

# A Uniform Dependency Language for Improving Data Quality

Wenfei Fan      Floris Geerts

*University of Edinburgh*

{wenfei, fgeerts}@inf.ed.ac.uk

## Abstract

*A variety of dependency formalisms have been studied for improving data quality. To treat these dependencies in a uniform framework, we propose a simple language, Quality Improving Dependencies (QIDs). We show that previous dependencies considered for data quality can be naturally expressed as QIDs, and that different enforcement mechanisms of QIDs yield various data repairing strategies.*

## 1 Introduction

Data quality has been a longstanding line of research for decades, and has become one of the most pressing challenges for data management. As an example, it is estimated that dirty data costs US companies alone 600 billion dollars each year [11]. With this comes the need for studying techniques for improving data quality.

A variety of dependency formalisms have recently been studied to specify data quality rules, *e.g.*, functional dependencies (FDs) [1], conditional functional dependencies (CFDs) [6, 14], order dependencies (ODs) [22, 29, 23], currency dependencies (CDs) [17], sequential dependencies (SDs) [25], matching dependencies (MDs) [18, 13, 4] and editing rules (eRs) [19]. These dependencies have proved useful in detecting and correcting inconsistencies in relational data. However, given these different formalisms, one would naturally ask which dependencies should be used in an application? Is a data quality rule expressible in one formalism but not in another? Is it more expensive to use one formalism than another? How should dependencies be enforced to repair data? These suggest that we study these formalisms in a comparative basis.

This paper takes a first step toward answering these questions. We propose a simple dependency language, referred to as Quality Improving Dependencies (QIDs), and show that all the dependencies mentioned above can be expressed as QIDs. That is, QIDs provide a uniform language to characterize these formalisms. We also introduce different mechanisms for enforcing QIDs as data quality rules, via revisions of the chase process [1]. We show that data repairing algorithms [5, 9, 19, 20] are just implementations of these mechanisms.

In the rest of the paper, we introduce QIDs (Section 2), and show that dependencies studied for data quality are special cases of QIDs (Section 3). In addition, we present data quality problems in terms of QIDs (Section 4), emphasizing data repairing as QID enforcement. We also give an overview of recent work on these issues.

## 2 Quality Improving Dependencies

We define QIDs in terms of a notion of comparison relations, which compare pairs of tuples based on some of their characteristics. Below we first introduce comparison relations, and then define QIDs.

---

*Copyright 2011 IEEE. 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.*

**Bulletin of the IEEE Computer Society Technical Committee on Data Engineering**

---

## 2.1 Comparison Relations

Consider a database schema  $\mathcal{R} = (R_1, \dots, R_m)$ , in which each relation schema  $R_i$  is specified with a set of attributes, denoted by  $R_i(\text{tid}, A_1, \dots, A_{n(i)})$ . Here  $\text{tid}$  is a *tuple identifier* and the domain of attribute  $A_j$  is denoted by  $\text{dom}(A_j)$ . We say that two attributes  $A_j$  and  $A_k$  are *compatible* if  $\text{dom}(A_j) = \text{dom}(A_k)$ . An instance  $\mathcal{I}$  of  $\mathcal{R}$  is  $(I_1, \dots, I_m)$ , where for  $i \in [1, m]$ ,  $I_i$  is finite set of tuples  $t$  with  $t[A_j] \in \text{dom}(A_j)$  for each  $A_j \in R_i$ .

A comparison relation is simply a binary relation  $\sim \subseteq \text{dom}(\text{tid}) \times \text{dom}(\text{tid})$ . In practice, it is often used to group tuples (identified by their  $\text{tid}$ 's) by comparing certain attributes (properties) of the tuples. More specifically, let  $X$  be a set of attributes  $\{A_1, \dots, A_k\}$ , and  $\mathcal{B}$  be a “binary” relation  $\mathcal{B} \subseteq \text{dom}(X) \times \text{dom}(X)$  that compares the  $X$  attributes of tuples, where  $\text{dom}(X) = \text{dom}(A_1) \times \dots \times \text{dom}(A_k)$ . A *comparison relation via  $(X, \mathcal{B})$* , denoted by  $\sim_{\mathcal{B}}^X$ , is defined such that for any tuples  $t_1$  and  $t_2$  that contain attributes  $X$ ,  $\sim_{\mathcal{B}}^X(t_1[\text{tid}], t_2[\text{tid}])$  iff  $\mathcal{B}(t_1[X], t_2[X])$ . Let  $\mathcal{I}$  be an instance of  $\mathcal{R}$  and  $I$  be the set of all tuples in  $\mathcal{I}$  that contain attributes  $X$ . Then

$$\sim_{\mathcal{B}}^X := \{(\text{tid}_1, \text{tid}_2) \mid t_i \in I, t_i[\text{tid}] = \text{tid}_i, i = 1, 2, \mathcal{B}(t_1[X], t_2[X])\}$$

for instance  $\mathcal{I}$ . That is,  $\sim_{\mathcal{B}}^X$  relates pairs of tuples by only considering the information present in their  $X$  attributes. Observe that a pair of tuples related by  $\sim_{\mathcal{B}}^X$  may possibly come from *different* relations in  $\mathcal{I}$ .

As an example, we consider comparison relations  $\sim_{\mathcal{B}}^A$  defined on an attribute  $A$ , where  $\mathcal{B}$  is a variation of the equality relation on  $A$ , *i.e.*, a binary relation in  $\text{dom}(A) \times \text{dom}(A)$  that is reflexive, symmetric and transitive.

**Example 1:** One can define  $\sim_{\mathcal{B}}^A$  to group tuples based on the equality of a certain attribute, via the following  $\mathcal{B}$ :

(1)  $\text{eq}_A = \{(a, a) \mid a \in \text{dom}(A)\}$ , *i.e.*, for any pair  $(t_1, t_2)$  of tuples,  $\sim_{\text{eq}_A}^A(t_1[\text{tid}], t_2[\text{tid}])$  as long as  $t_1$  and  $t_2$  have the same  $A$  attribute value. For a specific constant  $c \in \text{dom}(A)$ , one can define  $\mathcal{B}$  to be  $\text{eq}_c = \{(c, c)\}$ , such that  $\sim_{\text{eq}_c}^A$  collects all tuples  $t$  in which attribute  $t[A]$  is restricted to be the predefined constant  $c$ .

(2)  $\text{eq}_S = \cup_{c \in S} \text{eq}_c$  for a (finite) set  $S \subseteq \text{dom}(A)$ , *i.e.*,  $\sim_{\text{eq}_S}^A$  groups those tuples  $t$  in which  $t[A] \in S$ . Similarly one may define  $\mathcal{B}$  as  $\text{eq}_{\bar{S}} = (\text{eq}_A \setminus \text{eq}_S)$ , such that  $\sim_{\text{eq}_{\bar{S}}}^A$  groups tuples  $t$  in which  $t[A]$  does not belong to  $S$ .

(3) When the domain of  $A$  is ordered, we may define  $\mathcal{B}$  as  $\text{eq}_{S(\text{op } b)}$  where  $S(\text{op } b) = \{a \in \text{dom}(A) \mid a \text{ op } b\}$ , and  $\text{op} \in \{<, \leq, >, \geq, =, \neq\}$ ; *i.e.*,  $\sim_{\text{eq}_{S(\text{op } b)}}^A$  groups tuples  $t$  such that  $t[A] \text{ op } b$ .  $\square$

In practice, it is common to compare values of an attribute  $A$  based on their similarity rather than equality. In general, a similarity relation  $\mathcal{B} \subseteq \text{dom}(A) \times \text{dom}(A)$  if it is reflexive and symmetric, but not necessarily transitive.

**Example 2:** Similarity relations are often derived from metrics. Suppose that  $A$  is an attribute whose domain is string. We define the similarity relation  $\approx_{\theta}^{\text{dl}}$  based on the Damerau-Levenstein metric (with threshold  $\theta$ ):

$$\approx_{\theta}^{\text{dl}} = \{(s, t) \in \text{string} \times \text{string} \mid \text{edit distance}(s, t) \leq \theta\}.$$

In fact, any metric such as q-grams or Jaro distance can be turned into a similarity relation (see [10] for a survey on distance functions), from which a comparison relation (*e.g.*,  $\sim_{\approx_{\theta}^{\text{dl}}}^A$ ) can be readily derived.  $\square$

The comparison relations defined in terms of equality and similarity relations can be readily extended to a set  $X$  of attributes, *e.g.*,  $\sim_{\text{eq}_X}^X$  is defined as  $\bigcap_{A \in X} \sim_{\text{eq}_A}$ , and  $\sim_{\approx_X}^X$  as  $\bigcap_{A \in X} \sim_{\approx_A}$ .

A comparison relation  $\sim_{\mathcal{B}}^X$  via  $(X, \mathcal{B})$  is said to be *decomposable* if there exist  $\mathcal{B}_i \subseteq \text{dom}(A_i) \times \text{dom}(A_i)$  such that  $\sim_{\mathcal{B}}^X = \bigcap_{A_i \in X} \sim_{\mathcal{B}_i}^{A_i}$ . For instance, the comparison relations  $\sim_{\text{eq}_X}^X$  and  $\sim_{\approx_X}^X$  given above are clearly decomposable. Nevertheless, not all comparison relations are decomposable, as shown below.

**Example 3:** Let  $X = \{A_1, \dots, A_k\}$  be a set of ordered attributes, and let  $\bar{a} = (a_1, \dots, a_k)$  and  $\bar{b} = (b_1, \dots, b_k)$  be tuples in  $\text{dom}(A_1) \times \dots \times \text{dom}(A_k)$ . Define  $\bar{a}$  *lexico*  $\bar{b}$  iff either  $a_i = b_i$  for all  $i \in [1, k]$  or there exist a  $j \in [1, k]$  such that  $a_j <_{A_j} b_j$  while  $a_i = b_i$  for  $i \in [1, j-1]$ . Here  $<_{A_j}$  denotes an order on  $A_j$ . Then the corresponding comparison relation  $\sim_{\text{lexico}}^X$  is not decomposable. Indeed, observe that  $\sim_{\text{lexico}}^X$  can be written as

$$\left( \bigcap_{i \in [1, k]} \sim_{\text{eq}}^{A_i} \right) \cup \bigcup_{j \in [1, k]} \left( \sim_{<_{A_j}}^{A_j} \cap \bigcap_{i \in [1, j-1]} \sim_{\text{eq}}^{A_i} \right);$$

in other words, it can be written as a union of decomposable comparison relations, but not as  $\bigcap_{i \in [1, k]} \sim_{\mathcal{B}_i}^{A_i}$ . This example also shows that comparison relations can be defined in terms of Boolean operations.  $\square$

## 2.2 Quality Improving Dependencies

We next define quality improving dependencies in terms of comparison relations. Let  $R$  and  $R'$  be two relations in  $\mathcal{R}$ . Let  $X = \{A_1, \dots, A_k\}$  and  $X' = \{A'_1, \dots, A'_k\}$  be attributes in  $R$  and  $R'$ , respectively, such that  $A_i$  and  $A'_i$  are compatible for each  $i \in [1, k]$ ; similarly for attributes  $Y = \{E_1, \dots, E_\ell\}$  and  $Y' = \{E'_1, \dots, E'_\ell\}$ . Let  $\sim_{\mathcal{B}_X}^X$  and  $\sim_{\mathcal{B}_Y}^Y$  be comparison relations via  $(X, \mathcal{B}_X)$  and  $(Y, \mathcal{B}_Y)$ , respectively.

A *quality improving dependency* (QID) defined on  $(R, R')$  is of the form

$$\varphi = (\sim_{\mathcal{B}_X}^X (R(X), R'(X')) \rightarrow \sim_{\mathcal{B}_Y}^Y (R(Y), R'(Y'))).$$

An instance  $\mathcal{I}$  of  $\mathcal{R}$  satisfies  $\varphi$ , denoted by  $\mathcal{I} \models \varphi$ , if

$$\forall \text{tid} \in I, \text{tid}' \in I' : \sim_{\mathcal{B}_X}^X (\text{tid}, \text{tid}') \longrightarrow \sim_{\mathcal{B}_Y}^Y (\text{tid}, \text{tid}'),$$

where  $I$  is the instance of  $R$  in  $\mathcal{I}$ ,  $I'$  is the instance of  $R'$  in  $\mathcal{I}$  for which attributes of  $X'$  and  $Y'$  are renamed as their counterparts in  $X$  and  $Y$ , respectively, *i.e.*,  $A'_i$  as  $A_i$  for  $i \in [1, k]$  and  $B'_j$  as  $B_j$  for  $j \in [1, \ell]$ .

When  $\sim_{\mathcal{B}_X}^X$  and  $\sim_{\mathcal{B}_Y}^Y$  are decomposable as  $(\mathcal{B}_1^{A_1}, \dots, \mathcal{B}_k^{A_k})$  and  $(\mathcal{B}_1^{E_1}, \dots, \mathcal{B}_\ell^{E_\ell})$ , respectively, then  $\mathcal{I} \models \varphi$  if

$$\forall \text{tid} \in I, \text{tid}' \in I' : \bigwedge_{i \in [1, k]} \sim_{\mathcal{B}_i}^{A_i} (\text{tid}, \text{tid}') \longrightarrow \bigwedge_{j \in [1, \ell]} \sim_{\mathcal{B}_j}^{E_j} (\text{tid}, \text{tid}').$$

If all comparison relations in a QID  $\varphi$  are decomposable, we say that  $\varphi$  is *decomposable*.

When  $R = R'$ ,  $X = X'$  and  $Y = Y'$ , we say that  $\varphi$  is a *mono-QID*; otherwise we call it a *bi-QID*. In the sequel, we simply write a mono-QID as  $R(\sim_{\mathcal{B}_X}^X \rightarrow \sim_{\mathcal{B}_Y}^Y)$ . We only consider birelational QIDs that relate tuples in one relation  $R$  to tuples in another (possibly the same)  $R'$ . In practice, two relations often suffice. Nonetheless, one can readily extend QIDs to multiple relations. When developing algorithms for discovering and validating QIDs, and for detecting and correcting errors based on QIDs, it is important to know whether the QIDs under consideration are decomposable or not, and whether one or two relations are required (see Section 4).

## 3 Capturing Existing Dependency Formalisms with QIDs

As examples of QIDs, we show that a variety of dependencies studied for data quality can be expressed as QIDs, as summarized in Table 3. We first consider formalisms that are expressible as *decomposable mono-QIDs*.

**Functional dependencies (FDs)** [1]. A functional dependency on a relation  $R$  is defined as  $R(X \rightarrow Y)$ , where  $X$  and  $Y$  are sets of attributes. An instance  $I$  of  $R$  satisfies the FD if for any two tuples  $t_1$  and  $t_2$  in  $I$ ,  $t_1[Y] = t_2[Y]$  whenever  $t_1[X] = t_2[X]$ , *i.e.*, the  $X$  attributes of a tuple uniquely determine its  $Y$  attributes.

Every FD  $R(X \rightarrow Y)$  is a decomposable mono-QID  $R(\sim_{\text{eq}_X}^X \rightarrow \sim_{\text{eq}_Y}^Y)$ , where  $\sim_{\text{eq}_X}^X$  is defined in Section 2.

**Conditional functional dependencies (CFDs)** [14]. CFDs extend FDs by incorporating constant patterns. More specifically, a CFD is defined as  $R[(X \rightarrow Y) \mid t_p]$ , where  $R(X \rightarrow Y)$  is an FD, and  $t_p$  is a pattern tuple over  $X \cup Y$ . For each attribute  $A \in X \cup Y$ ,  $t_p[A] \in \text{dom}(A)$  or  $t_p[A]$  is wildcard, denoted by  $\sqcup$ . A tuple  $t$  in an instance  $I$  of  $R$  is said to *match*  $t_p$  in  $X$ , denoted by  $t[X] \asymp t_p[X]$ , if for each  $A \in X$ , either  $t_p[A] = \sqcup$  and  $t[A]$  can be an arbitrary value in  $\text{dom}(A)$ , or  $t[A] = t_p[A]$  when  $t_p[A]$  is a constant. An instance  $I$  of  $R$  satisfies the CFD if for any two tuples  $t_1$  and  $t_2$  in  $I$ , if  $t_1[X] = t_2[X] \asymp t_p[X]$  then also  $t_1[Y] = t_2[Y] \asymp t_p[Y]$ . For instance,  $R[(\text{CC}, \text{zip} \rightarrow \text{street}) \mid (44, \sqcup, \sqcup)]$  is a CFD, which states that when CC (country code) is 44 (indicating UK), zip code uniquely determines street.

Every CFD  $R[(X \rightarrow Y) \mid t_p]$  is a decomposable mono-QID  $R(\sim_{\text{eq}_C}^X \rightarrow \sim_{\text{eq}_C}^Y)$ , where  $\sim_{\text{eq}_C}^X$  is defined as  $\bigcap_{A \in X} \sim_{\text{eq}'_A}^A$ . Here  $\sim_{\text{eq}'_A}^A$  is  $\sim_{\text{eq}_A}^A$  when  $t_p[A] = \sqcup$ , and  $\sim_{\text{eq}_c}^A$  if  $t_p[A] = c$  (see Example 1); similarly for  $\sim_{\text{eq}_C}^Y$ .

**Extended CFDs (eCFDs)** [6]. This class of dependencies extends CFDs  $[R(X \rightarrow Y) \mid t_p]$  such that in the pattern tuples  $t_p$ , for each  $A \in X \cup Y$ ,  $t_p[A] = S$ ,  $t_p[A] = \bar{S}$  for a finite set of constants  $S \subseteq \text{dom}(A)$ , or  $t_p[A] = \sqcup$ . We say that  $t[X] \asymp t_p[X]$  if for each attribute  $A \in X$ ,  $t[A] \in S$  if  $t_p[A] = S$ ,  $t[A] \notin S$  if  $t_p[A] = \bar{S}$ , or  $t_p[A] = \sqcup$ . The semantics of eCFDs is the same as CFDs except for the revised matching operator.

Dependencies [References]	Quality Improving Dependencies		
	Mono/Bi	decomp.	QIDs (basic comparison relations used)
Functional dependencies (FD) [1]	Mono	Yes	$R(\sim_{\text{eq}_X}^X \rightarrow \sim_{\text{eq}_Y}^Y)$ (defined with $\sim_{\text{eq}_A}$ )
Conditional FDs (CFD) [14]	Mono	Yes	$R(\sim_{\text{eq}_C}^X \rightarrow \sim_{\text{eq}_C}^Y)$ (with $\sim_{\text{eq}_A}, \sim_{\text{eq}_a}$ )
Extended CFDs (eCFD) [6]	Mono	Yes	$R(\sim_{\text{eq}_S}^X \rightarrow \sim_{\text{eq}_S}^Y)$ ( $\sim_{\text{eq}_A}, \sim_{\text{eq}_S}, \sim_{\text{eq}_S}$ )
Similarity FD (SFD) [32, 34, 3, 30]	Mono	Yes	$R(\sim_{\approx_X}^X \rightarrow \sim_{\approx_Y}^Y)$ ( $\approx_X$ , reflexive and symmetric)
Domain dependencies (DDs) [12]	Mono	Yes	$R(\emptyset \rightarrow \sim_{\bar{S}}^A)$
Edits [21]	Mono	Yes	$R(\sim_{\bar{S}}^X \rightarrow \sim_{\bar{S}_k}^{A_k})$ (with $\sim_{\bar{S}}^A$ and $\sim_{\bar{S}}^A$ )
Association rules (AR) [35, 28]	Mono	Yes	$R(\sim_{\text{eq}_C}^X \rightarrow \sim_{\text{eq}_a}^A)$ (with $\sim_{\text{eq}_a}$ )
Order dependency (OD) [22, 29, 23]	Mono	Yes	$R(\sim_{\text{op}_X}^X \rightarrow \sim_{\text{op}_Y}^Y)$ (op: partial orders)
Currency dependencies (CD) [17]	Mono	Yes	$R(\sim_{B_X}^X \rightarrow \sim_{A}^A)$ ( $\sim_{A}^A$ a currency order)
Binary FDs (BFD) [33]	Bi	Yes	$\sim_{\text{eq}_X}^X (R(X), R'(X')) \rightarrow \sim_{\text{eq}_Y}^Y (R(Y), R'(Y'))$
Matching dependencies (MD) [18, 4]	Bi	Yes	$\sim_{\approx_X}^X (R(X), R'(X')) \rightarrow \sim_{\text{eq}_Y}^Y (R(Y), R'(Y'))$
Editing rules (eR) [19]	Bi	Yes	$\sim_{B_X}^X (R(X, X_p), R'(X', X_p)) \rightarrow \sim_{\text{eq}_B}^B (R(B), R'(B'))$
Sequential dependencies (SD) [25]	Mono	No	$R(\sim_{\text{succ}(<)}^X \rightarrow \sim_{\text{gap}}^A)$

Table 13: Different quality improving dependencies

Each eCFD  $[R(X \rightarrow Y) \mid t_p]$  is also a decomposable mono-QID  $R(\sim_{\text{eq}_S}^X \rightarrow \sim_{\text{eq}_S}^Y)$ . Here  $\sim_{\text{eq}_S}^X$  is defined as  $\bigcap_{A \in X} \sim_{\text{eq}'_A}^A$ , where  $\sim_{\text{eq}'_A}^A$  is  $\sim_{\text{eq}_A}^A$  when  $t_p[A] = \sqcup$ ,  $\sim_{\text{eq}_S}^A$  when  $t_p[A] = S$ , and it is  $\sim_{\text{eq}_S}^A$  when  $t_p[A] = \bar{S}$  (see Example 1 for the definitions of  $\sim_{\text{eq}_S}^A$  and  $\sim_{\text{eq}_S}^A$ ); similarly for  $\sim_{\text{eq}_S}^Y$ .

**Similarity functional dependencies (SFDs)** [32, 34, 3, 30]. Similarity FDs are a subclass of fuzzy FDs in which equality  $\text{eq}$  is relaxed to similarity relations that are reflexive and symmetric, but not necessarily transitive. More specifically, for each attribute  $A$  in  $R$ , a similarity relation  $\approx_A$  on  $\text{dom}(A)$  is defined. A similarity FD  $R(X \rightarrow Y)$  is specified the same as an FD. It is interpreted along the same lines as FDs, except that it is based on similarity, and its satisfaction may be a ‘‘truth degree’’ in the range  $[0, 1]$  rather than true or false.

Leveraging comparison relations  $\sim_{\approx_X}^X (\sim_{\approx_Y}^Y)$  defined in Section 2, one can verify that each SFD  $R(X \rightarrow Y)$  can be expressed as a decomposable mono-QID  $R(\sim_{\approx_X}^X \rightarrow \sim_{\approx_Y}^Y)$ .

**Domain dependencies (DDs)/Edits** [12, 21]. A DD [12] is of the form  $R(A, S)$ , where  $A$  is an attribute and  $S$  is a set of constants taken from  $\text{dom}(A)$ . It aims to assure that for each tuple  $t$  in an instance of  $R$ ,  $t[A] \in S$ . An edit  $e$  is of the form  $R([A_1 : S_1] \times \dots \times [A_k : S_k])$ , where  $S_i$  is a set of constants in  $\text{dom}(A_i)$  [21]. It states that for any tuple  $t$  in an instance of  $R$ ,  $t[A_1, \dots, A_k] \notin S_1 \times \dots \times S_k$ . Observe that edits can be expressed as eCFDs in which pattern tuples contain no wildcards. In contrast not every such eCFD is an edit rule.

To express DDs and edits as QIDs, we use a relation  $\succ_S = \{(a, b) \in \text{dom}(A) \times \text{dom}(A) \mid a, b \in S\}$ . Intuitively,  $\succ_S$  states that elements in  $S$  are indistinguishable. We also use  $\succ_{\bar{S}} = \text{dom}(A)^2 \setminus \succ_S$ , to identify elements not in  $S$ . We define comparison relation  $\sim_{\succ_S}^A$  and  $\sim_{\succ_{\bar{S}}}^A$  in terms of  $\succ_S$  and  $\succ_{\bar{S}}$ , respectively.

A DD  $R(A, S)$  can be expressed as a decomposable mono-QID  $R(\emptyset \rightarrow \sim_{\succ_S}^A)$ . An edit  $R([A_1 : S_1] \times \dots \times [A_k : S_k])$  can be specified as a decomposable mono-QID  $R(\sim_{\succ_S}^X \rightarrow \sim_{\succ_{\bar{S}_k}}^{A_k})$ , where  $\sim_{\succ_S}^X = \bigcap_{i \in [1, k-1]} \sim_{\succ_{S_i}}^{A_i}$ .

**Association rules with 100% confidence (ARs)** [35, 28]. An association rule (AR) is defined by means of an implication  $R([A_1 : a_1], \dots, [A_k : a_k]) \rightarrow [B : b]$ , where  $A_i$  and  $B$  are attributes, and  $a_i$  and  $b$  are constants from the corresponding domains. An instance  $I$  of  $R$  satisfies the AR with 100% confidence if for each tuple  $t$  in  $I$ , if  $t[A_1, \dots, A_k] = (a_1, \dots, a_k)$  then  $t[B] = b$ . Such association rules are a special case of CFDs in which the pattern tuples consist of constants only. Hence, ARs are decomposable mono-QID defined in terms of  $\sim_{\text{eq}_a}^A$ .

**Order dependencies (ODs)** [22, 29, 23]. Order dependencies (ODs) are specified as standard FDs in which equality may be replaced by a partial order relations  $<, \leq, >, \geq$  on ordered attributes, as well as their complement relations (e.g.,  $\leq^c$ ). More specifically, an OD is of the form  $R(\tilde{X} \rightarrow \tilde{Y})$  in which  $\tilde{X}$  and  $\tilde{Y}$  are sets of attributes that are either unmarked (for  $=$ ), or marked with one of the partial order relations. ODs can be expressed as decomposable mono-QIDs defined in terms of  $\sim_{\text{eq}_A}^A$  (if the attribute  $A$  is unmarked) and  $\sim_{\text{op}}^A$  (if  $A$  is marked).

Here for  $\text{op} \in \{<, \leq, >, \geq, \leq^c\}$ ,  $\sim_{\text{op}}^A$  is defined in terms of relation  $\text{op} = \{(a, b) \in \text{dom}(A) \times \text{dom}(A) \mid a \text{ op } b\}$ .

**Currency dependencies (CDs)** [17]. Currency dependencies (CDs) are constraints of the form

$$\forall t_1, \dots, t_k : R(\bigwedge_{j \in [1, k]} (t_1[\text{eid}] = t_j[\text{eid}] \wedge \psi) \longrightarrow t_u \prec_A t_v),$$

where  $u, v \in [1, k]$ ,  $t_u \prec_A t_v$  indicates that  $t_v[A]$  is more current than  $t_u[A]$ ; and  $\psi$  is a conjunction of predicates of the form  $t_j[A] \text{ op } t_h[A]$  with  $\text{op} \in \{=, \neq, <, \leq, >, \geq, \prec\}$ ,  $t_j[A] = a$  or  $t_j[A] \neq a$  for  $a \in \text{dom}(A)$ . Here  $\text{eid}$  denotes an entity id attribute, as introduced by Codd [8]. The CD is to assert that when  $t_i$ 's are tuples pertaining to the same entity, and if condition  $\psi$  is satisfied, then tuple  $t_v$  must be more up-to-date than  $t_u$  in attribute  $A$ .

When  $k = 2$ , CDs can be expressed as decomposable mono-QIDs of the form  $R(\sim_{B_X}^X \rightarrow \sim_{\prec_A}^A)$ . Here  $\sim_{\prec_A}^A$  is a comparison relation defined in terms of relation  $\prec_A$ , and  $\sim_{B_X}^X$  is an intersection of comparison relations  $\sim_{\text{op}}^E$  for  $\text{op} \in \{<, \leq, >, \geq, =, \neq\}$ , and  $\sim_{\text{eq}_S(\text{op}_a)}^E$  for  $\text{op} \in \{=, \neq\}$ , to express the condition  $\psi$ .

There are also *mono-QIDs* that are *not decomposable*.

**Sequential dependencies (SDs)** [25]. A sequential dependency (SD) is of the form  $R(X \rightarrow_g A)$ , where  $X$  is a set of ordered attributes,  $A$  is a numerical (linear ordered) attribute, and  $g \subseteq \text{dom}(A)$  is an interval. It is to assert that for a predefined permutation  $\pi$  of tuples in an instance  $I$  of  $R$  that are increasing in  $X$ , *i.e.*,  $t_{\pi(1)}[X] <_X t_{\pi(2)}[X] <_X \dots <_X t_{\pi(N)}[X]$ , we have that for any  $i \in [1, N - 1]$ ,  $t_{\pi(i+1)}[A] - t_{\pi(i)}[A] \in g$ . A conditional variant of sequential dependencies is studied [25] in which a sequential dependency is to hold only on a subset of tuples in  $I$ , where the subset is selected by means of admissible ranges of attribute values in  $X$ .

To express SDs as QIDs, we define the following two comparison relations: (1)  $\sim_{\text{gap}}^A$  in terms of relation  $\text{gap} = \{(a, b) \in \text{dom}(A) \times \text{dom}(A) \mid b - a \in g\}$ , which is neither reflexive, symmetric nor transitive; and (2)  $\sim_{\text{succ}(<)}^X = \sim_{\pi}^X \setminus (\sim_{\pi\text{-strict}}^X \cap (\sim_{\pi\text{-strict}}^X)^\circ)$ , where  $\sim_{\pi}^X$  is the comparison relation derived from the predefined order  $\pi$ ,  $\sim_{\pi\text{-strict}}^X = \sim_{\pi}^X \setminus \sim_{\text{eq}}^X$ , and  $S^\circ$  denotes the reversal operation on binary relation  $S$ ; observe that  $\sim_{\text{succ}(<)}^X$  is defined as a Boolean combination of comparison relations, and is not decomposable.

An SD  $R(X \rightarrow_g A)$  can be readily expressed as a non-decomposable mono-QID  $R(\sim_{\text{succ}(<)}^X \rightarrow \sim_{\text{gap}}^A)$ . In addition, for the conditional extension of SDs given in [25], one can construct an appropriate comparison relation  $\sim_{\text{range}}^X$  that restricts  $\sim_{\text{succ}(<)}^X$  and groups tuples together that fall in the admissible range.

We next consider *decomposable bi-QIDs*.

**Binary FDs (BFDs)** [27, 33]. A binary FD is an FD that relates attributes in two different relations. More specifically, a BFD is of the form  $R[X] = R'[X'] \rightarrow R[Y] = R'[Y']$  where  $X, X'$  and  $Y, Y'$  are lists of pairwise compatible attributes in  $R$  and  $R'$ , respectively. For an instance  $(I, I')$  of  $(R, R')$ , it is to assure that for any pair of tuples  $t \in I$  and  $t' \in I'$ , if  $t[X] = t'[X']$  then also  $t[Y] = t'[Y']$ .

We can express BFDs as decomposable bi-QIDs of the form  $\sim_{\text{eq}_X}^X (R(X), R'(X')) \rightarrow \sim_{\text{eq}_Y}^Y (R(Y), R'(Y'))$ , where  $\sim_{\text{eq}_X}^X$  and  $\sim_{\text{eq}_Y}^Y$  are defined in terms of  $\sim_{\text{eq}}^A$  as given earlier for the case of FDs, but across  $R$  and  $R'$ .

**Matching dependencies (MDs)** [18, 13, 4]. Matching dependencies (MDs) were proposed to specify record matching rules in [18]. MDs were originally defined in terms of a dynamic semantics. Here we attempt to redefine MDs as QIDs, and interpret their dynamic semantics in terms of their enforcement strategies in Section 4.

When specified as QIDs, MDs can be viewed as an extension of both SFDs and BFDs. They compare certain attributes of tuples (possibly across different relations) in terms of similarity, and if these attributes are similar enough to each other, then identify some other attributes of the tuples. More specifically, an MD is of the form

$$\left( \bigwedge_{i \in [1, k]} R[A_i] \asymp_i R'[A'_i] \right) \longrightarrow \left( \bigwedge_{j \in [1, \ell]} R[B_j] = R'[B'_j] \right).$$

For an instance  $(I, I')$  of  $(R, R')$ , it assure that for any tuples  $t \in I$  and  $t' \in I'$ , if  $t[A_i] \asymp_i t'[A'_i]$  for  $i \in [1, k]$ , where  $\asymp_i$  is a similarity relation, then  $t[B_j]$  and  $t'[B'_j]$  must be identified for  $j \in [1, \ell]$ , indicated by equality.

MDs can be written as decomposable bi-QIDs of the form  $\sim_{\approx_X}^X (R(X), R'(X')) \rightarrow \sim_{\text{eq}_Y}^Y (R(Y), R'(Y'))$ , where  $\sim_{\approx_X}^X$  and  $\sim_{\text{eq}_Y}^Y$  are given earlier for SFDs and FDs, respectively.

**Editing rules (eRs)** [19]. Editing rules (eRs) were introduced for correcting errors in a data relation ( $R$ ) with accurate values from a master relation ( $R'$ ). An eR is of the form  $[(R(X), R'(X') \rightarrow R(B), R'(B')) \mid t_p[X_p]]$ , where  $X$  and  $X'$  are lists of compatible attributes in schemas  $R$  and  $R'$ , respectively, and  $|X| = |X'|$ . Moreover,  $B$  is an attribute of  $R$  but  $B \notin X$ , and  $B'$  is an attribute of  $R'$  but is not in  $X'$ . Finally,  $t_p$  is a pattern tuple over a set of distinct attributes  $X_p$  in  $R$ , where for each  $A \in X_p$ ,  $t_p[A]$  is either  $\sqcup$  or a constant  $a$  drawn from  $\text{dom}(A)$ , similar to pattern tuples in CFDs. Intuitively, for an instance  $(I, I')$  of  $(R, R')$  and for any tuple  $t \in I$ , if  $t[X_p] \simeq t_p[X_p]$ , where  $\simeq$  is the match relation defined for CFDs, and moreover, if there exists  $t' \in I'$  such that  $t[X] = t'[X']$ , then we update  $t[B]$  by letting  $t[B] := t'[B]$ , *i.e.*, taking the value  $t'[B]$  from the master tuple  $t'$ .

Like MDs, eRs also have a dynamic semantics. Nevertheless, we can express eRs as QIDs in terms of a standard static semantics, by using equality to indicate updates. More specifically, an eR can be written as decomposable bi-QIDs of the form  $\sim_{B_X}^X (R(X, X_p), R'(X', X_p)) \rightarrow \sim_{\text{eq}_B}^B (R(B), R'(B'))$ , where  $\sim_{B_X}^X$  is defined as  $\bigcap_{A \in X} \sim_{\text{eq}_A}^A \cap \bigcap_{E \in X_p} \sim_{\text{eq}'_E}^E$ . Here  $\sim_{\text{eq}_A}^A$  is the equality relation on attributes of tuples across  $R$  and  $R'$ , and  $\sim_{\text{eq}'_E}^E$  is defined along the same lines as its counterpart for CFDs on  $R$  attributes only, to enforce patterns.

## 4 Data Quality Problems

We have proposed QIDs as data quality rules, to catch and fix errors. Below we address several central data quality issues associated with QIDs, and provide an overview of recent advances in the study of these issues.

**Discovering QIDs.** To use QIDs as data quality rules, it is necessary to have techniques in place that can *automatically discover* QIDs from real-life data. Indeed, it is often unrealistic to rely solely on human experts to design data quality rules via an expensive manual process. This suggests that we settle the following problem.

Given a database  $\mathcal{I}$ , the *profiling problem* is to find a *minimal cover* of all QIDs that  $\mathcal{I}$  satisfies, *i.e.*, a non-redundant set of QIDs that is logically equivalent to the set of all QIDs that hold on  $\mathcal{I}$ . We want to learn informative and interesting data quality rules from (possibly dirty) data, and prune away insignificant rules.

Several algorithms have been developed for discovering FDs (*e.g.*, [26]), CFDs [7, 15, 24] and MDs [31]. These algorithms allow a considerable pruning of the search space when QIDs are decomposable.

**Validating QIDs.** A given set  $\Sigma$  of QIDs, either automatically discovered or manually designed, may be inconsistent itself. In light of this we have to find those QIDs from  $\Sigma$  that are consistent and non-redundant, to be used as data quality rules. These highlight the need for studying the following classical problems for dependencies.

The *satisfiability problem* for QIDs is to determine, given a set  $\Sigma$  of QIDs defined on a relational schema  $\mathcal{R}$ , whether there exists a nonempty instance  $\mathcal{I}$  of  $\mathcal{R}$  such that  $\mathcal{I} \models \Sigma$ , *i.e.*,  $\mathcal{I} \models \varphi$  for *each* QID  $\varphi \in \Sigma$ .

The *implication problem* for QIDs is to decide, given a set  $\Sigma$  of QIDs and a single QID  $\varphi$  defined on a relational schema  $\mathcal{R}$ , whether  $\Sigma$  entails  $\varphi$ , *i.e.*, whether for all instances  $\mathcal{I}$  of  $\mathcal{R}$ , if  $\mathcal{I} \models \Sigma$  then  $\mathcal{I} \models \varphi$ . Effective implication analysis allows us to remove redundancies and deduce new data quality rules from  $\Sigma$ .

These problems have been studied for FDs [1], CFDs [14], SFDs [3], edits [21], ODs [29], BFDs [27], MDs [13] and editing rules [19]. These issues are, however, nontrivial. When CFDs are concerned, for instance, both problems are intractable [14]. This calls for the study of effective heuristic algorithms for their analyses.

**Error detection (localization).** After a consistent set of data quality rules is identified, the next question concerns how to effectively catch errors in a database by using the rules. Given a set  $\Sigma$  of QIDs and a database  $\mathcal{I}$ , we want to *detect inconsistencies* in  $\mathcal{I}$ , *i.e.*, to find all tuples and attributes in  $\mathcal{I}$  that violate some QID in  $\Sigma$ . Error detection can be done more efficiently for decomposable mono-QIDs due to the absence of intra-relational joins.

Error detection methods have been developed for centralized databases [14] and distributed data [16], based on CFDs. Nevertheless, it is rather challenging to precisely localize errors in the data. Dependencies typically help us catch a pair of tuples (or attributes) that are inconsistent with each other, but fall short of telling us which attribute is erroneous. As an example, consider a simple CFD  $R[(CC, AC \rightarrow \text{city}) \mid (44, 131, \text{Edi})]$ , which states that in the UK ( $CC = 44$ ), when the area code ( $AC$ ) is 131, then the city must be Edinburgh. A tuple  $t$ : ( $CC = 44$ ,

AC = 131, city = Ldn) violates this CFD: AC is 131 but the city is London. The CFD finds that  $t[\text{CC}]$ ,  $t[\text{AC}]$  or  $t[\text{city}]$  is incorrect, *i.e.*, an error is present in the tuple. However, it does not tell us which of the three attributes is wrong and to what value it should be changed. To rectify this problem, we need to adopt a stronger semantics when enforcing QIDs for data imputation, as will be seen below when we discuss editing rules [19].

**Data repairing (data imputation).** After the errors are detected, we want to automatically correct the errors and make the data consistent. This is known as *data repairing* [2]. Given a set  $\Sigma$  of QIDs and a database instance  $\mathcal{I}$  of  $\mathcal{R}$ , it is to find a candidate *repair* of  $\mathcal{I}$ , *i.e.*, another instance  $\mathcal{I}'$  of  $\mathcal{R}$  such that  $\mathcal{I}' \models \Sigma$  and  $\text{cost}(\mathcal{I}, \mathcal{I}')$  is minimum. Here  $\text{cost}()$  is a metric defined in terms of (a) the distance between  $\mathcal{I}$  and  $\mathcal{I}'$  and (b) the accuracy of attribute values in  $\mathcal{I}$  (see, *e.g.*, [5, 9]). When a value  $v$  in  $\mathcal{I}$  is changed to  $v'$  in  $\mathcal{I}'$  in data repairing, the more accurate the original value  $v$  is and the more distant the new value  $v'$  is from  $v$ , the higher the cost of the change is. That is, we want to avoid changing accurate values in  $\mathcal{I}$ , and want the repair  $\mathcal{I}'$  to *minimally differ* from  $\mathcal{I}$ .

No matter how important, the data repairing problem is nontrivial: it is intractable even when only a fixed set of FDs is considered [5]. In light of this, several heuristic algorithms have been developed, to repair data based on different strategies for enforcing FDs, CFDs, MDs and eRs [5, 9, 19, 20].

These QID enforcement strategies can be uniformly characterized in terms of a revised chase process (see, *e.g.*, [1] for details about chase). More specifically, assume that for each attribute  $A \in \mathcal{R}$ , a *fixing function*  $\sigma_A : \text{dom}(\text{tid}) \rightarrow \text{dom}(A)$  is defined. Given a tuple  $t$  in  $\mathcal{I}$  identified by  $\text{tid}$ , the function is used to change the values of  $t[A]$  to the value  $\sigma_A(\text{tid})$ . Let  $\bar{\sigma}$  denote a set of fixing functions  $\sigma_A$ , one for each  $A \in \mathcal{R}$ . Consider a QID  $\varphi = (\sim_{\mathcal{B}_X}^X (R(X), R'(X')) \rightarrow \sim_{\mathcal{B}_Y}^Y (R(Y), R'(Y')))$ , and two tuples  $t_1, t_2$  in  $\mathcal{I}$ . We say that an instance  $\mathcal{I}'$  of  $\mathcal{R}$  is the *result of enforcing  $\varphi$  on  $t_1$  and  $t_2$  in  $\mathcal{I}$  using  $\bar{\sigma}$* , denoted by  $\mathcal{I} \xrightarrow{\bar{\sigma}_{\varphi, (t_1, t_2)}} \mathcal{I}'$ , if

- in  $\mathcal{I}$ : we have that  $\sim_{\mathcal{B}_X}^X (t_1[\text{tid}], t_2[\text{tid}])$  holds, whereas  $\sim_{\mathcal{B}_Y}^Y (t_1[\text{tid}], t_2[\text{tid}])$  does not hold;
- in  $\mathcal{I}'$ :  $t_1[A_i] = \sigma_{A_i}(t_1[\text{tid}])$ ,  $t_2[A_i] = \sigma_{A_i}(t_2[\text{tid}])$  for all  $A_i \in X \cup Y$ , and both  $\sim_{\mathcal{B}_X}^X (t_1[\text{tid}], t_2[\text{tid}])$  and  $\sim_{\mathcal{B}_Y}^Y (t_1[\text{tid}], t_2[\text{tid}])$  hold; and finally,
- $\mathcal{I}$  and  $\mathcal{I}'$  agree on every other tuple and attribute value.

For a set  $\Sigma$  of QIDs, we say that  $\mathcal{I}'$  is a  $\Sigma$ -*repair* of  $\mathcal{I}$  if there is a finite *chasing sequence*  $\mathcal{I}_0 \dots, \mathcal{I}_k$  such that

$$\mathcal{I}_0 = \mathcal{I} \xrightarrow{\bar{\sigma}_{\varphi^1, (t_1^1, t_2^1)}} \mathcal{I}_1 \xrightarrow{\bar{\sigma}_{\varphi^2, (t_1^2, t_2^2)}} \dots \xrightarrow{\bar{\sigma}_{\varphi^k, (t_1^k, t_2^k)}} \mathcal{I}_k = \mathcal{I}'$$

and furthermore,  $\mathcal{I}' \models \Sigma$ , where  $\varphi^i \in \Sigma$  and  $\bar{\sigma}^i$  is a set of fixing functions for  $i \in [1, k]$ . A  $\Sigma$ -repair  $\mathcal{I}'$  of  $\mathcal{I}$  is minimal if  $\text{cost}(\mathcal{I}, \mathcal{I}')$  is minimum. We say that an enforcement strategy of QIDs is *Church-Rosser* if for any set  $\Sigma$  of QIDs, all  $\Sigma$ -repairs obtained via different finite chasing sequences result in the same instance of  $\mathcal{R}$ .

A repairing algorithm is essentially an implementation of the “chase”-process, determined by (a) how fixing functions are defined, (b) in what order the QIDs in  $\Sigma$  are enforced, and (c) for each QID  $\varphi$ , how it is enforced on the tuples in  $\mathcal{I}$  that violate  $\varphi$ . Below we outline how existing algorithms compute  $\Sigma$ -repairs by enforcing QIDs. Here decomposable QIDs again simplify the chase process because fixing functions can be defined independently for different attributes in a decomposable QID.

(Conditional) Functional dependencies [5, 9]. Given a set  $\Sigma$  of FDs and an instance  $\mathcal{I}$ , the repairing algorithm of [5] computes a  $\Sigma$ -repair by enforcing the FDs in  $\Sigma$  one by one, as follows. (1) For each FD  $\varphi = R(\sim_{\text{eq}_X}^X \rightarrow \sim_{\text{eq}_B}^B)$  in  $\Sigma$  and each tuple  $t$  in  $\mathcal{I}$ , it finds the set  $V_{\varphi, t}$  of all tuples  $t'$  in  $\mathcal{I}$  such that  $\sim_{\text{eq}_X}^X (t[\text{tid}], t'[\text{tid}])$  but not  $\sim_{\text{eq}_Y}^Y (t[\text{tid}], t'[\text{tid}])$ . (2) It defines fixing functions  $\bar{\sigma}$  such that for all  $B \in Y$  and for all  $t' \in V_{\varphi, t} \cup \{t\}$ ,  $\sigma_B(t[\text{tid}])$  is the same constant  $c$ , which is a value drawn from  $\{t''[B] \mid t'' \in V_{\varphi, t} \cup \{t\}\}$  such that the total cost of changing  $t'[B]$  to  $c$  is minimum. That is, it enforces  $\varphi$  on all tuples in  $V_{\varphi, t}$  consecutively, by changing the values of their attributes in the right-hand side of  $\varphi$  to the same value. Furthermore, in order for the repairing process to terminate, a *hard constraint* is imposed, which requires that once  $\sigma_B$  identifies attributes of a pair of tuples, these attributes remain identified throughout the repairing process. Note that this still allows adjustments to the values of those attributes at a later stage, *i.e.*, different fixing functions  $\bar{\sigma}^i$  may be used in the process. This enforcement strategy guarantees to find a  $\Sigma$ -repair for a set of FDs, but is not Church-Rosser.

This is extended by the algorithm of [9] for CFDs, in which fixing functions allow attribute values to be changed to a constant specified by a CFD, provided that the hard constraints imposed so far is not violated. For instance, tuple  $t : (CC = 44, AC = 131, city = Ldn)$  can be changed to  $(CC = 44, AC = 131, city = Edi)$  as specified by CFD  $R[(CC, AC \rightarrow city) \mid (44, 131, Edi)]$ , unless Ldn was assigned earlier by a fixing function  $\sigma_{city}$ . In the latter case,  $\sigma_{city}$  will not make any change to the right-hand side attribute, but instead a fixing function on the left-hand side attributes is applied; *e.g.*,  $t$  could be changed to  $(CC = 44, AC = 020, city = Ldn)$  by  $\sigma_{AC}$ . Termination is also assured, but the enforcement strategy is not Church-Rosser.

*Matching dependencies* [4, 18]. MDs were proposed in [18] to infer matching rules. An enforcement strategy was introduced in [4] for repairing data with MDs. It enforces a set  $\Sigma$  of MDs one by one on an instance  $\mathcal{I}$  as follows. Given an MD  $\varphi = \sim_{\approx_X}^X (R(X), R'(X')) \rightarrow \sim_{eq_Y}^Y (R(Y), R'(Y'))$ , for each  $B \in Y$ , it employs matching functions  $m_B : dom(B) \times dom(B) \rightarrow dom(B)$  that are idempotent ( $m_B(b, b) = b$ ), commutative ( $m_B(b, b') = m_B(b', b)$ ) and associative ( $m_B(b, m_B(b', b'')) = m_B(m_B(b, b'), b'')$ ). For a pair  $(t_1, t_2)$  of tuples in  $\mathcal{I}$  that violate  $\varphi$ , it defines fixing functions  $\sigma_B$  such that  $\sigma_B(t_1[tid]) = \sigma_B(t_2[tid]) = m_B(t_1[B], t_2[B])$ , and hence, computes  $\mathcal{I} \rightarrow_{\varphi, (t_1, t_2)}^{\bar{\sigma}} \mathcal{I}'$ , where  $\bar{\sigma}$  consists of  $\sigma_B$  for all  $B \in Y$ . It is shown in [4] that  $m_B$  imposes a bounded lattice-structure on  $dom(B)$  and that  $\mathcal{I}'$  is an instance that dominates  $\mathcal{I}$  in the corresponding bounded lattice on instances of  $\mathcal{R}$ . From this follows the termination of the process. The strategy is not Church-Rosser.

Employing both CFDs and MDs, a data repairing algorithm was developed in [20], with fixing functions defined for CFDs and MDs as above. It enforces a set of CFDs and MDs by interleaving data repairing and record matching operations; by leveraging their interaction, the algorithm improves the accuracy of repairs generated.

*Editing rules* [19]. As remarked earlier, dependencies typically only detect the presence of errors, but do not locate precisely what attributes are wrong and how to fix the errors. As a result, repairing algorithms often fail to correct errors and worse still, may even introduce new errors in the repairing process. To rectify this problem, a repairing strategy was proposed in [19] that aims to find a certain fix of a tuple, *i.e.*, a fix that is guaranteed correct, by using editing rules (eRs) and master data. Consider an eR  $\varphi = \sim_{B_X}^X (R(X, X_p), R'(X', X_p)) \rightarrow \sim_{eq_B}^B (R(B), R'(B'))$ , a tuple  $t_1$  in a data instance  $I$  of  $R$ , and a tuple  $t_2$  in a master relation of schema  $R'$  that provides consistent and correct data values. It defines a fixing function  $\sigma_B$  that equalizes  $t_1[B]$  and  $t_2[B]$  by always setting  $\sigma_B(t_1[tid])$  and  $\sigma_B(t_2[tid])$  to the master value  $t_2[B]$ . Furthermore, it allows  $\mathcal{I} \rightarrow_{\varphi, (t_1, t_2)}^{\sigma_B} \mathcal{I}'$  by enforcing  $\varphi$  only if (a)  $t_1[X, X_p]$  has been assured correct (*i.e.*, validated), either by users or via inference from fixes generated earlier; and (b) enforcing different eRs on  $t_1$  and master tuples yields the same repair. In addition, after an attribute  $t_1[B]$  is updated, it is no longer changed since it is already validated.

For instance, tuple  $t : (CC = 44, AC = 131, city = Ldn)$  is changed to  $(CC = 44, AC = 131, city = Edi)$  only if (a)  $t[CC, AC, city]$  is already validated, and (b) for all eRs  $\varphi$  in a given set  $\Sigma$  and all master tuples  $t'$  in  $I'$ , either  $\varphi$  does not apply to  $t$  and  $t'$ , or  $\varphi$  requires that  $t[city]$  and  $t'[city]$  be equalized and moreover,  $t'[city] = Edi$ .

From these it follows that the enforcement strategy guarantees the termination of repairing processes. Furthermore, it is per definition Church-Rosser and assures that the fixes generated are certain.

These algorithms are just examples of the revised chase. We expect that different QID enforcement strategies could be developed along the same lines, and yield data repairing algorithms for a wider range of applications.

**Acknowledgments.** Fan is supported in part by the RSE-NSFC Joint Project Scheme, an IBM scalable data analytics for a smarter planet innovation award, and the National Basic Research Program of China (973 Program) 2012CB316200.

## References

- [1] S. Abiteboul, R. Hull, and V. Vianu. *Foundations of Databases*. Addison-Wesley, 1995.
- [2] M. Arenas, L. E. Bertossi, and J. Chomicki. Consistent query answers in inconsistent databases. In *PODS*, 1999.
- [3] R. Belohlávek and V. Vychodil. Data dependencies in Codd's relational model with similarities. In *Handbook of Research on Fuzzy Information Processing in Databases*. IGI Global, 2008.



- [4] L. E. Bertossi, S. Kolahi, and L. V. S. Lakshmanan. Data cleaning and query answering with matching dependencies and matching functions. In *ICDT*, 2011.
- [5] P. Bohannon, W. Fan, M. Flaster, and R. Rastogi. A cost-based model and effective heuristic for repairing constraints by value modification. In *SIGMOD*, 2005.
- [6] L. Bravo, W. Fan, F. Geerts, and S. Ma. Increasing the expressivity of conditional functional dependencies without extra complexity. In *ICDE*, 2008.
- [7] F. Chiang and R. Miller. Discovering data quality rules. In *VLDB*, 2008.
- [8] E. F. Codd. Extending the database relational model to capture more meaning. *TODS*, 4(4):397–434, 1979.
- [9] G. Cong, W. Fan, F. Geerts, X. Jia, and S. Ma. Improving data quality: Consistency and accuracy. In *VLDB*, 2007.
- [10] M. M. Deza and E. Deza. *Encyclopedia of Distances*. Springer, 2009.
- [11] W. W. Eckerson. Data quality and the bottom line: Achieving business success through a commitment to high quality data. The Data Warehousing Institute, 2002.
- [12] R. Fagin. A normal form for relational databases that is based on domains and keys. *TODS*, 6:387–415, 1981.
- [13] W. Fan, H. Gao, X. Jia, J. Li, and S. Ma. Dynamic constraints for record matching. *The VLDB Journal*, 2011.
- [14] W. Fan, F. Geerts, X. Jia, and A. Kementsietsidis. Conditional functional dependencies for capturing data inconsistencies. *TODS*, 33(2), 2008.
- [15] W. Fan, F. Geerts, J. Li, and M. Xiong. Discovering conditional functional dependencies. *TKDE*, 23(5):683–698, 2011.
- [16] W. Fan, F. Geerts, S. Ma, and H. Müller. Detecting inconsistencies in distributed data. In *ICDE*, 2010.
- [17] W. Fan, F. Geerts, and J. Wijsen. Determining the currency of data. In *PODS*, 2011.
- [18] W. Fan, X. Jia, J. Li, and S. Ma. Reasoning about record matching rules. *PVLDB*, 2(1):407–418, 2009.
- [19] W. Fan, J. Li, S. Ma, N. Tang, and W. Yu. Towards certain fixes with editing rules and master data. *PVLDB*, 3(1):173–184, 2010.
- [20] W. Fan, J. Li, S. Ma, N. Tang, and W. Yu. Interaction between record matching and data repairing. In *SIGMOD*, 2011.
- [21] I. P. Fellegi and D. Holt. A systematic approach to automatic edit and imputation. *Journal of the American Statistical Association*, 71(353):17–35, 1976.
- [22] S. Ginsburg and R. Hull. Order dependency in the relational model. *TCS*, 26(1-2):149 – 195, 1983.
- [23] S. Ginsburg and R. Hull. Sort sets in the relational model. *J. ACM*, 33:465–488, 1986.
- [24] L. Golab, H. Karloff, F. Korn, D. Srivastava, and B. Yu. On generating near-optimal tableaux for conditional functional dependencies. In *VLDB*, 2008.
- [25] L. Golab, H. J. Karloff, F. Korn, A. Saha, and D. Srivastava. Sequential dependencies. *PVLDB*, 2(1):574–585, 2009.
- [26] Y. Huhtala, J. Kärkkäinen, P. Porkka, and H. Toivonen. TANE: An efficient algorithm for discovering functional and approximate dependencies. *Comput. J.*, 42(2):100–111, 1999.
- [27] E. Komissartschik. Restructuring and dependencies in databases. In *MFDBS*, 1989.
- [28] M. Luxenburger. Implications partielles dans un contexte. *Mathématiques, Informatique et Sciences Humaines*, 29(113):35–55, 1991.
- [29] W. Ng. Ordered functional dependencies in relational databases. *Inf. Syst.*, 24(7):535 – 554, 1999.
- [30] K. V. S. V. N. Raju and A. K. Majumdar. Fuzzy functional dependencies and lossless join decomposition of fuzzy relational database systems. *TODS*, 13:129–166, 1988.
- [31] S. Song and L. Chen. Discovering matching dependencies. In *CIKM*, 2009.
- [32] M. I. Sözat and A. Yazici. A complete axiomatization for fuzzy functional and multivalued dependencies in fuzzy database relations. *Fuzzy Sets Syst.*, 117:161–181, 2001.
- [33] J. Wang. Binary equality implication constraints, normal forms and data redundancy. *Inf. Process. Lett.*, 101(1):20–25, 2007.
- [34] A. Yazici, E. Gocmen, B. Buckles, R. George, and F. Petry. An integrity constraint for a fuzzy relational database. In *FUZZ*, 1993.
- [35] M. J. Zaki. Mining non-redundant association rules. *Data Min. Knowl. Discov.*, 9(3):223–248, 2004.