will play a fundamental role to schema design for uncertain data.

Recall that an FD  $X \rightarrow Y$  over relation schema  $R$  is satisfied by a relation  $r$  over  $R$  whenever every pair of tuples in  $r$  that have matching values on all the attributes in  $X$  have also matching values on all the attributes in  $Y$ . For example, the FD  $\sigma_5 : \text{Manager}, \text{Room} \rightarrow \text{Time}$  is not satisfied by any relation  $r_1, \dots, r_4$ . The FD  $\sigma_4 : \text{Project} \rightarrow \text{Manager}$  is satisfied by  $r_1$ , but not by  $r_2$  and therefore not by  $r_3$  and  $r_4$ . The FD  $\sigma_3 : \text{Project}, \text{Time} \rightarrow \text{Manager}$  is satisfied by  $r_1$  and  $r_2$ , but not by  $r_3$  and therefore not by  $r_4$ . The FD  $\sigma_2 : \text{Time}, \text{Room} \rightarrow \text{Project}$  is satisfied by  $r_1, r_2$ , and  $r_3$ , but not by  $r_4$ . Finally, the FD  $\sigma_1 : \text{Manager}, \text{Time} \rightarrow \text{Room}$  is satisfied by all relations  $r_1, \dots, r_4$ .

Naturally, the p-degrees of tuples that define an uncertain relation also define degrees of certainty with which FDs hold in the uncertain relation. Intuitively, since  $\sigma_1$  is satisfied in every possible world, it is fully certain ( $\beta_1$ ) to hold in  $r$ . As  $\sigma_2$  is only violated in a somewhat possible world  $r_4$ , it is quite certain ( $\beta_2$ ) to hold in  $r$ . Since  $\sigma_3$  is only violated in a medium possible world  $r_3$ , it is medium certain ( $\beta_3$ ) to hold in  $r$ . As  $\sigma_4$  is only violated in a quite possible world  $r_2$ , it is somewhat certain ( $\beta_4$ ) to hold in  $r$ . Finally, as  $\sigma_5$  is violated in the fully possible world  $r_1$ , it is not certain at all ( $\beta_5$ ) to hold in  $r$ .

In summary, the certainty with which an FD holds in an uncertain relation corresponds to the possibility of the smallest possible world in which the FD is violated. Therefore, similar to a scale  $\mathcal{S}$  of possibility degrees for tuples we use a scale  $\mathcal{S}^T$  of certainty degrees (c-degrees) for FDs. We commonly use subscripted versions of the Greek letter  $\beta$  to denote c-degrees associated with FDs. Formally, the correspondence between p-degrees in  $\mathcal{S}$  and the c-degrees in  $\mathcal{S}^T$  can be defined by the mapping  $\alpha_i \mapsto \beta_{k+2-i}$  for  $i = 1, \dots, k+1$ . Formally then, the certainty  $C_r(X \rightarrow Y)$  with which the FD  $X \rightarrow Y$  holds in the uncertain relation  $r$  is the top degree  $\beta_1$  whenever  $X \rightarrow Y$  is satisfied by  $r_k$ , or otherwise the minimum amongst the c-degrees  $\beta_{k+2-i}$  that correspond to possible worlds  $r_i$  in which  $X \rightarrow Y$  is violated, that is,  $C_r(X \rightarrow Y)$

$$= \begin{cases} \beta_1 & , \text{ if } r_k \text{ satisfies } X \rightarrow Y \\ \min\{\beta_{k+2-i} \mid \not\models_{r_i} X \rightarrow Y\} & , \text{ otherwise} \end{cases}.$$

We can now define the syntax and semantics of pFDs.

**DEFINITION 1.** A possibilistic FD (pFD) over a p-schema  $(R, \mathcal{S})$  is an expression  $(X \rightarrow Y, \beta)$  where  $X, Y \subseteq R$  and  $\beta \in \mathcal{S}^T$ . A p-relation  $(r, \text{Poss}_r)$  over  $(R, \mathcal{S})$  is said to satisfy the pFD  $(X \rightarrow Y, \beta)$  if and only if  $C_r(X \rightarrow Y) \geq \beta$ .

The p-relation  $(r, \text{Poss}_r)$  of our running example satisfies:  $(\sigma_1, \beta_1)$ ,  $(\sigma_2, \beta_2)$ ,  $(\sigma_3, \beta_3)$ ,  $(\sigma_4, \beta_4)$ , and  $(\sigma_5, \beta_5)$ . It violates the pFD  $(\sigma_3, \beta_2)$  since  $C_r(\sigma_3) = \beta_3 < \beta_2$ .

## 6. SCALED DESIGN THEORY

Classical normalization is founded on the theory of FDs. Consequently, we now establish a design theory of pFDs as a foundation for normalizing uncertain data.

### 6.1 The Magic of $\beta$ -Cuts

We link the implication of pFDs and FDs. Let  $\Sigma \cup \{\varphi\}$  denote a set of pFDs over a p-schema  $(R, \mathcal{S})$ . We say that  $\Sigma$  implies  $\varphi$ , denoted by  $\Sigma \models \varphi$ , if every p-relation  $(r, \text{Poss}_r)$  over  $(R, \mathcal{S})$  that satisfies every pFD in  $\Sigma$  also satisfies  $\varphi$ .

**EXAMPLE 1.** Let  $\Sigma$  consist of the following four pFDs  $(\text{Manager}, \text{Time} \rightarrow \text{Room}, \beta_1)$ ,  $(\text{Room}, \text{Time} \rightarrow \text{Project}, \beta_2)$ ,  $(\text{Project}, \text{Time} \rightarrow \text{Manager}, \beta_3)$ ,  $(\text{Project} \rightarrow \text{Manager}, \beta_4)$  over  $(\text{MEETING}, \{\alpha_1, \dots, \alpha_5\})$ . Further, let  $\varphi$  denote the pFD  $(\text{Room}, \text{Time} \rightarrow \text{Manager}, \beta_2)$ . Then  $\Sigma$  does not imply  $\varphi$  as the following p-relation witnesses:

Project	Time	Manager	Room	Poss. degree
Lion	Wed, 3pm	Gill	Cyan	$\alpha_1$
Lion	Wed, 3pm	Robert	Cyan	$\alpha_3$

For a set  $\Sigma$  of pFDs on some p-schema  $(R, \mathcal{S})$  and c-degree  $\beta \in \mathcal{S}^T$  where  $\beta > \beta_{k+1}$ , let

$$\Sigma_\beta = \{X \rightarrow Y \mid (X \rightarrow Y, \beta') \in \Sigma \text{ and } \beta' \geq \beta\}$$

be the  $\beta$ -cut of  $\Sigma$ . The major strength of our framework is engraved in the following result.

**THEOREM 1.** Let  $\Sigma \cup \{(X \rightarrow Y, \beta)\}$  be a set of pFDs over a p-schema  $(R, \mathcal{S})$  where  $\beta > \beta_{k+1}$ . Then  $\Sigma \models (X \rightarrow Y, \beta)$  if and only if  $\Sigma_\beta \models X \rightarrow Y$ .  $\square$

**EXAMPLE 2.** Let  $\Sigma$  and  $\varphi$  be as in Example 1, in particular  $\Sigma$  does not imply  $\varphi$ . Theorem 1 says that  $\Sigma_{\beta_2}$  does not imply  $\text{Room}, \text{Time} \rightarrow \text{Manager}$ . Indeed, the possible world  $r_3$  of the p-relation from Example 1

Project	Time	Manager	Room
Lion	Wed, 3pm	Gill	Cyan
Lion	Wed, 3pm	Robert	Cyan

satisfies the two classical FDs  $\text{Manager}, \text{Time} \rightarrow \text{Room}$  and  $\text{Room}, \text{Time} \rightarrow \text{Project}$  that form  $\Sigma_{\beta_2}$ , and violates the FD  $\text{Room}, \text{Time} \rightarrow \text{Manager}$ .

### 6.2 Scaled Armstrong Axioms

The semantic closure  $\Sigma^* = \{\varphi \mid \Sigma \models \varphi\}$  contains all dependencies implied by  $\Sigma$ . One can utilize a syntactic approach to compute  $\Sigma^*$  by applying inference rules of the form  $\frac{\text{premise}}{\text{conclusion}}$  condition, where rules without a premise are called *axioms*. For a set  $\mathfrak{R}$  of inference rules let  $\Sigma \vdash_{\mathfrak{R}} \varphi$  denote the inference of  $\varphi$  from  $\Sigma$  by  $\mathfrak{R}$ . That is, there is some sequence  $\sigma_1, \dots, \sigma_n$  such that  $\sigma_n = \varphi$  and every  $\sigma_i$  is an element of  $\Sigma$  or is the conclusion that results from an application of an inference rule in  $\mathfrak{R}$  to some premises in  $\{\sigma_1, \dots, \sigma_{i-1}\}$ . Let  $\Sigma_{\mathfrak{R}}^+ = \{\varphi \mid \Sigma \vdash_{\mathfrak{R}} \varphi\}$  denote the syntactic closure of  $\Sigma$  under inferences by  $\mathfrak{R}$ .  $\mathfrak{R}$  is sound (complete) if for every p-schema  $(R, \mathcal{S})$  and for every set  $\Sigma$  we have  $\Sigma_{\mathfrak{R}}^+ \subseteq \Sigma^*$  ( $\Sigma^* \subseteq \Sigma_{\mathfrak{R}}^+$ ). The (finite) set  $\mathfrak{R}$  is said to be a (finite) axiomatization if  $\mathfrak{R}$  is both sound and complete.

PFDs enjoy the axiomatization  $\mathfrak{P}$  from Table 2. It subsumes Armstrong's axioms [1] for FDs as the special case where the scale  $\mathcal{S}^T$  consists of just two c-degrees. In these rules all attribute sets  $X, Y, Z$  are subsets of the arbitrarily given relation schema  $R$ , and the c-degrees  $\beta$  and  $\beta'$  are elements of the arbitrarily given certainty scale  $\mathcal{S}^T$ . In particular,  $\beta_{k+1}$  denotes the bottom c-degree.



$\frac{}{(XY \rightarrow Y, \beta)}$ (reflexivity, $\mathcal{R}$ )	$\frac{(X \rightarrow Y, \beta)}{(X \rightarrow XY, \beta)}$ (extension, $\mathcal{E}$ )
$\frac{(X \rightarrow Y, \beta) \quad (Y \rightarrow Z, \beta)}{(X \rightarrow Z, \beta)}$ (transitivity, $\mathcal{T}$ )	
$\frac{}{(X \rightarrow Y, \beta_{k+1})}$ (bottom, $\mathcal{B}$ )	$\frac{(X \rightarrow Y, \beta)}{(X \rightarrow Y, \beta')} \beta' < \beta$ (weakening, $\mathcal{W}$ )

**Table 2: Axiomatization  $\mathfrak{P} = \{\mathcal{R}, \mathcal{E}, \mathcal{T}, \mathcal{B}, \mathcal{W}\}$  of pFDs**

**THEOREM 2.** *The set  $\mathfrak{P} = \{\mathcal{R}, \mathcal{E}, \mathcal{T}, \mathcal{B}, \mathcal{W}\}$  forms a finite axiomatization for the implication of pFDs.  $\square$*

**EXAMPLE 3.** *Let  $\Sigma$  and  $\varphi$  be as in Example 1, and  $\varphi' = (\text{Room, Time} \rightarrow \text{Manager}, \beta_3)$ . We show an inference of  $\varphi'$  from  $\Sigma$  by  $\mathfrak{P}$ . Attributes are abbreviated by their first letters.*

$$\frac{\frac{(RT \rightarrow P, \beta_2)}{\mathcal{W} : (RT \rightarrow P, \beta_3)} \quad \frac{\mathcal{E} : (RT \rightarrow RTP, \beta_3) \quad \mathcal{R} : (RTP \rightarrow PT, \beta_3)}{\mathcal{T} : (RT \rightarrow PT, \beta_3)}}{\mathcal{T} : (RT \rightarrow M, \beta_3)} \quad (PT \rightarrow M, \beta_3)$$

Of course,  $\varphi$  cannot be inferred from  $\Sigma$  by  $\mathfrak{P}$ , as Example 1 and the soundness of  $\mathfrak{P}$  show.

### 6.3 Scaled Decision Algorithm

In practice it is often unnecessary to compute  $\Sigma^*$ . Rather frequent is the problem of deciding if  $\Sigma$  implies  $\varphi$ .

PROBLEM: IMPLICATION	
INPUT:	Relation schema $R$ , Scale $\mathcal{S}$ with $k + 1$ possibility degrees, Set $\Sigma \cup \{\varphi\}$ of pFDs over $(R, \mathcal{S})$
OUTPUT:	Yes, if $\Sigma \models \varphi$ , and No, otherwise

Computing  $\Sigma^*$  to check if  $\varphi \in \Sigma^*$  is hardly efficient nor makes use of the input  $\varphi$ . We can exploit Theorem 1 to derive an algorithm that decides the implication problem for pFDs in time linear in the input. Given a pFD set  $\Sigma \cup \{(X \rightarrow Y, \beta)\}$  we return *true* if  $\beta = \beta_{k+1}$ , otherwise we check if  $\Sigma_\beta \models X \rightarrow Y$ . The latter test can be done in linear time by computing the attribute set closure  $X_{\Sigma_\beta}^+ = \{A \in R \mid \Sigma \vdash_{\mathfrak{A}} X \rightarrow A\}$  of  $X$  for  $\Sigma_\beta$  [2]. We use  $||\Sigma||$  for the total number of attributes in  $\Sigma$ , and  $|X|$  for the cardinality of  $X$ . We obtain the following result.

**THEOREM 3.** *The implication problem  $\Sigma \models \varphi$  of pFDs can be decided in time  $\mathcal{O}(|\Sigma \cup \{\varphi\}|)$ .  $\square$*

**EXAMPLE 4.** *Let  $\Sigma$  and  $\varphi$  be as in Example 1, and  $\varphi' = (\text{Room, Time} \rightarrow \text{Manager}, \beta_3)$ . It follows that  $\varphi$  is not implied by  $\Sigma$  as  $\text{CLOSURE}(\text{MEETING}, RT, \Sigma_{\beta_2}) = RTP$ . Further,  $\varphi'$  is implied by  $\Sigma$  as  $\text{CLOSURE}(\text{MEETING}, RT, \Sigma_{\beta_3}) = RTPM$ .*

## 7. SCALED NORMAL FORMS

In relational databases, Boyce-Codd normal form (BCNF) syntactically characterizes relation schemata that are guaranteed to be **free of data redundancy** in any relations over the schema, in terms of FDs [27]. Third normal form (3NF) syntactically characterizes relation schemata that are guaranteed to have the **least amount of data redundancy** in their relations amongst all schemata on which all FDs can be enforced locally [19]. In relations, data redundancy is treated uniformly for all data, and the elimination of all data redundancy requires normalization for all FDs.

In p-relations, different data can occur with different p-degrees. Therefore, **data redundancy may occur with different p-degrees, too**. Intuitively, the smaller the p-degree for which data redundancy is to be eliminated, the smaller the normalization effort will be. We will introduce notions of data redundancy that are tailored to the p-degree of tuples in which they occur. This results in a whole range of semantic normal forms by which data redundancy of growing p-degrees are eliminated. We then characterize each of these semantic normal forms by a corresponding syntactic normal form, and establish strong correspondences with BCNF and 3NF in relational databases.

### 7.1 Scaled Data Redundancy

Building on Section 2 we now introduce notions of data redundancy that are tailored to the p-degree of tuples. We exploit the classical proposal by Vincent [27]. Let  $R$  denote a relation schema,  $A$  an attribute of  $R$ ,  $t$  a tuple over  $R$ , and  $\Sigma$  a set of FDs over  $R$ . A *replacement* of  $t(A)$  is a tuple  $\bar{t}$  over  $R$  such that: i) for all  $\bar{A} \in R - \{A\}$  we have  $\bar{t}(\bar{A}) = t(\bar{A})$ , and ii)  $\bar{t}(A) \neq t(A)$ . For a relation  $r$  over  $R$  that satisfies  $\Sigma$  and  $t \in r$ , the data value occurrence  $t(A)$  in  $r$  is *redundant* for  $\Sigma$  iff for every replacement  $\bar{t}$  of  $t(A)$ ,  $\bar{r} := (r - \{t\}) \cup \{\bar{t}\}$  violates some FD in  $\Sigma$ . A relation schema  $R$  is in *Redundancy-Free normal form* (RFNF) for a set  $\Sigma$  of FDs iff there are no relation  $r$  over  $R$  that satisfies  $\Sigma$ , tuple  $t \in r$ , and attribute  $A \in R$  such that the data value occurrence  $t(A)$  is redundant for  $\Sigma$  [27].

**DEFINITION 2.** *Let  $(R, \mathcal{S})$  denote a p-schema,  $\Sigma$  a set of pFDs over  $(R, \mathcal{S})$ ,  $A \in R$  an attribute,  $(r, \text{Poss}_r)$  a p-relation over  $(R, \mathcal{S})$  that satisfies  $\Sigma$ , and  $t$  a tuple in  $r_i$ . The data value occurrence  $t(A)$  is  $\alpha_i$ -redundant iff  $t(A)$  is redundant for  $\Sigma_{\alpha_i} = \{X \rightarrow Y \mid (X \rightarrow Y, \beta) \in \Sigma \text{ and } \beta \geq \beta_{k+1-i}\}$ .  $\square$*

In our motivating example the occurrences of *Lava*, *Lion*, and *Gill* are  $\alpha_4$ -,  $\alpha_3$ - and  $\alpha_1$ -redundant, respectively. Importantly,  $\alpha_i$ -redundant data value occurrences can only be caused by pFDs  $(X \rightarrow Y, \beta)$  that apply to the world of the occurrence, i.e. where  $\beta \geq \beta_{k+1-i}$ . Hence,  $\alpha_1$ -redundancy can be caused by pFDs with any c-degree  $\beta_1, \dots, \beta_k$ , while  $\alpha_k$ -redundancy can only be caused by pFDs with c-degree  $\beta_1$ . We have arrived at the following definition.

**DEFINITION 3.** *A p-schema  $(R, \mathcal{S})$  is in  $\alpha_i$ -Redundancy-Free Normal Form ( $\alpha_i$ -RFNF) for a set  $\Sigma$  of pFDs over  $(R, \mathcal{S})$  iff there do not exist a p-relation  $(r, \text{Poss}_r)$  over  $(R, \mathcal{S})$  that satisfies  $\Sigma$ , an attribute  $A \in R$ , and a tuple  $t \in r_i$  such that  $t(A)$  is  $\alpha_i$ -redundant.  $\square$*

$(\text{MEETING}, \mathcal{S})$  is not in  $\alpha_4$ -RFNF,  $\alpha_3$ -RFNF, nor  $\alpha_1$ -RFNF, but it is in  $\alpha_2$ -RFNF for  $\Sigma$ . The negative results follows

from Figure 2, but the satisfaction of the  $\alpha_2$ -RFNF condition is not obvious. The next result shows that  $\alpha_i$ -RFNF characterizes p-schemata that permit only p-relations whose possible world  $r_i$  is free from data redundancy caused by the classical FDs that apply to it.

**THEOREM 4.**  *$(R, S)$  is in  $\alpha_i$ -RFNF for  $\Sigma$  if and only if  $R$  is in RFNF for  $\Sigma_{\alpha_i}$ .*

## 7.2 Scaled BCNF

We now characterize  $\alpha$ -RFNF, which is a semantic normal form, purely syntactically. This is achieved by scaling the classical BCNF condition. A relation schema  $R$  is in **Boyce-Codd normal form (BCNF)** for an FD set  $\Sigma$  over  $R$  iff for all  $X \rightarrow Y \in \Sigma_{\alpha}^+$  where  $Y \not\subseteq X$ , we have  $X \rightarrow R \in \Sigma_{\alpha}^+$ . Here,  $\Sigma_{\alpha}^+$  is the syntactic closure of  $\Sigma$  for the Armstrong axioms  $\mathfrak{A}$  [1]. While  $\alpha$ -RFNF is defined in terms of the semantics in the possible world with p-degree  $\alpha$ , scaled BCNF is defined in terms of the syntax of the given pFDs with their c-degrees.

**DEFINITION 4.** A p-schema  $(R, S)$  is in  $\beta$ -Boyce-Codd Normal Form for a set  $\Sigma$  of pFDs over  $(R, S)$  iff for every pFD  $(X \rightarrow Y, \beta) \in \Sigma_{\beta}^+$  where  $Y \not\subseteq X$ , we have  $(X \rightarrow R, \beta) \in \Sigma_{\beta}^+$ .  $\square$

Sets  $\Sigma$  and  $\Theta$  are covers for one another iff  $\Sigma^* = \Theta^*$  holds. Being in  $\beta$ -BCNF for  $\Sigma$  is independent of a cover for  $\Sigma$ . That is, for any cover  $\Sigma'$ ,  $(R, S)$  is in  $\beta$ -BCNF for  $\Sigma$  iff  $(R, S)$  is in  $\beta$ -BCNF for  $\Sigma'$ . The  $\beta$ -BCNF condition on a pFD set  $\Sigma$  is equivalent to the BCNF condition on the FD set  $\Sigma_{\beta}$ .

**THEOREM 5.**  *$(R, S)$  is in  $\beta$ -BCNF for a set  $\Sigma$  if and only if  $R$  is in BCNF for  $\Sigma_{\beta}$ .*

We are now in a position to characterize the semantic  $\alpha_i$ -RFNF by the syntactic  $\beta_{k+1-i}$ -BCNF.

**THEOREM 6.** *For  $i = 1, \dots, k$ ,  $(R, S)$  with  $|S| = k + 1$  is in  $\alpha_i$ -RFNF for  $\Sigma$  iff  $(R, S)$  is in  $\beta_{k+1-i}$ -BCNF for  $\Sigma$ .*

One may wonder about the efficiency of checking if a given p-schema  $(R, S)$  is in  $\beta$ -BCNF for a set  $\Sigma$ . Indeed, it suffices to check some pFDs in  $\Sigma$  instead of checking all pFDs in  $\Sigma_{\beta}^+$ .

**THEOREM 7.** *A p-schema  $(R, S)$  is in  $\beta$ -BCNF for a set  $\Sigma$  of pFDs over  $(R, S)$  iff for every pFD  $(X \rightarrow Y, \beta') \in \Sigma$  where  $\beta' \geq \beta$  and  $Y \not\subseteq X$  we have  $(X \rightarrow R, \beta) \in \Sigma_{\beta}^+$ .*

**EXAMPLE 5.** Let  $(\text{MEETING}, S)$  and  $\Sigma$  be as in Example 1. Using Theorem 7 we can observe that the schema is neither in  $\beta_1$ -, nor  $\beta_2$ -, nor  $\beta_4$ -BCNF for  $\Sigma$ , but it is in  $\beta_3$ -BCNF for  $\Sigma$ . By Theorem 6 we conclude that the schema is neither in  $\alpha_4$ -, nor  $\alpha_3$ -, nor  $\alpha_1$ -RFNF for  $\Sigma$ , but it is in  $\alpha_2$ -RFNF for  $\Sigma$ . By Theorem 5 it follows that MEETING is neither in BCNF for  $\Sigma_{\beta_1}$ , nor  $\Sigma_{\beta_2}$ , nor  $\Sigma_{\beta_4}$ , but it is in BCNF for  $\Sigma_{\beta_3}$ . Finally, by Theorem 4, it follows that MEETING is neither in RFNF for  $\Sigma_{\alpha_4}$ , nor  $\Sigma_{\alpha_3}$ , nor  $\Sigma_{\alpha_1}$ , but it is in RFNF for  $\Sigma_{\alpha_2}$ .

## 7.3 Scaled Third Normal Form

We now introduce  $\beta$ -Third normal form (3NF), whose goal is to ensure that all FDs can be enforced locally, without the need of joining relations to check for consistency of updates. Recall the classical 3NF condition: A relation schema  $R$  is

in 3NF for a given FD set  $\Sigma$  iff for every FD  $X \rightarrow A \in \Sigma_{\alpha}^+$  where  $A \notin X$ , we have  $X \rightarrow R \in \Sigma_{\alpha}^+$  or  $A$  is a prime attribute. An attribute  $A$  is *prime* iff it occurs in some minimal key for  $\Sigma$ . An attribute subset  $X$  of  $R$  is a *key* of  $R$  for  $\Sigma$  iff  $X \rightarrow R \in \Sigma_{\alpha}^+$ . A key  $X$  of  $R$  is *minimal* for  $\Sigma$  iff every proper subset  $Y$  of  $X$  is not a key of  $R$  for  $\Sigma$ .

We now introduce analogous concepts for the possibilistic setting. Given a p-schema  $(R, S)$ , a c-degree  $\beta \in S^T$ , and a set  $\Sigma$  of pFDs over  $(R, S)$ , a subset  $X$  of  $R$  is a  $\beta$ -key of  $R$  for  $\Sigma$  iff  $(X \rightarrow R, \beta) \in \Sigma_{\beta}^+$ . A  $\beta$ -key  $X$  of  $R$  for  $\Sigma$  is a  $\beta$ -minimal key iff every proper subset  $Y$  of  $X$  is not a  $\beta$ -key of  $R$  for  $\Sigma$ . An attribute  $A \in R$  is  $\beta$ -prime iff it is contained in some  $\beta$ -minimal key  $X$  of  $R$  for  $\Sigma$ .

**DEFINITION 5.** A p-schema  $(R, S)$  is in  $\beta$ -Third Normal Form (3NF) for a set  $\Sigma$  of pFDs over  $(R, S)$  iff for every pFD  $(X \rightarrow A, \beta) \in \Sigma_{\beta}^+$  where  $A \notin X$ , we have  $(X \rightarrow R, \beta) \in \Sigma_{\beta}^+$  or  $A$  is a  $\beta$ -prime attribute.  $\square$

Theorem 5 characterized  $\beta$ -BCNF for the pFD set  $\Sigma$  in terms of classical BCNF for the  $\beta$ -cut  $\Sigma_{\beta}$ . An analogous result holds for 3NF.

**THEOREM 8.**  *$(R, S)$  is in  $\beta$ -3NF for a set  $\Sigma$  if and only if  $R$  is in 3NF for  $\Sigma_{\beta}$ .*

Deciding if a schema is in  $\beta$ -3NF can be done by checking some pFDs in  $\Sigma$  instead of checking all pFDs in  $\Sigma_{\beta}^+$ .

**THEOREM 9.** *A p-schema  $(R, S)$  is in  $\beta$ -3NF for a set  $\Sigma$  of pFDs over  $(R, S)$  iff for every pFD  $(X \rightarrow Y, \beta') \in \Sigma$  where  $\beta' \geq \beta$  and  $Y \not\subseteq X$ , we have  $(X \rightarrow R, \beta) \in \Sigma_{\beta}^+$  or every attribute  $A \in Y - X$  is a  $\beta$ -prime.*

**EXAMPLE 6.** Let  $(\text{MEETING}, S)$  and  $\Sigma$  be as in Example 1. By Theorem 9, the schema is in neither  $\beta_1$ - nor  $\beta_2$ -3NF for  $\Sigma$ , but it is in  $\beta_3$ -3NF and in  $\beta_4$ -3NF for  $\Sigma$ . Finally, by Theorem 8, it follows that MEETING is neither in 3NF for  $\Sigma_{\beta_1}$  nor  $\Sigma_{\beta_2}$ , but it is in 3NF for  $\Sigma_{\beta_3}$  and in 3NF for  $\Sigma_{\beta_4}$ .

## 8. SCALED NORMALIZATION

We establish algorithms to design relational schemata for applications with uncertain data. For that purpose, we normalize a given p-schemata  $(R, S)$  for the given pFD set  $\Sigma$ . Our strategy is to fix a c-degree  $\beta \in S^T$  that determines which possible worlds we normalize for which FDs. We pursue BCNF decompositions to obtain lossless decompositions free from any data redundancy but potentially not dependency-preserving (that is, some FDs may require validation on the join of some relations). We also pursue 3NF-synthesis to obtain lossless, dependency-preserving decompositions where the level of data redundancy is minimal for all decompositions that are dependency-preserving. Applying our strategy to different c-degrees provides organizations with a whole scale of normalized schemata, each targeted at different levels of data integrity, data losslessness, and the efficiency of different updates and queries.

### 8.1 Scaled BCNF Decomposition

We recall classical terminology. A decomposition of relation schema  $R$  is a set  $\mathcal{D} = \{R_1, \dots, R_n\}$  of relation schemata such that  $R_1 \cup \dots \cup R_n = R$ . For  $R_j \subseteq R$  and FD set  $\Sigma$  over  $R$ ,  $\Sigma[R_j] = \{X \rightarrow Y \mid X \rightarrow Y \in \Sigma_{\alpha}^+ \text{ and } X, Y \subseteq R_j\}$  is the *projection* of  $\Sigma$  onto  $R_j$ . A decomposition  $\mathcal{D}$

of a relation schema  $R$  with FD set  $\Sigma$  is *lossless* iff every relation  $r$  over  $R$  that satisfies  $\Sigma$  is the join of its projections on the elements of  $\mathcal{D}$ , that is,  $r = \bowtie_{R_j \in \mathcal{D}} r[R_j]$ . Here,  $r[R_j] = \{t(R_j) \mid t \in r\}$ . A BCNF decomposition of a relation schema  $R$  with FD set  $\Sigma$  is a decomposition  $\mathcal{D}$  of  $R$  where every  $R_j \in \mathcal{D}$  is in BCNF for  $\Sigma[R_j]$ . Theorem 5 motivates Definition 6.

**DEFINITION 6.** An  $\alpha_{k+1-i}$ -lossless BCNF decomposition of a  $p$ -schema  $(R, \{\alpha_1, \dots, \alpha_{k+1}\})$  for the pFD set  $\Sigma$  is a lossless BCNF decomposition of  $R$  for  $\Sigma_{\beta_i}$ .  $\square$

Instrumental to Definition 6 is the following decomposition theorem. It covers the classical decomposition theorem [23] as the special case of having just one possible world.

**THEOREM 10.** Let  $(X \rightarrow Y, \beta_i)$  be a pFD that satisfies the  $p$ -relation  $(r, \text{Poss}_r)$  over the  $p$ -schema  $(R, \mathcal{S})$ . Then  $r_{k+1-i} = r_{k+1-i}[XY] \bowtie r_{k+1-i}[X(R - Y)]$ , that is, the possible world  $r_{k+1-i}$  of  $r$  is the lossless join of its projections on  $XY$  and  $X(R - Y)$ .

Therefore, an  $\alpha_{k+1-i}$ -lossless BCNF decomposition for a pFD set  $\Sigma$  can be obtained by performing a classical lossless BCNF decomposition for the  $\beta_i$ -cut  $\Sigma_{\beta_i}$  of  $\Sigma$ . This suggests a simple lossless BCNF decomposition strategy.

<b>PROBLEM:</b> Scaled BCNF Decomposition	
<b>INPUT:</b>	P-Relation Schema $(R, \mathcal{S})$ Set $\Sigma$ of pFDs over $(R, \mathcal{S})$ Certainty degree $\beta_i \in \mathcal{S}^T$
<b>OUTPUT:</b>	$\alpha_{k+1-i}$ -lossless BCNF decomposition of $(R, \mathcal{S})$ for $\Sigma$
<b>METHOD:</b>	Perform a lossless BCNF decomposition of $R$ for $\Sigma_{\beta_i}$

We illustrate the decomposition on our running example.

**EXAMPLE 7.** Let  $(\text{MEETING}, \mathcal{S})$  and  $\Sigma$  be as in Example 1. As  $(\text{MEETING}, \mathcal{S})$  is not in  $\beta_2$ -BCNF for  $\Sigma$ , we perform an  $\alpha_3$ -lossless BCNF decomposition for  $\Sigma_{\beta_2}$ . The result is the decomposition  $D_3$  from Figure 3:  $R_1 = \{\text{Project}, \text{Room}, \text{Time}\}$  with projected FD set  $\Sigma_{\beta_2}[R_1] = \{\text{Room}, \text{Time} \rightarrow \text{Project}\}$ , and  $R_2 = \{\text{Manager}, \text{Room}, \text{Time}\}$  with projected FD set  $\Sigma_{\beta_2}[R_2] = \{\text{Manager}, \text{Time} \rightarrow \text{Room}\}$ . Every FD in  $\Sigma_{\beta_2}$  is implied by  $\Sigma_{\beta_2}[R_1] \cup \Sigma_{\beta_2}[R_2]$ .  $\square$

The last example is rather special, since one cannot hope to preserve all FDs in the BCNF decomposition process. This is illustrated with another example.

**EXAMPLE 8.** Let  $(\text{MEETING}, \mathcal{S})$  and  $\Sigma$  be as in Example 1. As  $(\text{MEETING}, \mathcal{S})$  is not in  $\beta_4$ -BCNF for  $\Sigma$ , an  $\alpha_1$ -lossless BCNF decomposition for  $\Sigma_{\beta_4}$  results in  $D_1$  from Figure 3:  $R_1 = \{\text{Manager}, \text{Project}\}$  with projected FD set  $\Sigma_{\beta_4}[R_1] = \{\text{Project} \rightarrow \text{Manager}\}$ , and  $R_2 = \{\text{Project}, \text{Room}, \text{Time}\}$  with projected FD set  $\Sigma_{\beta_4}[R_2] = \{\text{Room}, \text{Time} \rightarrow \text{Project}; \text{Project}, \text{Time} \rightarrow \text{Room}\}$ . The FD  $\text{Manager}, \text{Time} \rightarrow \text{Room}$  is not implied by  $\Sigma_{\beta_4}[R_1] \cup \Sigma_{\beta_4}[R_2]$ .  $\square$

A decomposition  $\mathcal{D}$  of relation schema  $R$  with FD set  $\Sigma$  is *dependency-preserving* if and only if  $\Sigma_{\mathcal{A}}^+ = (\bigcup_{R_j \in \mathcal{D}} \Sigma[R_j])_{\mathcal{A}}^+$ .

**DEFINITION 7.** A  $\beta$ -dependency-preserving decomposition of a  $p$ -schema  $(R, \mathcal{S})$  for the pFD set  $\Sigma$  is a dependency-preserving decomposition of  $R$  for  $\Sigma_{\beta}$ .  $\square$

The  $\alpha_3$ -lossless BCNF decomposition from Example 7 is  $\beta_2$ -dependency-preserving, but the  $\alpha_1$ -lossless BCNF decomposition from Example 8 is not  $\beta_4$ -dependency-preserving. In fact, for  $(\text{MEETING}, \mathcal{S})$  and  $\Sigma$  as in Example 1 there is no  $\beta_4$ -dependency-preserving,  $\alpha_1$ -lossless BCNF decomposition of  $\Sigma$ . In practice, lost dependencies can only be validated by joining relations after inserts or modification. For example, to validate the FD  $\text{Manager}, \text{Time} \rightarrow \text{Room}$  after an update, one would have to join  $R_1$  and  $R_2$  from Example 8. This can be prohibitively expensive.

## 8.2 Scaled 3NF Synthesis

3NF synthesis guarantees dependency-preservation, but cannot guarantee the elimination of all data redundancy in terms of FDs. Recently, Third normal form was shown to pay the smallest possible price, in terms of data redundancy, for achieving dependency-preservation [19]. We will now equip our schema design framework for uncertain data with an appropriate 3NF synthesis strategy.

A 3NF decomposition of a relation schema  $R$  for an FD set  $\Sigma$  is a decomposition  $\mathcal{D}$  of  $R$  where every  $R_j \in \mathcal{D}$  is in 3NF for  $\Sigma[R_j]$ . Theorem 8 motivates the following definition.

**DEFINITION 8.** A  $\beta_i$ -dependency-preserving,  $\alpha_{k+1-i}$ -lossless 3NF decomposition of a  $p$ -schema  $(R, \{\alpha_1, \dots, \alpha_{k+1}\})$  for the pFD set  $\Sigma$  is a dependency-preserving, lossless 3NF decomposition of  $R$  for  $\Sigma_{\beta_i}$ .  $\square$

Due to Theorem 10 a  $\beta_i$ -dependency-preserving,  $\alpha_{k+1-i}$ -lossless 3NF synthesis for a pFD set  $\Sigma$  can simply be obtained by performing a classical dependency-preserving lossless 3NF synthesis for the  $\beta_i$ -cut  $\Sigma_{\beta_i}$  of  $\Sigma$ .

<b>PROBLEM:</b> Scaled 3NF Synthesis	
<b>INPUT:</b>	P-relation schema $(R, \mathcal{S})$ Set $\Sigma$ of pFDs over $(R, \mathcal{S})$ Certainty degree $\beta_i \in \mathcal{S}^T$
<b>OUTPUT:</b>	$\beta_i$ -dependency-preserving, $\alpha_{k+1-i}$ -lossless 3NF decomposition of $(R, \mathcal{S})$ for $\Sigma$
<b>METHOD:</b>	Perform a dependency-preserving, lossless 3NF synthesis of $R$ for $\Sigma_{\beta_i}$

We illustrate the synthesis on our running example.

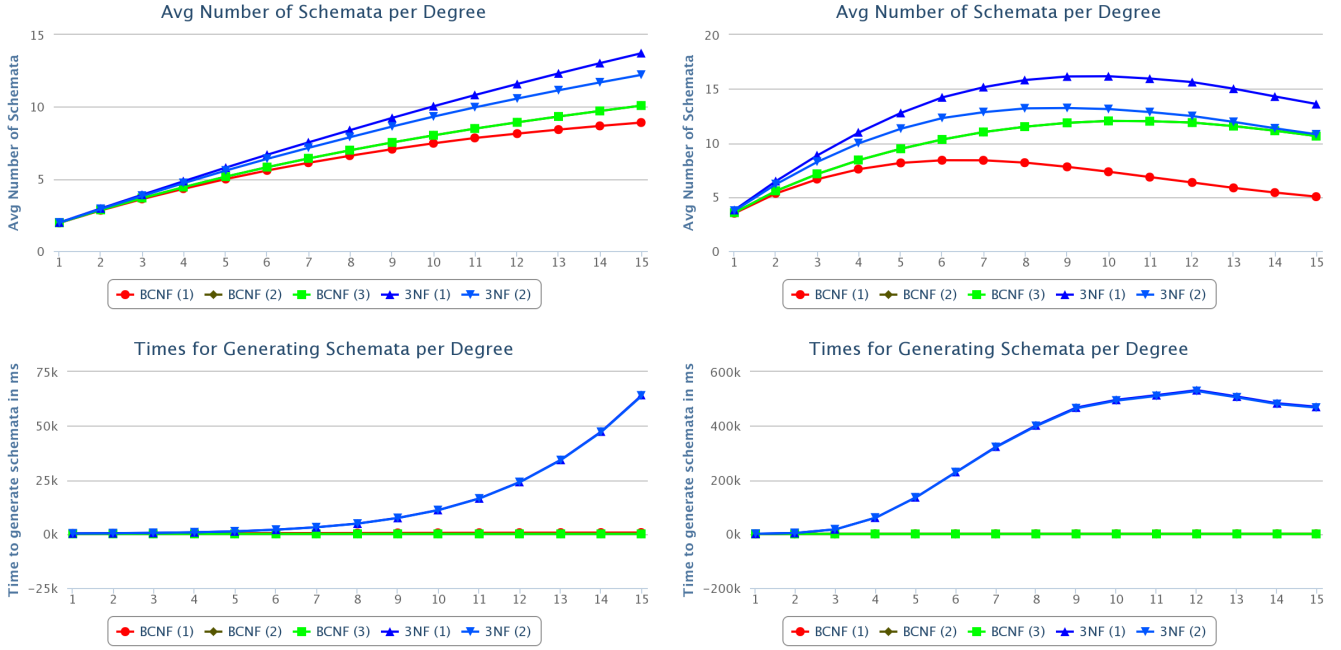
**EXAMPLE 9.** Let  $(\text{MEETING}, \mathcal{S})$  and  $\Sigma$  be as in Example 1. As  $(\text{MEETING}, \mathcal{S})$  is not in  $\beta_2$ -3NF for  $\Sigma$ , we perform an  $\alpha_3$ -lossless,  $\beta_2$ -dependency-preserving 3NF synthesis. The result consists of  $R_1 = \{\text{Project}, \text{Room}, \text{Time}\}$  with projected FD set  $\Sigma_{\beta_2}[R_1] = \{\text{Room}, \text{Time} \rightarrow \text{Project}\}$ , and  $R_2 = \{\text{Manager}, \text{Room}, \text{Time}\}$  with projected FD set  $\Sigma_{\beta_2}[R_2] = \{\text{Manager}, \text{Time} \rightarrow \text{Room}\}$ . Note that  $R_1$  is in BCNF for  $\Sigma_{\beta_2}[R_1]$ , and  $R_2$  is in BCNF for  $\Sigma_{\beta_2}[R_2]$ , as observed in Example 7.  $\square$

## 9. EXPERIMENTS

We analyze the results of our experiments for the scaled versions of BCNF decompositions and 3NF syntheses. This also includes new insights into the traditional problems.

Our core analysis concerns the trade-offs between the elimination of data redundancy achieved by BCNF decompositions and the preservation of FDs achieved by 3NF syntheses, as well as between query and update efficiency. The simplest general measure for the query and update efficiency





**Figure 4: Average size of decompositions and average times to generate them, taken over 5,000 runs each, where the input relation schemata have  $n = 15$  attributes, and the maximum pFD set size is  $n$  in the left column and  $5n$  in the right column**

of a decomposition is the *number of its schemata* - to which we refer as the size of the decomposition. Intuitively, the more schemata a decomposition contains the more updates and the less queries it can support efficiently. In particular, less schemata require less joins, the main source of complexity in query evaluation. The normalization effort can be measured either in view of the given  $c$ -degree  $\beta$  (possibilistic view), or in the size of its  $\beta$ -cut (classical view). Our analysis should be read with both views in mind.

For our analysis we compare five normalization algorithms: BCNF(1) is the classical BCNF decomposition algorithm that takes exponential time and space, BCNF(2) is Tsou and Fischer’s P-time algorithm, BCNF(3) is the variant of BCNF(2) where the projection of original FD sets is only determined for the output schemata, 3NF(1) is the 3NF synthesis algorithm that first computes a non-redundant, L-reduced cover in which all FDs have a singleton attribute on the right-hand side, and 3NF(2) is the 3NF synthesis algorithm that first computes a canonical cover.

The top row of Figure 4 shows the graphs of two experiments. For each of 15  $c$ -degrees  $\beta_1, \dots, \beta_{15}$ , 5,000 sets of at most  $a \cdot n$  pFDs over a relation schema with  $n = 15$  attributes were created, and  $\beta_i$ -BCNF decompositions and  $\beta_i$ -3NF syntheses created, using the 5 algorithms. The graph with  $a = 1$  is shown on the left and  $a = 5$  on the right.

The experiments show nearly uniform results: The smallest-sized decomposition is BCNF(1), followed by BCNF(2,3) which agree, followed by 3NF(2), and then 3NF(1). The results are intuitive as 3NF produces more schemata to accommodate the preservation of all FDs and thereby sacrificing the elimination of some data redundancy. The graphs show the price for producing a BCNF decomposition in P-time, which results from not having to check

the BCNF condition of intermediate results and thereby producing superfluous schemata not required by BCNF(1). It is interesting that this price is still mostly smaller than that to accommodate dependency-preservation. However, we will see later that the superfluous schemata in BCNF(2,3) are unlikely to help with dependency-preservation. Finally, 3NF(2) generates less schemata than 3NF(1) as a canonical cover contains less FDs than the other cover considered. The difference between BCNF(1) and 3NF(2) tells us how many more schemata are required to guarantee dependency-preservation. The results suggest to prefer BCNF decomposition in terms of query efficiency and, obviously, in terms of elimination of data redundancy. They do not mean to prefer BCNF decomposition in terms of update efficiency. Indeed, BCNF decompositions are only more efficient for updates of redundant data value occurrences caused by FDs preserved during the decomposition, while update inefficiencies result from having to join schemata to validate FDs that were not preserved during the decomposition.

The graphs further show that the bigger the size of the input relation schema, the bigger the size of the decomposition - observing the graphs on the left and that on the right. However, the story is different in terms of growing FD size. When there are about as many given FDs as underlying attributes, then the average size of decompositions grows linearly. When there are significantly more given FDs than underlying attributes, then a global maximum is observable in the average size of decompositions. This is explained by the fact that there will be a point when there are sufficiently many FDs that turn attribute sets into keys, in which case further decompositions become unnecessary. Not surprisingly, this saturation point comes earlier for BCNF(1) than for any of the other algorithms considered. For larger FD

	$\alpha_{k+1-i}$ -lossless	Free from $\alpha_{k+1-i}$ data redundancy	$\beta_i$ -dependency-preserving
$\beta_i$ -BCNF	✓	✓	
$\beta_i$ -3NF	✓		✓

**Table 3: Achievements of  $\beta$ -BCNF and  $\beta$ -3NF**

sets, higher levels of data integrity achieve the same levels of query efficiency as significantly lower levels of data integrity.

The bottom of Figure 4 shows the average times to generate the decompositions for our two experiments. As observed for the average size of decompositions, both experiments show nearly uniform results for the average times: BCNF(2) takes the lowest time, followed by BCNF(3), followed by BCNF(1), and followed by 3NF(1,2). These results are intuitive, too. BCNF(2) demonstrates its major strength in comparison to the other algorithms, with a maximum average time of 8ms for  $a = 1$  and 13ms for  $a = 5$ . BCNF(3) requires slightly more time with a maximum average time of 12ms for  $a = 1$  and 35ms for  $a = 5$ , resulting from the projections of the original FD sets to the final output schemata at the end. BCNF(1) requires more time with a maximum average time of 630ms for  $a = 1$  and 627ms for  $a = 5$ , as a result of projecting the original FD sets to all intermediate and final schemata. Finally, 3NF(1,2) require significantly more time with a maximum average time of 64s for  $a = 1$  and 530s for  $a = 5$ .

Two main observations illustrate inherent trade-offs. For BCNF decompositions firstly, smaller-sized decompositions come at the price of requiring more time to be generated. In other words, polynomial-time BCNF decompositions result from not having to spend time on validating the BCNF condition, thereby generating additional schemata that are superfluous in terms of update efficiency but may significantly add to the inefficiency of query evaluation. For 3NF decompositions secondly, the price to pay for dependency-preservation is a higher complexity in terms of the size of the decompositions and the time taken to generate them.

## 10. CONCLUSION AND FUTURE WORK

We developed a full framework to design relational schemata for applications with uncertain data. The fact that FDs are closed downwards, i.e. hold on every sub-relation of a relation that satisfies them, enabled us to introduce different levels of data redundancy. We established a customized segmentation of the well-known Boyce-Codd and Third normal forms, extending all of their properties to each segment, Table 3. Certainty degrees of FDs provide a mechanism for organizations to control the levels of data integrity, data losslessness, query efficiency, and update efficiency, Figure 1.

Several problems warrant future research. As relational schema design for certain data depends on the identification of FDs that are meaningful for the given application, schema design for uncertain data depends on the identification of meaningful pFDs. This motivates research about Armstrong relations [11] in the possibilistic setting. Further normal forms should be customized to their use for uncertain data, including 4NF [10, 27] as well as Inclusion Dependency normal form [20]. Recently, the relational normalization toolbox has been extended to incomplete relations [15, 16, 17, 18]. It would be interesting to develop a normalization framework for possibilistic incomplete relations.

## 11. REFERENCES

- [1] W. W. Armstrong. Dependency structures of data base relationships. In *IFIP Congress*, pages 580–583, 1974.
- [2] C. Beeri and P. Bernstein. Computational problems related to the design of normal form relational schemas. *ACM TODS*, 4(1):30–59, 1979.
- [3] O. Benjelloun, A. D. Sarma, A. Y. Halevy, M. Theobald, and J. Widom. Databases with uncertainty and lineage. *VLDB J.*, 17(2):243–264, 2008.
- [4] P. A. Bernstein. Synthesizing third normal form relations from functional dependencies. *ACM TODS*, 1(4):277–298, 1976.
- [5] N. A. Chaudhry, J. R. Moyne, and E. A. Rundensteiner. An extended database design methodology for uncertain data management. *Inf. Sci.*, 121(1-2):83–112, 1999.
- [6] E. F. Codd. A relational model of data for large shared data banks. *Commun. ACM*, 13(6):377–387, 1970.
- [7] E. F. Codd. Further normalization of the database relational model. In *Courant Computer Science Symposia 6: Data Base Systems*, pages 33–64, 1972.
- [8] N. N. Dalvi and D. Suciu. Management of probabilistic data: foundations and challenges. In *PODS*, pages 1–12, 2007.
- [9] D. Dubois and H. Prade. Possibility theory. In R. A. Meyers, editor, *Computational Complexity*, pages 2240–2252. Springer New York, 2012.
- [10] R. Fagin. Multivalued dependencies and a new normal form for relational databases. *ACM TODS*, 2(3):262–278, 1977.
- [11] R. Fagin. Horn clauses and database dependencies. *J. ACM*, 29(4):952–985, 1982.
- [12] N. Hall, H. Köhler, S. Link, H. Prade, and X. Zhou. Cardinality constraints on qualitatively uncertain data. *Data Knowl. Eng.*, 99:126–150, 2015.
- [13] A. K. Jha and D. Suciu. Probabilistic databases with markoViews. *PVLDB*, 5(11):1160–1171, 2012.
- [14] H. Köhler, U. Leck, S. Link, and H. Prade. Logical foundations of possibilistic keys. In *JELIA*, pages 181–195, 2014.
- [15] H. Köhler, U. Leck, S. Link, and X. Zhou. Possible and certain keys for SQL. *VLDB J.*, 25(4):571–596, 2016.
- [16] H. Köhler and S. Link. SQL schema design: Foundations, normal forms, and normalization. In *SIGMOD*, pages 267–279, 2016.
- [17] H. Köhler, S. Link, and X. Zhou. Possible and certain SQL keys. *PVLDB*, 8(11):1118–1129, 2015.
- [18] H. Köhler, S. Link, and X. Zhou. Discovering meaningful certain keys from incomplete and inconsistent relations. *IEEE Data Eng. Bull.*, 39(2):21–37, 2016.
- [19] S. Kolahi and L. Libkin. An information-theoretic analysis of worst-case redundancy in database design. *ACM TODS*, 35(1), 2010.
- [20] M. Levene and M. W. Vincent. Justification for inclusion dependency normal form. *IEEE TKDE*, 12(2):281–291, 2000.
- [21] S. Link and H. Prade. Relational database schema design for uncertain data. Technical Report CDMTCS-469, The University of Auckland, 2014.
- [22] S. Link and H. Prade. Possibilistic functional dependencies and their relationship to possibility theory. *IEEE Trans. Fuzzy Systems*, 24(3):757–763, 2016.
- [23] J. Rissanen. Independent components of relations. *ACM TODS*, 2(4):317–325, 1977.
- [24] A. D. Sarma, J. D. Ullman, and J. Widom. Schema design for uncertain databases. In *AMW*, 2009.
- [25] D. Suciu, D. Olteanu, C. Ré, and C. Koch. *Probabilistic Databases*. Morgan & Claypool Publishers, 2011.
- [26] D.-M. Tsou and P. C. Fischer. Decomposition of a relation scheme into Boyce-Codd normal form. *SIGACT News*, 14(3):23–29, 1982.
- [27] M. Vincent. Semantic foundations of 4NF in relational database design. *Acta Inf.*, 36(3):173–213, 1999.