

Information-Based Distance Measures and the Canonical Reflection of View Updates

Stephen J. Hegner
Umeå University
Department of Computing Science
SE-901 87 Umeå, Sweden
hegner@cs.umu.se
<http://www.cs.umu.se/~hegner>

Abstract

For the problem of reflecting an update on a database view to the main schema, the constant-complement strategies are precisely those which avoid all update anomalies, and so define the gold standard for well-behaved solutions to the problem. However, the families of view updates which are supported under such strategies are limited, so it is sometimes necessary to go beyond them, albeit in a systematic fashion. In this work, an investigation of such extended strategies is initiated for relational schemata. The approach is to characterize the information content of a database instance, and then require that the optimal reflection of a view update to the main schema embody the least possible change of information. The key property is identified to be *strong monotonicity* of the view, meaning that view insertions may always be reflected as insertions to the main schema, and likewise for deletions. In that context it is shown that for insertions and deletions, an optimal update, entailing the least change of information, exists and is unique up to isomorphism for wide classes of constraints.

1 Introduction

1.1 The limitations of constant complement The problem of reflecting view updates to the main schema of a database system is a difficult one whose solution invariably involves compromise. The constant-complement approach [BS81] avoids all so-called update anomalies [Heg04], and so is the gold standard for well-behaved strategies. On the other hand, it is also quite conservative regarding the updates which it admits. A short example will help to illustrate the scope of this approach. Let \mathbf{E}_0 be the relational schema consisting of the single relation symbol $R[ABC]$, constrained by the join dependency $\bowtie [AB, BC]$, and let $\Pi_{AB}^{\mathbf{E}_0}$ be the view whose single relation symbol is $R_{AB}[AB]$ and whose morphism $\pi_{AB}^{\mathbf{E}_0}$ is the projection of $R[ABC]$ onto AB . In the constant-complement strategy, all updates to $\Pi_{AB}^{\mathbf{E}_0}$ must hold a so-called *complementary view* fixed. The natural complement to $\Pi_{AB}^{\mathbf{E}_0}$ is the view $\Pi_{BC}^{\mathbf{E}_0}$, defined by the projection of $R[ABC]$ onto $R_{BC}[BC]$. It is easy to see that the updates to $\Pi_{AB}^{\mathbf{E}_0}$ which hold $\Pi_{BC}^{\mathbf{E}_0}$ fixed are precisely those which hold the projection onto B fixed. Thus, for example, if

$\{R(a_0, b_0, c_0), R(a_1, b_1, c_1)\}$ is the current instance of \mathbf{E}_0 , so that the instance of the view schema is $\{R_{AB}(a_0, b_0), R_{AB}(a_1, b_1)\}$, then insertion of $R_{AB}(a_2, b_1)$ is realized by inserting $R(a_2, b_1, c_1)$ into the instance of the main schema. Unfortunately, even so simple an update as inserting $R_{AB}(a_2, b_2)$ into $\Pi_{AB}^{\mathbf{E}_0}$ is not supported, since the projection onto B cannot be held fixed under such an update.

Of course, there is a reason for this limitation. In order to insert $R_{AB}(a_2, b_2)$ into the view, it is necessary to insert some tuple of the form $R(a_2, b_2, c_x)$ into the main schema, with information about c_x not visible within the view. Such an insertion would violate the principle that views be encapsulated with respect to the updates which are allowed, in the sense that the effect of all such updates be contained entirely within the view itself. It is precisely the constant-complement strategy which guarantees this sort of encapsulation [Heg04, Sec. 1.2]. Nevertheless, there are certainly situations in which it is desirable, if not necessary, to lift this limitation in a controlled manner. The goal of this paper is to develop an extension to the constant-complement strategy which admits a wider class of view updates while preserving as many of the desirable properties of the original strategy. In particular, the following three properties are regarded as uncompromisable.

Invariance of admissibility: The admissibility of a view update must depend only upon the current view instance, and not upon the instance of the main schema which gave rise to it.

Canonicity of reflections: All allowable reflections of a view update to the main schema must be equivalent up to some natural notion of isomorphism.

Reflection of monotonicity: If the view mapping is monotonic, then every insertion (resp. deletion) on the view must be reflected as an insertion (resp. deletion) on the main schema.

These conditions are all natural extensions of that which is expected from a constant complement strategy. Invariance of admissibility is always satisfied by a constant-complement strategy; see (cc:1) of [Heg04, Sec. 1.2]. While the reflection defined by a constant-complement strategy may in principle depend upon the choice of complement, it has been shown that it is in fact independent of that choice when the view morphism is monotonic [Heg04, Thm. 4.3] [Heg08c, Cor. 4.24]. Reflection of monotonicity is similarly guaranteed in wide variety of circumstances; see [Heg04, Def. 3.1 (upt:6)] and [Heg08c].

Much of the existing work on this problem, such as [DB82], [Kel85], [Lan90], [BL97], and [BL98], focuses upon translation of view updates via the relational algebra. As such, while they provide useful insight into commonly occurring problem instances, and all support the simple update problem of inserting $R_{AB}(a_2, b_2)$ sketched above, they do not provide a unified theory of how and under which circumstances view updates could be reflected. While such a detailed comparison is a very interesting topic, it must be left to a different paper.

1.2 Database repairs, distance measures, and information content More closely related to the approach developed here is one which has been developed in the logic-programming community – database repair. Roughly, in the repair problem, a database M is given which is inconsistent with respect to a set Ψ of constraints, the task being to “repair” M to a consistent version M' which obeys the constraints in Ψ . It is straightforward to recast a view-update problem in this context, at least in principle. Extending the example of 1.1,

let \mathbf{E}'_0 be the schema which augments \mathbf{E}_0 with the relation symbol $R_{AB}[AB]$ together with its defining constraint $(\forall x)(\forall y)(R_{AB}(x, y) \Leftrightarrow (\exists z)(R(x, y, z)))$. Thus, the current instance of this schema is $M_{\text{old}} = \{R(a_0, b_0, c_0), R(a_1, b_1, c_1), R_{AB}(a_0, b_0), R_{AB}(a_1, b_1)\}$. The “defective” new instance which includes the inserted view tuple $R_{AB}(a_2, b_2)$ but not the corresponding update to the main schema is $M'_{\text{def}} = M_{\text{old}} \cup \{R_{AB}(a_2, b_2)\}$. The task is to repair M'_{def} to be a legal instance, *i.e.*, one which satisfies the integrity constraints of \mathbf{E}_0 , subject to the condition that the instance of R_{AB} is held constant. A ranking function measures the quality of the of various solutions. For any two sets S_1 and S_2 , let $\text{SymDiff}\langle S_1, S_2 \rangle$ denote their symmetric difference $(S_1 \setminus S_2) \cup (S_2 \setminus S_1)$, and let $\text{Card}(S_i)$ denote the cardinality of S_i . Two principal measures are *subset ranking* in which M_{new} is preferred to M'_{new} if $\text{SymDiff}\langle M_{\text{old}}, M_{\text{new}} \rangle \subseteq \text{SymDiff}\langle M_{\text{old}}, M'_{\text{new}} \rangle$ and the *count ranking* in which M_{new} is preferred to M'_{new} if $\text{Card}(\text{SymDiff}\langle M_{\text{old}}, M_{\text{new}} \rangle) \leq \text{Card}(\text{SymDiff}\langle M_{\text{old}}, M'_{\text{new}} \rangle)$. Under either of these measures, the optimal solutions are precisely those of the form $\text{Insert}\langle R(a_2, b_2, v) \rangle$, with v any allowable value for the domain C . See [ADNB06] for further details, as well as a comprehensive list of other papers which employ related approaches.

A drawback of both of these measures is that tuple similarity is an all-or-nothing affair — all tuples which are not identical are equally different from one another. Consequently, $R(a_0, b_0, c_1)$ is just as different from $R(a_0, b_0, c_0)$ as is $R(a_1, b_1, c_1)$. Recently, more sophisticated distance measures have been proposed [ADB07], some of which are based upon (pseudo-)distance measures of individual tuples, such as those proposed in [Hut97] or [NC97]. These tuple-based measures may then be extended to sets of tuples via measures such as that of Eiter and Mannila [EM97], which defines the distance between database instances M_1 and M_2 in terms of the distances between tuples to be

$$\text{Dist}\langle M_1, M_2 \rangle = \frac{1}{2} \cdot \left(\sum_{t_1 \in M_1} \min_{t_2 \in M_2} \text{Dist}\langle t_1, t_2 \rangle + \sum_{t_2 \in M_2} \min_{t_1 \in M_1} \text{Dist}\langle t_1, t_2 \rangle \right).$$

Using a model of distance between tuples such as that of [NC97], which defines such distance in terms of how much each term of one tuple differs from the corresponding term in the other, and assuming that distinct terms have positive distance from one another, it follows that $\text{Dist}\langle R(a_2, b_2, c_1), R(a_1, b_1, c_1) \rangle < \text{Dist}\langle R(a_2, b_2, c_2), R(a_1, b_1, c_1) \rangle$ and so $\text{Insert}\langle R(a_2, b_2, c_1) \rangle$ is preferred to $\text{Insert}\langle R(a_2, b_2, c_2) \rangle$ in support of the view update request $\text{Insert}\langle R_{AB}(a_2, b_2) \rangle$. Similarly, $\text{Insert}\langle R(a_2, b_2, c_0) \rangle$ is preferred to $\text{Insert}\langle R(a_2, b_2, c_2) \rangle$, while $\text{Insert}\langle R(a_2, b_2, c_0) \rangle$ and $\text{Insert}\langle R(a_2, b_2, c_1) \rangle$ are of equal preference.

Despite the obvious attractiveness of such tuple-based metrics from both a mathematical and an aesthetic point of view, in this paper it is argued that quite a different type of metric is more appropriate — one which in fact prefers $\text{Insert}\langle R(a_2, b_2, c_2) \rangle$ to both $\text{Insert}\langle R(a_2, b_2, c_1) \rangle$ and $\text{Insert}\langle R(a_2, b_2, c_0) \rangle$ — precisely the opposite of that which the above tuple-based metric advises. The idea is to measure the information content of database instances relative to a set of sentences, and to prefer updates which involve less change of information. More precisely, relative to a set Φ of sentences, the *information content* $\text{Info}\langle M, \Phi \rangle$ of the database instance M is the set of all sentences in Φ which are satisfied by M . The *information distance* relative to Φ between database instances M_1 and M_2 is then $\Delta\langle (M_1, M_2), \Phi \rangle = \text{SymDiff}\langle \text{Info}\langle M_1, \Phi \rangle, \text{Info}\langle M_2, \Phi \rangle \rangle$. The information distance between M_1 and M_2 is thus not a number but rather the set of formulas of Φ on which M_1 and M_2 differ. Preference of repairs is then defined in the obvious way, with changing M_1 to M_2 preferred to changing M_1 to M_3 if $\Delta\langle (M_1, M_2), \Phi \rangle \subseteq \Delta\langle (M_1, M_3), \Phi \rangle$.

The utility of this approach depends, of course, upon a suitable choice for Φ . If Φ chosen to be the set of all atoms for the schema, then information distance reduces to subset ranking as defined above. The most suitable choice in many situations is to let Φ be the set of all sentences in the language of the schema which are existential (no universal quantification), positive (no negation of any kind), and conjunctive (no disjunction), and which employ at most the constant symbols which occur in the formulation of the update itself — those which occur in the current instances of the main schema and the view as well as those which occur in the proposed new instance of the view. In the context of repairs, if the proposed update is an insertion, then this amounts to just the constant symbols which occur in the database to be repaired. On the other hand, if deletion of tuples is also allowed, then the constant symbols in tuples to be deleted must also be included. Updates will be formalized in detail in the core of this paper; for now, it suffices to write $\text{ConstSym}(\mathbf{u})$ for the set of all constant symbols which occur in the update request \mathbf{u} . Then, using the notation to be introduced in the next section, the sentences of importance are $\text{WFS}(\mathbf{D}, \exists\wedge+, \text{ConstSym}(\mathbf{u}))$, with \mathbf{D} the database schema under consideration.

Returning to the example based upon \mathbf{E}'_0 , which augments \mathbf{E}_0 with the relation symbol $R_{AB}[AB]$ together with its defining constraint, insertion of $R_{AB}(a_2, b_2)$ into the view requires that the sentence $(\exists z)(R(a_2, b_2, z))$ be true in the main schema. This may be satisfied by binding z to any available constant and adding the corresponding tuple. However, if it is bound to a constant in $\text{ConstSym}(M_{\text{def}'}) = \{a_0, a_1, a_2, b_0, b_1, b_2, c_0, c_1\}$, then additional sentences will appear in the information content of the new instance. For example, if z is bound to c_1 , then the sentence $R_{AB}(a_2, b_2, c_1)$ will also be in the information content of the new instance, as would the sentence $(\exists z)(R(a_1, b_1, z) \wedge R(a_2, b_2, z))$. Neither would be true were $R(a_2, b_2, c_2)$ inserted instead. Thus, to add the least possible information relative to sentences which only involve constants which are already used, z must be bound to a constant which does not occur in $\text{ConstSym}(M_{\text{def}'})$. Intuitively, this corresponds to binding it to a generic constant, and not one which also plays some other rôle in the database.

By itself, this information-based approach does not ensure tuple minimality. The insertion of both $R(a_2, b_2, c_2)$ and $R(a_2, b_2, c_3)$ into M_{old} produces exactly the same added information as does the insertion of either alone, since each tuple adds precisely $(\exists z)(R(a_2, b_2, z))$ to the information content. To remedy this, tuple minimality (*i.e.*, minimality with respect to subset ranking) is also required. Thus, an optimal repair, whenever it exists, must be tuple minimal as well as minimal with respect to information change.

1.3 Further example To illustrate the ideas of information-based optimization of updates more completely, a slightly more complex example is presented. Let \mathbf{E}_1 be the relational schema with relations $R[ABC]$ and $S[CD]$, constrained by the inclusion dependency $R[C] \sqsubseteq S[C]$. Regard a database as a set of ground atoms over the associated logic. For example, $M_{00} = \{R(a_0, b_0, c_0), R(a_1, b_1, c_1), S(c_0, d_0), S(c_1, d_1)\}$ is such a database. Now, let K be a set of constants in the underlying logical language, regarded as domain elements for this schema. For information content, the base set Φ of sentences is $\text{WFS}(\mathbf{E}_1, \exists\wedge+, K)$, the set of all positive (*i.e.*, no negation, explicit or implicit), existential, and conjunctive sentences in the language of the schema \mathbf{E}_1 which involve at most the constant symbols in K . Relative to this set, the information content of M is the set of all sentences in $\text{WFS}(\mathbf{E}_1, \exists\wedge+, K)$ which are implied by M . The central step is to choose K properly. Using the notation to be introduced in 3.2, this

information content is denoted $\text{Info}\langle M, \text{WFS}(\mathbf{E}_1, \exists\wedge+, K) \rangle$.

A *cover* for this information content is a subset $\Psi \subseteq \text{Info}\langle M, \text{WFS}(\mathbf{E}_1, \exists\wedge+, K) \rangle$ such that Ψ and $\text{Info}\langle M, \text{WFS}(\mathbf{E}_1, \exists\wedge+, K) \rangle$ are logically equivalent. For $K_{00} = \{a_0, a_1, b_0, b_1, c_0, c_1, d_0, d_1\}$, the set of all constant symbols of M_{00} , the set M_{00} itself is clearly a cover for $\text{Info}\langle M_{00}, \text{WFS}(\mathbf{E}_3, \exists\wedge+, K_{00}) \rangle$. On the other hand, with $K'_{00} = \{a_0, a_1, b_0, b_1, c_0, d_0\}$, a cover for $\text{Info}\langle M_{00}, \text{WFS}(\mathbf{E}_1, \exists\wedge+, K'_{00}) \rangle$ is $\{R(a_0, b_0, c_0), S(c_0, d_0), (\exists x)(\exists y)(R(a_1, b_1, x) \wedge S(x, y))\}$. Note that the constants in $K_{00} \setminus K'_{00}$ have been replaced by existentially quantified variables.

To see how this idea is useful in the context of view updates, let $\Pi_{R_{AB}[AB]}^{\mathbf{E}_1} = (R_{AB}[AB], \pi_{R_{AB}[AB]}^{\mathbf{E}_1})$ be the view of \mathbf{E}_1 which projects $R[ABC]$ onto $R_{AB}[AB]$ and which drops the relation S entirely. Consider M_{00} to be the initial instance of schema \mathbf{E}_1 ; its image instance in the view is then $N_{00} = \{R_{AB}(a_0, b_0), R_{AB}(a_1, b_1)\}$. Now, suppose that the view update $\text{Insert}\langle R_{AB}(a_2, b_2) \rangle$ is requested, so that $N_{01} = N_{00} \cup \{R_{AB}(a_2, b_2)\}$ is the desired new view instance, and consider $M_{01} = M_{00} \cup \{R(a_2, b_2, c_2), S(c_2, d_2)\}$ as a proposed reflection to the main schema \mathbf{E}_1 . Relative to its entire set $K_{01} = \{a_0, a_1, a_2, b_0, b_1, b_2, c_0, c_1, c_2, d_0, d_1, d_2\}$ of constant symbols, a cover for $\text{Info}\langle M_{01}, \text{WFS}(\mathbf{E}_1, \exists\wedge+, K_{01}) \rangle$ is just M_{01} itself. Similarly, for $M_{02} = M_{00} \cup \{R(a_2, b_2, c_3), S(c_3, d_3)\}$ with $K_{02} = \{a_0, a_1, a_2, b_0, b_1, b_2, c_0, c_1, c_3, d_0, d_1, d_3\}$ a cover for $\text{Info}\langle M_{02}, \text{WFS}(\mathbf{E}_1, \exists\wedge+, K_{02}) \rangle$ is just M_{02} itself. However, relative to K_{00} , which consists of the constant symbols found in elements of M_{00} , $\text{Info}\langle M_{01}, \text{WFS}(\mathbf{E}_1, \exists\wedge+, K_{00}) \rangle = \text{Info}\langle M_{02}, \text{WFS}(\mathbf{E}_1, \exists\wedge+, K_{00}) \rangle = M_{00} \cup \{(\exists x)(\exists y)(R(a_2, b_2, x) \wedge S(x, y))\}$. Denote this set of sentences by I_1 . This recaptures formally that the proposed updates M_{01} and M_{02} are identical up to a renaming of the new constants. The utility of information measure is that it provides a means to recapture this idea formally.

Now, consider the alternative solution $M_{03} = M_{00} \cup \{R(a_2, b_2, c_3), S(c_3, d_1)\}$ to this view-update problem. A cover for $\text{Info}\langle M_{03}, \text{WFS}(\mathbf{E}_1, \exists\wedge+, K_{00}) \rangle$ is $I_3 = M_{00} \cup \{(\exists x)(R(a_2, b_2, x) \wedge S(x, d_1))\}$, which is strictly stronger than I_1 , *i.e.*, $I_3 \models I_1$, since $(\exists x)(R(a_2, b_2, x) \wedge S(x, d_1)) \models (\exists x)(\exists y)(R(a_2, b_2, x) \wedge S(x, y))$, but not conversely. Thus, relative to the information measure defined by K_{00} , M_{03} adds more information to M_{00} than does M_{01} or M_{02} . Similarly, $M_{04} = M_{00} \cup \{R(a_2, b_2, c_0)\}$ adds more information than does M_{01} or M_{02} , since a cover for its information content is just M_{04} itself, which is stronger than I_1 , since $R(a_2, b_2, c_0) \wedge S(c_0, d_0) \models (\exists x)(\exists y)(R(a_2, b_2, x) \wedge S(x, y))$, but not conversely.

The first and primary measure of quality of a reflected update is the change of information content which it induces. Under this measure, M_{01} and M_{02} are equivalent, and both are superior to either of M_{03} or M_{04} . However, this measure alone is not quite adequate. Rather, there is an additional measure of quality which must be taken into account. To illustrate, consider the proposed solution $M_{05} = M_{01} \cup M_{02} = M_{00} \cup \{R(a_2, b_2, c_2), R(a_2, b_2, c_3), S(c_2, d_2), S(c_3, d_3)\}$ to this update problem. It has the same information content, I_1 , relative to K_{00} , as do M_{01} and M_{02} . The information measure cannot distinguish the insertion of two new tuples with completely new constants from the insertion of just one. However, it is clear that M_{05} should be considered inferior to both M_{01} and M_{02} as a solution to the given update problem, since it is a proper superset of each. Therefore, a second criterion of quality is invoked; namely that no solution whose set of changes is a proper superset of those of another can be considered to be superior. In the terminology introduced earlier in this section, the update must be minimal under subset ranking, and not just under count ranking. For example, consider again the proposed solution M_{04} . From a strict counting point of view, M_{04} involves fewer changes than do M_{01} or M_{02} . However, neither M_{01} nor M_{02} is a superset of M_{04} . Thus, the superiority of M_{01}

and M_{02} is not contradicted. In other words, only solutions which are *tuple minimal*, in the sense that no proper subset of the changes is also an admissible solution, are permitted.

The main modelling premise of this paper is that the quality of a view update can be measured by the amount of change in information content which it induces, and so an optimal reflection of a view update request is one which is both tuple minimal and which induces the least amount of change of information content. Under this premise, both M_{01} and M_{02} are superior to each of M_{03} , M_{04} , and M_{05} . Furthermore, since M_{01} and M_{02} induce the same change in information content, they are equivalent. In Section 3, it is established that, under suitable conditions, all such optimal solutions are equivalent, up to a renaming of the constant symbols. In Section 4, it is established, again under suitable conditions, that for insertions, a minimal solution (in terms of change of information content) must be optimal. These conditions include in particular schemata constrained by a very wide class of dependencies called generalized Horn dependencies.

In summary, there are two conditions which must be met for optimality of a proposed update reflection u . First, it must be *tuple minimal*, in that there can be no other solution whose set of changes is a proper subset of those of u . Second, it must induce a *least change of information* relative to a specific set of sentences. This approach applies also to deletions and updates which involve both insertion and deletion, and this generality is incorporated into the formalism which is presented.

1.4 Further issues Despite the connection to database repair just presented, the primary focus of this paper is not to present yet another measure for such repairs, but rather to present a unified approach to the support of updates on traditional relational schemata which pays particular attention to the requirements of invariance of admissibility, canonicity of reflections, and reflection of monotonicity. Such an approach requires that careful choices be made regarding the class of schemata and views which are supported. Foremost, schemata or view mappings which allow disjunction preclude canonicity of reflections in most cases. For example, let \mathbf{E}_2 have the three unary relation symbols $R[A]$, $S[A]$, and $T[A]$, subject to the constraint $(\forall x)(T(x) \Leftrightarrow (R(x) \vee S(x)))$, and let Ω_T be the view which contains only $T[A]$. It is easy to see that canonicity of reflections can never be satisfied. Consider the database $M_{\text{old}} = \{R(a_0), S(a_0), T(a_0)\}$, with the insertion request $\text{Insert}\langle T(a_1) \rangle$. A minimal solution would insert either $\{R(a_1), T(a_1)\}$ or else $\{S(a_1), T(a_1)\}$; to insert $\{R(a_1), S(a_1), T(a_1)\}$ would not be minimal. For this reason, attention in this work is restricted to relational schemata which are restricted by Horn constraints, such as the XEIDs of Fagin [Fag82].

Unfortunately, even within contexts which involve at most functional dependencies (FDs) and projections, two sorts of problems may occur for insertions. First of all, an attempted reflection to the main schema of a view update may introduce new tuples in the main schema, called *orphan tuples*, whose images are visible in the view. This phenomenon is illustrated via a concrete example in 4.8. Secondly, an insertion to a given view instance may be possible for some instances of the main schema which map to it, but not others. This phenomenon is illustrated via a concrete example in 4.12. It is shown that these anomalies may be remedied by requiring that the view be *strongly monotonic* — that is, that every deletion to the view may be reflected as a deletion to the main schema, and every insertion to the view may be reflected as an insertion to the main schema. Simple conditions which guarantee strong monotonicity for projections of relations constrained by FDs and inclusion dependencies are developed.

At first thought, it might appear that the management of deletions would be simpler than that of insertions, since generic values need not be created. However, there is a quite different type of complication which arises. Specifically, Horn-style dependencies of the form $A_1 \wedge A_2 \wedge \dots \wedge A_n \Rightarrow B$ are disjunction free with respect to insertions, but not deletions. Roughly, to delete B minimally, it is necessary to delete one of the A_i , but generally not all. One might therefore be lead to propose *weak optimal realizations*, as illustrated in the example of 6.4 and defined formally in 6.5, in which all of the contributing A_i 's are deleted. Unfortunately, as illustrated in 6.9, it is not even possible in general to delete them all and obtain a correct solution. Therefore, attention is focused upon schemata whose tuple-generating dependencies are of the form $A \Rightarrow B$, with just one assertion in the head of the rule. Fortunately, even with such a restriction, many common situations, particularly schemata constrained by functional and inclusion dependencies, and views defined by projection, are supported.

The most difficult cases surround updates which involve both insertions and deletions. In general, the information-based approach forwarded here does not provide optimal solutions to such requests, and so that topic must remain a subject for future work.

This article is based upon [Heg08a], but has been completely reworked to address some shortcomings in that preliminary version.

2 The Relational Model

The results of this paper are formulated within the relational model, and familiarity with its standard notions, as presented in references such as [PDGV89] and [AHV95], is assumed. Nevertheless, there are aspects which must be formulated with particular care. Most important are the need to take all relational schemata over the same domain, with the same constant symbols, and the need to express databases themselves as sets of ground atoms. For this reason, the central ideas which are special to this formulation are presented in this section.

2.1 Relational contexts and constant interpretations A relational context contains the logical information which is shared amongst the schemata and database mappings. Formally, a relational context \mathcal{D} consists of a finite nonempty set $\mathcal{A}_{\mathcal{D}}$ of attribute names, a countable set $\text{Vars}(\mathcal{D})$ of variables, and for each $A \in \mathcal{A}_{\mathcal{D}}$, a countable set $\text{Const}_{\mathcal{D}}(A)$ of constant symbols, with $\text{Const}(\mathcal{D}) = \bigcup \{\text{Const}_{\mathcal{D}}(A) \mid A \in \mathcal{A}_{\mathcal{D}}\}$. The variables in $\text{Vars}(\mathcal{D})$ are further partitioned into two disjoint sets; a countable set $\text{GenVars}(\mathcal{D}) = \{x_0, x_1, x_2, \dots\}$ of *general variables*, and special $\mathcal{A}_{\mathcal{D}}$ -indexed set $\text{AttrVars}(\mathcal{D}) = \{x_A \mid A \in \mathcal{A}_{\mathcal{D}}\}$ of *attribute variables*. The latter is used in the definition of interpretation mappings; see 2.6 for details. Lowercase letters at the end of the alphabet, such as v, w, x, y , and z , as well as subscripted instances using these names, will also be used as general variables.

Databases are represented as ground atoms, as elaborated in 2.2 below. Therefore, it is necessary that each domain element, in the sense of a logical structure for a first-order language, [Mon76, Def. 11.1], be bound to a unique constant symbol. Formally, a *constant interpretation* for the relational context \mathcal{D} is a pair $\mathcal{J} = (\text{Dom}_{\mathcal{J}}, \text{IntFn}_{\mathcal{J}})$ in which $\text{Dom}_{\mathcal{J}}$ is a countably infinite set, called the *domain* of \mathcal{J} , and $\text{IntFn}_{\mathcal{J}} : \text{Const}(\mathcal{D}) \rightarrow \text{Dom}_{\mathcal{J}}$ is a bijective function, called the *interpretation function* of \mathcal{J} . This effectively stipulates the following two well-known conditions [GN87, p. 120]:

Domain closure: $(\forall x \in \text{Vars}(\mathcal{D}))(\bigvee_{a \in \text{Const}(\mathcal{D})} x = a)$ (DCA(\mathcal{D}))

Unique naming: $(\neg(a = b))$ for distinct $a, b \in \text{Const}(\mathcal{D})$ (UNA(\mathcal{D}))

Since there are countably many constant symbols, the domain-closure axiom is not a finite disjunction. This is not a problem however, since it is never used in a context in which a first-order constraint is necessary. Rather, the assignment of a constant to each variable is taken to be part of the context in which this work is carried out.

As a notational convention, from this point on, unless stated otherwise, fix a relational context \mathcal{D} and a constant interpretation $\mathcal{J} = (\text{Dom}_{\mathcal{J}}, \text{IntFn}_{\mathcal{J}})$ for it.

2.2 Tuples and databases An *unconstrained relational schema* over $(\mathcal{D}, \mathcal{J})$ is a pair $\mathbf{D} = (\text{Rels}(\mathbf{D}), \text{Ar}_{\mathbf{D}})$ in which $\text{Rels}(\mathbf{D})$ is a finite set of relational symbols and $\text{Ar}_{\mathbf{D}} : \text{Rels}(\mathbf{D}) \rightarrow \mathbf{2}^{\mathcal{A}_{\mathcal{D}}}$ a function which assigns an *arity*, a set of distinct attributes from $\mathcal{A}_{\mathcal{D}}$, to each $R \in \text{Rels}(\mathbf{D})$.

An *R-atom* is a function $t : \text{Ar}_{\mathbf{D}}(R) \rightarrow \text{Const}(\mathcal{D}) \cup \text{Vars}(\mathcal{D})$ with the property that $t[A] \in \text{Const}_{\mathcal{D}}(A) \cup \text{GenVars}(\mathcal{D}) \cup \{x_A\}$; in other words, all terms, constant and variable, are of the appropriate type. A *ground R-atom* has the additional property that it contains no variables; *i.e.*, $t[A] \in \text{Const}_{\mathcal{D}}(A)$. The set of all *R*-atoms (resp. ground *R*-atoms) is denoted $\text{Atoms}(R, \mathbf{D})$ (resp. $\text{GrAtoms}(R, \mathbf{D})$).

A *D-atom* is an *R*-atom for some $R \in \text{Rels}(\mathbf{D})$; the set of all such atoms is denoted $\text{Atoms}(\mathbf{D})$. A *ground atom* is defined to be a ground *R*-atom for some $R \in \text{Rels}(\mathbf{D})$, with the set of all such atoms denoted $\text{GrAtoms}(\mathbf{D})$. An *atom database for D* is a finite subset of $\text{GrAtoms}(\mathbf{D})$, with the set of all atom databases for \mathbf{D} denoted $\text{DB}(\mathbf{D})$. In this work, ground atoms are also called *tuples*.

It is convenient to be able to recover the associated relation name from a tuple, and so *tagging* is employed, in which tuples are marked with the relation name. Formally, this is accomplished by introducing a new attribute $\text{RName} \notin \mathcal{A}_{\mathcal{D}}$, and then regarding an *R*-atom not as a function t just on $\text{Ar}_{\mathbf{D}}(R)$, but rather as one on $\{\text{RName}\} \cup \text{Ar}_{\mathbf{D}}(R)$ with the property that $t[\text{RName}] = R$. Tagging of *R*-atoms will be used from this point on throughout the paper. Nevertheless, in writing such atoms, the more conventional notation $R(a_1, a_2, \dots, a_n)$ will be used in lieu of the technically more correct $(R, a_1, a_2, \dots, a_n)$, although tags will be used in formal constructions. To be completely pedantic, this entails introducing a new attribute name $\text{RNames} \in \mathcal{A}_{\mathcal{D}}$ with $\text{Const}_{\mathcal{D}}(\text{RNames}) = \text{Rels}(\mathbf{D})$, and these constant values used only for the attribute RNames . Furthermore, in any atom, the value for the RNames attribute must be a constant, never a variable. Since this is a logically inessential tactic whose full formal treatment is tedious but routine, the details will not be elaborated further.

There is a third type of atom which will be of use in defining constraints, the equality atom. Formally, an *equality atom* is of one of the forms $(x_i = x_j)$, $(x_i = a_j)$, or $(a_i = a_j)$, for $x_i, x_j \in \text{GenVars}(\mathcal{D})$ and $a_i, a_j \in \text{Const}(\mathcal{D})$. The set of all equality \mathcal{D} -atoms is denoted $\text{EqAtoms}(\mathcal{D})$. Equality atoms which equate two constants; *e.g.*, $(a_i = a_j)$ are called *ground equality atoms*; note that the truth value of such atoms is predetermined by the unique naming assumption. All other equality atoms; *e.g.*, those of the forms $(x_i = x_j)$ or $(x_i = a_j)$, are called *variable equality atoms*. The set of all variable equality atoms is denoted $\text{VarEqAtoms}(\mathcal{D})$. Note that the definitions of equality atoms depend only upon the relational context \mathcal{D} , and not upon the specific schema \mathbf{D} .

2.3 Formulas and constraint classes The first-order language associated with the relational schema \mathbf{D} is defined in the natural way; however, it is useful to introduce some notation which identifies particular sets of formulas. Define $\text{WFF}(\mathbf{D})$ to be the set of all well-formed first-order formulas with equality in the language whose set of relational symbols is $\text{Rels}(\mathbf{D})$, whose set of constant symbols is $\text{Const}(\mathcal{D})$, and which contains no non-nullary function symbols. The variables are those of \mathcal{D} ; these formulas are typed to the extent that for $A \in \mathcal{A}_{\mathcal{D}}$, a term in a position of type A must lie in $\text{GenVars}(\mathcal{D}) \cup \{x_A\} \cup \text{Const}_{\mathcal{D}}(A)$.

A *constraint class* \mathcal{C} identifies a subset of $\text{WFF}(\mathbf{D})$, denoted $\text{WFF}(\mathbf{D}, \mathcal{C})$. For this paper, the two most important constraint classes are $\exists\wedge+$ and GrAtoms , defined as follows.

- $\text{WFF}(\mathbf{D}, \exists\wedge+)$ is the subset of $\text{WFF}(\mathbf{D})$ in which only existential quantification is allowed, and the only logical connective which is allowed is conjunction (\wedge). These formulas define the so-called *conjunctive queries* [CGT90, Sec. 4.2]. It is not necessary to allow the equality predicate in such formulas, since equality can always be expressed by simply using the same name for the two atoms which are equated.
- $\text{WFF}(\mathbf{D}, \text{GrAtoms})$ is just $\text{GrAtoms}(\mathbf{D})$.

$\text{WFF}(\mathbf{D}, \mathcal{C})$ may be trimmed further by limiting the constant symbols which may occur in it. Specifically, if $S \subseteq \text{Const}(\mathcal{D})$, then $\text{WFF}(\mathbf{D}, \mathcal{C}, S)$ denotes the formulas in $\text{WFF}(\mathbf{D})$ which involve only constant symbols from S .

Each of these classes may be limited to sentences; *i.e.*, formulas without free variables.

$\text{WFS}(\mathbf{D})$

(resp. $\text{WFS}(\mathbf{D}, \mathcal{C})$, resp. $\text{WFS}(\mathbf{D}, \mathcal{C}, S)$) denotes the subset of $\text{WFF}(\mathbf{D})$ (resp. $\text{WFF}(\mathbf{D}, \mathcal{C})$, resp. $\text{WFF}(\mathbf{D}, \mathcal{C}, S)$) consisting of sentences.

Let $\Phi \subseteq \Psi \subseteq \text{WFS}(\mathbf{D})$. The *closure* of Φ in Ψ , denoted $\text{Closure}\langle\Phi, \Psi\rangle$, is $\{\varphi \in \Psi \mid \Phi \models \varphi\}$. A *cover* for Φ relative to Ψ is a subset $\Phi' \subseteq \Phi$ with $\text{Closure}\langle\Phi', \Psi\rangle = \text{Closure}\langle\Phi, \Psi\rangle$. A *minimal cover* Ψ' has the property that none of its proper subsets is itself a cover.

Finally, the symbol \perp will be used to denote the sentence which is always false.

2.4 Atomic models Even though databases are represented as sets of ground atoms, and not as interpretations in the usual logical sense, it is still essential to have an appropriate notion of model for a given sentence. This is relatively straightforward; a model for a sentence φ is a database which is consistent with both φ and the unique-naming axioms. There is one complication, however. In representing a database as a set of \mathbf{D} -atoms, the closed-world assumption is implicit. On the other hand, to express what it means for such a representation to satisfy an arbitrary sentence in $\text{WFS}(\mathbf{D})$, it is necessary to state explicitly which atoms are not true as well. Formally, for $M \in \text{DB}(\mathbf{D})$, define the *diagram* of M to be $\text{Diagram}_{\mathbf{D}}(M) = M \cup \{\neg t \mid t \in \text{GrAtoms}(\mathbf{D}) \setminus M\}$. Now, say that $M \in \text{DB}(\mathbf{D})$ is an *atomic \mathcal{J} -model* of $\varphi \in \text{WFS}(\mathbf{D})$ if $\text{Diagram}_{\mathbf{D}}(M) \cup \{\varphi\} \cup \text{UNA}(\mathcal{D})$ is consistent. $\text{AtMod}_{\mathcal{J}}(\varphi)$ denotes the set of all atomic \mathcal{J} -models of φ , with $\text{AtMod}_{\mathcal{J}}(\Phi) = \bigcap\{\text{AtMod}_{\mathcal{J}}(\varphi) \mid \varphi \in \Phi\}$ for $\Phi \subseteq \text{WFS}(\mathbf{D})$. Since only atomic \mathcal{J} -models will be considered in this paper, the simple term *model* will be used as a synonym for *atomic \mathcal{J} -model*.

2.5 Schemata with constraints and constrained databases To obtain full relational schemata, constraints are added to the unconstrained schemata of 2.2. Formally, a *relational*

schema over $(\mathcal{D}, \mathcal{J})$ is a triple $\mathbf{D} = (\text{Rels}(\mathbf{D}), \text{Ar}_{\mathbf{D}}, \text{Constr}(\mathbf{D}))$ in which $(\text{Rels}(\mathbf{D}), \text{Ar}_{\mathbf{D}})$ is an unconstrained relational schema over $(\mathcal{D}, \mathcal{J})$ and $\text{Constr}(\mathbf{D}) \subseteq \text{WFS}(\mathbf{D})$ is the set of *dependencies* or *constraints* of \mathbf{D} .

Define the *legal* (or *constrained*) *databases* $\text{LDB}(\mathbf{D})$ of \mathbf{D} to be $\text{AtMod}_{\mathcal{J}}(\text{Constr}(\mathbf{D}))$.

Although $\text{Constr}(\mathbf{D})$ is allowed to be an infinite set, it will always be assumed that $\text{Constr}(\mathbf{D})$ is *constant finite*; that is, that all of the sentences in $\text{Constr}(\mathbf{D})$ together contain only a finite number of distinct constant symbols.

2.6 Database morphisms and views Let \mathbf{D}_1 and \mathbf{D}_2 be relational schemata over $(\mathcal{D}, \mathcal{J})$. There are two fundamental ways to represent a database morphism $f : \mathbf{D}_1 \rightarrow \mathbf{D}_2$ in the relational context. On the one hand, such a morphism may be represented as a function $f : \text{DB}(\mathbf{D}_1) \rightarrow \text{DB}(\mathbf{D}_2)$, using expressions from the relational algebra. On the other hand, by providing an interpretation formula $f^R \in \text{WFF}(\mathbf{D}_1)$ for each $R \in \text{Rels}(\mathbf{D}_2)$, the morphism may be represented using the relational calculus [JAK82]. The equivalence of these two representations is one of the classical results of relational database theory [PDGV89, Sec. 2.4-2.6]. The interpretation formulation is taken as the base definition for views in this work. Formally, given $R \in \text{Rels}(\mathbf{D}_2)$, an *interpretation* for R into \mathbf{D}_1 is a $\varphi \in \text{WFF}(\mathbf{D}_1)$ in which precisely the variables $\{x_A \mid A \in \text{Ar}_{\mathbf{D}}(R)\}$ are free, and in which x_A is used to mark the position in the formula which is bound to attribute A . Since each position in the view relation is associated with a distinct attribute, one variable per attribute suffices. The set of all interpretations of R into \mathbf{D}_1 is denoted $\text{Interp}(R, \mathbf{D}_1)$. A *syntactic morphism* $f : \mathbf{D}_1 \rightarrow \mathbf{D}_2$ is a family $f = \{f^R \mid R \in \text{Rels}(\mathbf{D}_2) \text{ and } f^R \in \text{Interp}(R, \mathbf{D}_1)\}$. The morphism f is said to be *of class* $\exists\wedge+$ if $f^R \in \text{WFF}(\mathbf{D}, \exists\wedge+)$ for each $R \in \text{Rels}(\mathbf{D}_2)$.

Let $t \in \text{Atoms}(R, \mathbf{D}_2)$. The *substitution* of t into f , denoted $\text{Substf}\langle f, t \rangle$, is the formula in $\text{WFF}(\mathbf{D}_1)$ obtained by substituting, into f^R , $t[A]$ for x_A , for each $A \in \text{Ar}_{\mathbf{D}}(R)$. Note that if t is a ground atom, then $\text{Substf}\langle f, t \rangle \in \text{WFS}(\mathbf{D}_1)$.

For $M \in \text{DB}(\mathbf{D}_1)$, define $f(M) = \{t \in \text{GrAtoms}(\mathbf{D}_2) \mid M \in \text{AtMod}_{\mathcal{J}}(\text{Substf}\langle f, t \rangle)\}$. f is called an *LDB-morphism* if it maps legal databases to legal databases; formally, an LDB-morphism has the property that $f(M) \in \text{LDB}(\mathbf{D}_2)$ for each $M \in \text{LDB}(\mathbf{D}_1)$. When no qualification is given, *database morphism* will always mean LDB-morphism.

Let \mathbf{D} be a relational schema over $(\mathcal{D}, \mathcal{J})$. A (*relational*) *view* of \mathbf{D} is a pair $\Gamma = (\mathbf{V}, \gamma)$ in which \mathbf{V} is a relational schema over $(\mathcal{D}, \mathcal{J})$ and $\gamma : \mathbf{D} \rightarrow \mathbf{V}$ is an LDB-morphism which is furthermore *LDB-surjective* in the sense that for every $N \in \text{LDB}(\mathbf{V})$, there is an $M \in \text{LDB}(\mathbf{D})$ with $\gamma(M) = N$. Surjectivity is required because the instance of the view must always be determined by the instance of the main schema \mathbf{D} . The view $\Gamma = (\mathbf{V}, \gamma)$ is said to be *of class* $\exists\wedge+$ precisely in the case that γ has that property.

In order to illustrate these ideas, a simple example is in order. Consider again the schema \mathbf{E}_0 and the view $\Pi_{AB}^{\mathbf{E}_0}$ of 1.1. The view mapping $\pi_{AB}^{\mathbf{E}_0}$ is expressed as an interpretation via the formula $(\pi_{AB}^{\mathbf{E}_0})^{R_{AB}} = (\exists z)(R(x_A, x_B, z))$. Note in particular how x_A and x_B are used to mark the appropriate attributes. For $t = R_{AB}(a_0, b_0)$, $\text{Substf}\langle \pi_{AB}^{\mathbf{E}_0}, t \rangle = (\exists z)(R(a_0, b_0, z))$, while for $t = R_{AB}(x_0, x_1)$, $\text{Substf}\langle \pi_{AB}^{\mathbf{E}_0}, t \rangle = (\exists z)(R(x_0, x_1, z))$.

Occasionally, it will be useful to separate the quantifiers from the rest of the formula of an interpretation γ^R . To this end, define $\overline{\gamma}^R$ to be that which is left when the quantifier prefix is removed from γ^R . For example, in the above, $(\overline{\pi_{AB}^{\mathbf{E}_0}})^{R_{AB}} = (R(x_A, x_B, z))$.

2.7 Notation — extracting constant symbols and variables For X an entity (for example, an atom, a formula, a database, etc.), or a set of entities, $\text{ConstSym}(X)$ denotes the set of all $a \in \text{Const}(\mathcal{D})$ which occur in X . Furthermore, for \mathbf{D} a database schema, $\text{ConstSym}_{\mathbf{D}}(X) = \text{ConstSym}(X \cup \text{Constr}(\mathbf{D}))$, and for $\Gamma = (\mathbf{V}, \gamma)$ a view of \mathbf{D} , $\text{ConstSym}_{\Gamma}(X) = \text{ConstSym}(X \cup \text{Constr}(\mathbf{D}) \cup \text{Constr}(\mathbf{V})) \cup \text{ConstSym}(\gamma)$, where $\text{ConstSym}(\gamma)$ is the set of all constant symbols which occur in the defining interpretation formulas associated with γ .

Similarly, $\text{Vars}(X)$ denotes the set of all variables which occur in X . This will not be formalized further, but the meaning should always be unambiguous.

2.8 Notation for inclusion dependencies It is assumed that the reader is familiar with the relational model and the standard dependencies which have been studied in that context. Here only some notation and terminology is clarified. First, $R[X] \sqsubseteq S[Y]$ (note the squared subset symbol) will be used to denote the inclusion dependency (IND) which states that the projection onto attributes X of relation R is a subset of the projection onto attributes Y of relation S . Second, a *unary inclusion dependency*, abbreviated *UIND*, is one in which each of X and Y consist of a single attribute.

3 Information and Canonical Models

The theory of support for view updates which is forwarded in this paper is based upon a duality between a set of sentences defining information content and canonical models for such information. In this section, that duality is developed in detail.

3.1 Notational convention Throughout the rest of this paper, unless stated specifically to the contrary, take \mathbf{D} to be a relational schema over $(\mathcal{D}, \mathcal{J})$. The notation Υ will be used as an abbreviation for $\text{WFS}(\mathbf{D}, \exists\wedge+)$, and Υ_K will be used as an abbreviation for $\text{WFS}(\mathbf{D}, \exists\wedge+, K)$. Furthermore, in the context of a set of the form $\text{WFS}(\mathbf{D}, \exists\wedge+, K)$, if no further information is given, K will be taken to be an arbitrary subset of $\text{Const}(\mathcal{D})$.

3.2 Information content and Φ -equivalence Let $\Phi \subseteq \text{WFS}(\mathbf{D})$ and let $M \in \text{DB}(\mathbf{D})$. The *information content* of M relative to Φ is the set of all sentences in Φ which are true for M . More precisely, $\text{Info}\langle M, \Phi \rangle = \{\varphi \in \Phi \mid M \in \text{AtMod}_{\mathcal{J}}(\varphi)\}$. For $\varphi \in \text{WFS}(\mathbf{D})$, $\text{Info}\langle M, \varphi \rangle$ denotes $\text{Info}\langle M, \{\varphi\} \rangle$. M_1 and M_2 are Φ -*equivalent* if they have the same information content relative to Φ ; *i.e.*, $\text{Info}\langle M_1, \Phi \rangle = \text{Info}\langle M_2, \Phi \rangle$.

The semantics of conventional databases are based upon the closed-world assumption — all assertions which cannot be established to be true are taken to be false. Thus, intuitively, information content should be monotone; that is, for any $M_1, M_2 \in \text{DB}(\mathbf{D})$ if $M_1 \subseteq M_2$, then $\text{Info}\langle M_1, \Phi \rangle \subseteq \text{Info}\langle M_2, \Phi \rangle$. However, this is manifestly false for most choices of Φ . Indeed, if φ holds in M_2 but not M_1 , then $\neg\varphi$ holds in M_1 but not M_2 . Thus, if there is some $\varphi \in \Phi$ for which $\neg\varphi \in \Phi$ as well, Φ cannot be information monotone.

Formally, it is best to begin by defining monotonicity for individual sentences. To be precise, the sentence $\varphi \in \text{WFS}(\mathbf{D})$ is *information monotone* if for any $M_1, M_2 \in \text{DB}(\mathbf{D})$ if $M_1 \subseteq M_2$, then $\text{Info}\langle M_1, \varphi \rangle \subseteq \text{Info}\langle M_2, \varphi \rangle$. The set $\Phi \subseteq \text{WFS}(\mathbf{D})$ is then said to be *information monotone* if each $\varphi \in \Phi$ has this property. It is easy to see that any $\varphi \in \text{WFS}(\mathbf{D})$ which

does not involve negation, when expressed entirely in terms of the connectives \wedge , \vee , and \neg , is information monotone. Connectives such as \Rightarrow involve negation implicitly, and so sentences which involve implication need not be (and usually are not) information monotone.

In this work, there are two families of information-monotone sentences which are of central importance. The first is $\text{WFS}(\mathbf{D}, \text{GrAtoms})$. It is easy to see that $\text{Info}\langle M, \text{WFS}(\mathbf{D}, \text{GrAtoms}) \rangle = M$ for any $M \in \text{DB}(\mathbf{D})$, so that the information content of a database relative to $\text{WFS}(\mathbf{D}, \text{GrAtoms})$ is just that database itself. Although trivial in its characterization, this case is nonetheless important. The second important family of information monotone sentences is $\text{WFS}(\mathbf{D}, \exists\wedge+, K)$ for a given $K \subseteq \text{Const}(\mathcal{D})$, and is far less trivial in its characterization.

3.3 Tuple-minimal models Let $\Phi \subseteq \text{WFS}(\mathbf{D})$. and let $M \in \text{AtMod}_J(\Phi)$. M is a *tuple-minimal model* of Φ if for any $M' \in \text{AtMod}_J(\Phi)$ with $M' \subseteq M$, it must be that $M' = M$. The set of all tuple-minimal models of Φ is denoted $\text{MinAtMod}_J(\Phi)$. For $\varphi \in \text{WFS}(\mathbf{D})$, $\text{MinAtMod}_J(\varphi)$ is shorthand for $\text{MinAtMod}_J(\{\varphi\})$.

For $\Phi \subseteq \text{WFS}(\mathbf{D})$ and $\varphi \in \text{WFS}(\mathbf{D})$, say that Φ *minimally entails* φ , written $\Phi \models_{\text{min}} \varphi$, if $\text{MinAtMod}_J(\Phi) \subseteq \text{AtMod}_J(\varphi)$. In other words, Φ minimally entails φ if every tuple-minimal model of Φ is also a model (not necessarily minimal) of φ .

3.4 Fully Reduced $\exists\wedge+$ -families The concept of a minimal cover for a set Φ of sentences is well known and has already been recalled in 2.3. In the context of sentences in Υ , there is a stronger notion which is central to the development of the ideas presented here. To motivate this idea, let $\Xi = \{R(a_1, a_2), R(a_2, a_3), (\exists x_1)(\exists x_2)(\exists x_3)(R(x_1, x_2) \wedge R(x_2, a_3) \wedge R(a_3, x_3))\}$. It is easy to see that Ξ is a minimal cover of itself, in that none of its proper subsets is equivalent to it. However, it is also clear that $\xi_0 = (\exists x_1)(\exists x_2)(\exists x_3)(R(x_1, x_2) \wedge R(x_2, a_3) \wedge R(a_3, x_3))$ may be replaced with $(\exists x_3)(R(a_3, x_3))$ while retaining logical equivalence. In other words, conjuncts may be removed from one of the sentences while preserving information content.

To formalize this notion, let $\varphi = (\exists x_1) \dots (\exists x_m)(A_1 \wedge \dots \wedge A_n) \in \Upsilon$ have at least two conjuncts, and for any i , $1 \leq i \leq n$, define $\text{Reduction}\langle \varphi, A_i \rangle$ to be the sentence obtained by removing A_i as a conjunct from φ , and removing any quantifier term which is no longer used as well. For example, $\text{Reduction}\langle \xi_0, R(x_1, x_2) \rangle = (\exists x_2)(\exists x_3)(R(x_2, a_3) \wedge R(a_3, x_3))$ and $\text{Reduction}\langle \text{Reduction}\langle \xi_0, R(x_1, x_2) \rangle, R(x_2, a_3) \rangle = (\exists x_3)(R(a_3, x_3))$. Call $\Phi \subseteq \Upsilon$ *conjunct reduced* if for no $\varphi \in \Phi$ with at least two conjuncts is there a conjunct A_i of φ with $(\Phi \setminus \varphi) \cup \{\text{Reduction}\langle \varphi, A_i \rangle\}$ logically equivalent to Φ . Call Φ *fully reduced* if it is both conjunct reduced and a minimal cover of itself. In the above example, $\{R(a_1, a_2), R(a_2, a_3), (\exists x_3)(R(a_3, x_3))\}$ is fully reduced.

The goal is to establish that by substituting distinct constants for each variable in a finite, fully reduced family of sentences, a canonical model of those sentences is obtained. Thus, in the above example, $\{R(a_1, a_2), R(a_2, a_3), (R(a_3, b_1))\}$ would be such a model, with b_1 a “generic” constant. To render all of this formal, some additional notions are necessary.

3.5 Armstrong models in an information-monotone context Let $\Psi \subseteq \text{WFS}(\mathbf{D})$ and let $\Phi \subseteq \Psi$. Informally, an Armstrong model for Φ relative to Ψ is a model of Φ which satisfies only those constraints of Ψ which are implied by Φ . More formally, an *Armstrong model* for Φ relative to Ψ is an $M \in \text{AtMod}_J(\Phi)$ with the property that for any $\psi \in \Psi$, if $M \in \text{AtMod}_J(\psi)$,

then $\text{AtMod}_{\mathcal{J}}(\Phi) \subseteq \text{AtMod}_{\mathcal{J}}(\psi)$. A *tuple-minimal Armstrong model* for Φ relative to Ψ is an Armstrong model with the property that no proper subset is an Armstrong model for Φ relative to Ψ . In general, a tuple-minimal Armstrong model M of Φ relative to Ψ need not be a tuple-minimal model of Φ , since there may be an $M' \subsetneq M$ which is a non-Armstrong model of Φ . However, if Ψ is information monotone, it is easy to see that this cannot happen, so every tuple-minimal Armstrong model must in fact be a minimal model. Armstrong models have been studied extensively for database dependencies; see, for example, [Fag82] and [FV83].

In the current context, for a given finite, fully reduced set $\Phi \subseteq \Upsilon$, a suitably constructed Armstrong model for Φ relative to Υ_K for a given K with $\text{ConstSym}(\Phi) \subseteq K$ will serve as a canonical representation for insertions with generic constants, as sketched in the introduction. To proceed further, a special representation is useful.

3.6 Representation of $\exists\wedge+$ -sentences as sets of D-atoms There is an alternative syntactic representation for formulas in Υ which will be used in that which follows. Specifically, for $\varphi \in \Upsilon$ define $\text{AtRep}(\varphi)$ to be the set of all atoms which occur as conjuncts in φ . For example, if $\varphi = (\exists x_1)(\exists x_2)(\exists x_3)(R(x_1, a) \wedge R(x_1, b) \wedge S(x_2, a) \wedge T(x_2, x_3))$ then $\text{AtRep}(\varphi) = \{R(x_1, a), R(x_1, b), S(x_2, a), T(x_2, x_3)\}$.

This representation is dual to that used in theorem-proving contexts in classical artificial intelligence [GN87, 4.1]. Here the variables are existentially quantified and the atoms are conjuncts of one another; in the AI setting the atoms are disjuncts of one another and the variables are universally quantified.

3.7 Substitutions Let $V = \{x_1, x_2, \dots, x_n\} \subseteq \text{GenVars}(\mathcal{D})$. A *substitution* for V (in \mathcal{D}) is a function $s : V \rightarrow \text{Const}(\mathcal{D}) \cup \text{GenVars}(\mathcal{D})$. If $s(x_i) = \tau_i$ for $i \in \{1, 2, \dots, n\}$, following (somewhat) standard notation this substitution is $\{\tau_1/x_1, \tau_2/x_2, \dots, \tau_n/x_n\}$ [CL73, Sec. 5.3] and will be used here, although the reader is cautioned that some authors write $\{x_1/\tau_1, x_2/\tau_2, \dots, x_n/\tau_n\}$ instead [GN87, 4.2].

Let $\varphi \in \Upsilon$ with $\text{Vars}(\varphi) \subseteq V$. Call s *correctly typed* for φ if for each $t \in \text{AtRep}(\varphi)$ and each $A \in \text{Ar}_{\mathbf{D}}(t[\text{RName}])$, if $t[A] \in \text{Vars}(\mathbf{D})$ then $s(t[A]) \in \text{Const}_{\mathcal{D}}(A) \cup \text{GenVars}(\mathcal{D})$. Define $\text{Subst}(\varphi, s)$ to be the set of atoms obtained by substituting $s(x_i)$ for x_i in $\text{AtRep}(\varphi)$. For example, with $s = \{a_1/x_1, a_2/x_2, a_3/x_3\}$ and $\text{AtRep}(\varphi) = \{R(x_1, a), R(x_1, b), S(x_2, a), T(x_2, x_3)\}$, $\text{Subst}(\varphi, s) = \{R(a_1, a), R(a_1, b), S(a_2, a), T(a_2, a_3)\}$.

If $s(x_i) \in \text{Const}(\mathcal{D})$ for each $x_i \in V$, s is called a *constant substitution*. In this case, $\text{Subst}(\varphi, s)$ is a set of ground atoms.

Now let $\Phi \subseteq \Upsilon$ be a finite set. A *constant substitution set* for Φ is a Φ -indexed set $S = \{s_{\varphi} \mid \varphi \in \Phi\}$ of substitutions, with s_{φ} a constant substitution for $\text{Vars}(\varphi)$. For K a finite set with $\text{ConstSym}(\Phi) \subseteq K \subseteq \text{Const}(\mathcal{D})$, S is *free for $\langle \Phi, K \rangle$* if each s_{φ} is correctly typed for φ , injective, $s_{\varphi}(x_i) \notin K$ for any $\varphi \in \Phi$ and $x_i \in V$, and, furthermore, for any distinct $\varphi_1, \varphi_2 \in \Phi$, $s_{\varphi_1}(\text{Vars}(\varphi_1)) \cap s_{\varphi_2}(\text{Vars}(\varphi_2)) = \emptyset$.

With S free for $\langle \Phi, K \rangle$, the *Armstrong model defined by $\langle \Phi, K, S \rangle$* is obtained by applying the substitution s_{φ} to φ for each $\varphi \in \Phi$. Formally, $\text{ArmMod}\langle \Phi, K, S \rangle = \bigcup \{\text{Subst}(\varphi, s_{\varphi}) \mid \varphi \in \Phi\}$. Of course, this terminology is a bit presumptuous, as it has not yet been established that $\text{ArmMod}\langle \Phi, K, S \rangle$ is in fact an Armstrong model of anything; this will be rectified in 3.9 below.

3.8 Constant endomorphisms Informally, an endomorphism on \mathcal{D} is a function which renames constants. More formally, an *endomorphism* on \mathcal{D} is a function $h : \text{Const}(\mathcal{D}) \rightarrow \text{Const}(\mathcal{D})$ which preserves attribute types, in the precise sense that for each $A \in \mathcal{A}_{\mathcal{D}}$ and each $a \in \text{Const}_{\mathcal{D}}(A)$, $h(a) \in \text{Const}_{\mathcal{D}}(A)$. If h is additionally a bijection, then it is called an *automorphism* of \mathcal{D} . For $K \subseteq \text{Const}(\mathcal{D})$, call h *K-invariant* if $h(a) = a$ for all $a \in K$.

Given a database schema \mathbf{D} , an endomorphism on \mathcal{D} induces a mapping from $\text{GrAtoms}(\mathbf{D})$ to itself given by sending $t \in \text{GrAtoms}(\mathbf{D})$ to the tuple t' with $t'[\text{RName}] = t[\text{RName}]$ and $t'[A] = t[h(A)]$ for all $A \in \text{Ar}_{\mathbf{t}[\text{RName}]}$. This mapping on atoms will also be represented by h , as will the induced mapping from $\text{DB}(\mathbf{D})$ to itself given by $M \mapsto \{h(t) \mid t \in M\}$.

The following theorem establishes that $\text{ArmMod}\langle\Phi, K, S\rangle$ is a weak sort of initial model for Φ , in the sense that for any other database M which satisfies Φ , there is an endomorphism $h : \text{ArmMod}\langle\Phi, K, S\rangle \rightarrow M$ which holds K constant. On the other hand, it is not an initial model for Φ in the traditional categorical sense [HS73, §7], since h need not be unique.

3.9 Theorem — Characterization of tuple-minimal Armstrong models *Let $\Phi \subseteq \Upsilon$ be finite and fully reduced, let K be a finite set with $\text{ConstSym}(\Phi) \subseteq K \subseteq \text{Const}(\mathcal{D})$, and let S be a constant substitution set which is free for Φ . Then the following hold.*

- (a) *For any $M \in \text{DB}(\mathbf{D}) \cap \text{AtMod}_{\mathcal{J}}(\Phi)$, there is a K -invariant endomorphism h on \mathcal{D} with $h(\text{ArmMod}\langle\Phi, K, S\rangle) \subseteq M$.*
- (b) *$\text{ArmMod}\langle\Phi, K, S\rangle$ is a tuple-minimal Armstrong model for Φ relative to Υ_K .*
- (c) *If $M \in \text{DB}(\mathbf{D})$ is any other tuple-minimal Armstrong model for Φ relative to Υ_K , then there is a $\text{ConstSym}(\Phi)$ -invariant automorphism h on \mathcal{D} with $h(\text{ArmMod}\langle\Phi, K, S\rangle) = M$.*

PROOF: To establish (a), let $M \in \text{Mod}_{\mathcal{J}}(\Phi)$, and for each $\varphi \in \Phi$, let M_{φ} be a minimal subset of M with $M_{\varphi} \in \text{Mod}_{\mathcal{J}}(\varphi)$. Let V_{φ} denote the set of variables of $s_{\varphi} \in S$. It is easy to see that there must be a constant substitution s'' with $\text{Vars}(s'') = V_{\varphi}$ and $\text{Subst}(\varphi, s'') = M_{\varphi}$. Indeed, there is trivially a constant substitution with $\text{Subst}(\varphi, s'') \subseteq M_{\varphi}$, but if the subset inclusion were proper, M_{φ} would not be tuple minimal.

Now define $h : s_{\varphi}(V_{\varphi}) \rightarrow s''(V_{\varphi})$ by $a \mapsto s''(s_{\varphi}^{-1}(a))$. Since s_{φ} is injective, h is well defined. Since $s_{\varphi_1}(\text{Vars}(\varphi_1)) \cap s_{\varphi_2}(\text{Vars}(\varphi_2)) = \emptyset$ for distinct $\varphi_1, \varphi_2 \in \Phi$, there are no conflicts in this definition of h . Finally, extend h to be the identity on all $a \in \text{Const}(\mathcal{D})$ which are not covered by the above definition. The result is a endomorphism on \mathcal{D} which satisfies $h(\text{ArmMod}\langle\Phi, K, S\rangle) \subseteq M$.

For (b), first observe that $\text{ArmMod}\langle\Phi, K, S\rangle$ is a model of Φ just by construction. It is furthermore easy to see that since Φ is fully reduced, it is tuple minimal. Indeed, if any tuple $t \in \text{ArmMod}\langle\Phi, K, S\rangle$ could be removed, then the corresponding conjunct could be removed from the $\varphi \in \Phi$ associated with t , contradicting the fact that Φ is fully reduced. To show that $\text{ArmMod}\langle\Phi, K, S\rangle$ is an Armstrong model, let $\psi \in \Upsilon$ with $\text{ArmMod}\langle\Phi, K, S\rangle \in \text{Mod}_{\mathcal{J}}(\psi)$, and let $M \in \text{Mod}_{\mathcal{J}}(\Phi)$. In view of (a), there is an endomorphism h on \mathcal{D} with $h(\text{ArmMod}\langle\Phi, K, S\rangle) \subseteq M$. In view of Lyndon's theorem [Mon76, Thm. 25.22], which states that satisfaction of sentences not involving negation is closed under endomorphic images, it follows that $M \in \text{Mod}_{\mathcal{J}}(\psi)$ also. Hence, $\Phi \models \psi$, and so $\text{ArmMod}\langle\Phi, K, S\rangle$ is an Armstrong model of Φ .

To show (c), let M be any other tuple-minimal Armstrong model for Φ relative to Υ_K . In the above construction for the proof of (a), the resulting h must be surjective (else M would

not be tuple minimal), and it must be injective (since there must also be an endomorphism in the opposite direction, and both $\text{ArmMod}\langle\Phi, K, S\rangle$ and M are finite, by assumption). Hence, h is an automorphism. \square

It is easy to see that the endomorphism h need not be unique. For example, if \mathbf{D} has the single unary relation symbol $R[A]$, and $\Phi = \{(\exists x)(R(x))\}$, then $M_1 = \{R(a)\}$ is a minimal Armstrong model, while $M_2\{R(b), R(c)\}$ is an Armstrong model which is not tuple minimal. There are two endomorphisms from M_1 to M_2 , $h_1 : a \mapsto b$ and $h_2 : a \mapsto c$.

In some ways, the construction given above is similar to the construction of the universal solutions of [FKMP05, Def. 2.4], in that both are based upon similar notions of endomorphism (there termed *homomorphism*). However, those universal solutions are not required to be tuple minimal. On the other hand, they are not limited to positive sentences, but rather apply to the more general class of XEIDs [Fag82].

The existence of a (finite) Armstrong model for a set of Φ is guaranteed under fairly simple circumstances; all that is necessary is that Φ have a finite cover.

3.10 Lemma *Let $\Phi \subseteq \Upsilon_K$. Then Φ admits a tuple-minimal Armstrong model with respect to Υ_K iff Φ admits a finite cover relative to Υ_K .*

PROOF: If Φ admits a finite cover, then Φ admits a tuple-minimal Armstrong model with respect to Υ_K by 3.9(b). Conversely, if Φ is not finite and has no finite cover, then for any positive integer n , there is a $\varphi \in \Phi$ which contains at least n distinct conjuncts and which is not equivalent to any finite subset of Υ_K , each of whose elements contains fewer than n conjuncts. An Armstrong model must thus contain at least n tuples. Since n may be chosen arbitrarily large, it follows that no such finite model can exist. \square

3.11 Canonical models Let K be a finite subset of $\text{Const}(\mathcal{D})$. In (a)-(c) and (e) below, take $\Phi \subseteq \Upsilon_K$ as well.

(a) Φ is \mathbf{D} -consistent if $\text{AtMod}_j(\Phi) \cup \text{LDB}(\mathbf{D}) \neq \emptyset$.

Thus, Φ is \mathbf{D} -consistent if there is some legal database which satisfies Φ . Such a database must also satisfy the sentences in $\text{Constr}(\mathbf{D})$; the total set of sentences which it must satisfy is the extended information, expressed formally as follows.

(b) Define the *extended information* of Φ with respect to Υ_K to be $\text{XInfo}_{\mathbf{D}}\langle\Phi, \Upsilon_K\rangle = \{\varphi \in \Upsilon_K \mid \Phi \cup \text{Constr}(\mathbf{D}) \models \varphi\}$.

Note that if Φ is not \mathbf{D} -consistent, then $\text{XInfo}_{\mathbf{D}}\langle\Phi, \Upsilon_K\rangle = \Upsilon_K$. Also note that, equivalently, $\text{XInfo}_{\mathbf{D}}\langle\Phi, \Upsilon_K\rangle = \{\varphi \in \Upsilon_K \mid (\forall M \in \text{LDB}(\mathbf{D}))((M \in \text{AtMod}_j(\Phi)) \Rightarrow (M \in \text{AtMod}_j(\varphi)))\}$ whenever Φ is \mathbf{D} -consistent. In other words, $\text{XInfo}_{\mathbf{D}}\langle\Phi, \Upsilon_K\rangle$ is the set of all sentences in Υ_K which are true in every $M \in \text{LDB}(\mathbf{D}) \cap \text{AtMod}_j(\Phi)$.

Since the databases of this paper are finite, consistency is not enough. Rather, Φ together with $\text{Constr}(\mathbf{D})$ must admit a finite model. In view of 3.10, this property is equivalent to Φ having a finite cover. Formally, this is recaptured as follows.

(c) Φ *extends finitely* to \mathbf{D} with respect to Υ_K if $\text{XInfo}_{\mathbf{D}}\langle\Phi, \Upsilon_K\rangle$ has a finite cover with respect to Υ_K .

- (d) The schema \mathbf{D} *admits finite extensions* with respect to Υ_K if every finite and \mathbf{D} -consistent $\Phi \subseteq \Upsilon_K$ extends finitely to \mathbf{D} with respect to Υ_K .
- (e) A *canonical database* for Φ in \mathbf{D} with respect to Υ_K is a tuple-minimal Armstrong model M for $\mathbf{XInfo}_{\mathbf{D}}\langle\Phi, \Upsilon_K\rangle$ with respect to Υ_K .

Observe that, in view of 3.9(c), canonical databases are unique up to automorphism. \mathbf{D} admits canonical databases conditionally if there is a canonical model whenever the extended information is finite, with unconditional extension requiring further that this finiteness condition always be satisfied.

- (f) The schema \mathbf{D} *admits canonical models conditionally* with respect to Υ_K if for every $\Phi \subseteq \Upsilon_K$ which extends finitely to \mathbf{D} with respect to Υ_K , every canonical database with respect to Υ_K is in $\mathbf{LDB}(\mathbf{D})$.
- (g) The schema \mathbf{D} *admits canonical models unconditionally* if \mathbf{D} admits finite extensions with respect to Υ_K and every canonical database with respect to Υ_K is in $\mathbf{LDB}(\mathbf{D})$.

These existence conditions are characterized precisely in the following lemma.

3.12 Lemma *Continue with the notation of 3.11 above.*

- (a) \mathbf{D} *admits canonical models conditionally with respect to Υ_K iff for every \mathbf{D} -consistent $\Phi \subseteq \Upsilon_K$, $\mathbf{XInfo}_{\mathbf{D}}\langle\Phi, \Upsilon_K\rangle \models_{\min} \varphi$ for every $\varphi \in \mathbf{Constr}(\mathbf{D})$,*
- (b) \mathbf{D} *admits canonical models unconditionally with respect to Υ_K iff it admits canonical models conditionally and $\mathbf{XInfo}_{\mathbf{D}}\langle\Phi, \Upsilon_K\rangle$ has a finite cover relative to Υ_K .*

PROOF: Both parts follow immediately from 3.10. \square

3.13 Example — Canonical models conditionally but not unconditionally It is not the case that every schema which admits canonical models conditionally admits them unconditionally. For example, let the schema \mathbf{E}_2 have three relational symbols $R_1[A]$, $R_2[AB]$, and $R_3[AB]$, with the inclusion dependencies $R_1[A] \subseteq R_2[A]$, $R_2[A] \subseteq R_3[A]$, and $R_3[B] \subseteq R_2[B]$. Let $M_1 = \{R_1(a_0), R_2(a_0, b_0), R_3(a_0, b_0), R_1(a_1)\}$, let $K = \{a_0, a_1, b_0\}$, and note that $M_1 \subseteq \mathbf{WFS}(\mathbf{E}_3, \exists\wedge+, K)$, since databases are taken to be sets of ground atoms.

In $\mathbf{XInfo}_{\mathbf{E}_3}\langle M_1, \mathbf{WFS}(\mathbf{E}_3, \exists\wedge+, K)\rangle$, a tuple of the form $R_2(a_1, b_1)$ must be present, which implies that one of the form $R_3(a_2, b_1)$ must be present as well, which in turn implies that one of the form $R_2(a_2, b_2)$ must be present, and so forth. If the constant symbols which are introduced to satisfy the dependencies, are always new ones which have not been used previously, then this construction proceeds indefinitely. In terms of the construction of the extended information $\mathbf{XInfo}_{\mathbf{D}}\langle M_1, \mathbf{WFS}(\mathbf{E}_3, \exists\wedge+, K)\rangle$, it is not difficult to see that an infinite increasing sequence $\langle\varphi_0, \varphi_1, \varphi_2, \dots, \varphi_i, \dots\rangle$ of sentences arises, as shown in Fig. 1, with φ_{i+1} strictly longer than φ_i and furthermore not a consequence of $\{\varphi_0, \varphi_1, \dots, \varphi_i\}$. Thus, $\mathbf{XInfo}_{\mathbf{E}_3}\langle M_1, \mathbf{WFS}(\mathbf{E}_3, \exists\wedge+, K)\rangle$ cannot have a finite cover. If the sequence is terminated, by choosing, say, $b_2 = b_1$, then an additional sentence beyond those in $\mathbf{XInfo}_{\mathbf{E}_2}\langle M_1, \mathbf{WFS}(\mathbf{E}_2, \exists\wedge+, K)\rangle$ is included, and so the resulting database is not Armstrong with respect to $\mathbf{WFS}(\mathbf{E}_3, \exists\wedge+, K)$. On the other hand, the database $M_2 = \{R_1(a_0), R_2(a_0, b_0), R_3(a_0, b_0)\}$ already satisfies every constraint in $\mathbf{XInfo}_{\mathbf{E}_3}\langle M_2, \mathbf{WFS}(\mathbf{E}_2, \exists\wedge+, K)\rangle$, and so is a

$$\begin{aligned}
& (\exists x_1)(R_2(a_1, x_1)) \\
& (\exists x_1)(\exists x_2)(R_2(a_1, x_1) \wedge R_3(x_1, x_2)) \\
& (\exists x_1)(\exists x_2)(\exists x_3)(R_2(a_1, x_1) \wedge R_3(x_1, x_2) \wedge R_2(x_3, x_2)) \\
& (\exists x_1)(\exists x_2)(\exists x_3)(\exists x_4)(R_2(a_1, x_1) \wedge R_3(x_1, x_2) \wedge R_2(x_3, x_2) \wedge R_3(x_3, x_4)) \\
& \quad \vdots
\end{aligned}$$

Figure 1: A strictly increasing sequence of information sentences

tuple-minimal Armstrong model of itself. Thus, the terminology *conditionally* is justified; \mathbf{E}_2 admits canonical models for some sets of sentences, but not for others. In 3.22, conditions under which canonical models are always admitted unconditionally are identified.

3.14 Example — Canonical models and positive disjunction While the definitions of 3.11 apply to any relational database schema, further restrictions must be imposed to render them meaningful. Consider again the schema \mathbf{E}_2 of Section 1, with the three unary relation symbols $R[A]$, $S[A]$, and $T[A]$, subject to the constraint $(\forall x)(R(x) \Leftrightarrow (S(x) \vee T(x)))$. For $M_1 = \{R(a_0)\}$ and $K = \{a_0\}$, $\text{XInfo}_{\mathbf{E}_2} \langle M_1, \text{WFS}(\mathbf{E}_2, \exists \wedge +, K) \rangle = M_1$, yet $M_1 \notin \text{LDB}(\mathbf{E}_2)$. The problem is that the “full” extended information, relative to $\text{WFS}(\mathbf{E}_2)$ is $\{R(a_0), S(a_0) \vee T(a_0)\}$, but the disjunction $S(a_0) \vee T(a_0)$ does not lie in $\text{WFS}(\mathbf{E}_2, \exists \wedge +)$. Hence the canonical database for M_1 with respect to $\text{WFS}(\mathbf{E}_2, \exists \wedge +)$ is not in $\text{LDB}(\mathbf{E}_2)$. It is clear that the notion of a canonical database is not really meaningful in the presence of such disjunctions. Rather, attention must be restricted to Horn dependencies, which avoid such positive disjunction and which are described below.

3.15 Generalized Horn dependencies The vast majority of relational database dependencies which have been considered over the years belong to a general class of logical formulas called *Horn clauses*. Originally presented as a characterization of formulas which are true under direct products [Hor51], they are more generally central to the modelling canonical instances in computer science [Mak87]. In the context of database dependencies, the following form is used, with each A_i and each B_i an atom.

$$(GHD) \quad (\forall x_1)(\forall x_2) \dots (\forall x_m)((A_1 \wedge A_2 \wedge \dots \wedge A_n) \Rightarrow (\exists y_1)(\exists y_2) \dots (\exists y_r)(B_1 \wedge B_2 \wedge \dots \wedge B_s))$$

In this work, such dependencies will be allowed in their most general form, which will be called *generalized Horn dependencies*, or *GHDs*. The only restrictions are the following.

- (ghd-1) Each GHD is in fact a sentence, so that each variable lies within the scope of a quantifier.
- (ghd-2) $\{x_1, x_2, \dots, x_m\} \cap \{y_1, y_2, \dots, y_r\} = \emptyset$.
- (ghd-3) Each x_i occurs in some A_j ; no universally quantified variable occurs only in a B_j .
- (ghd-4) $A_i \in \text{Atoms}(\mathbf{D})$ for each i .
- (ghd-5) If $s > 0$, then either each $B_i \in \text{Atoms}(\mathbf{D})$, in which the sentence is *tuple generating* (also called a *TGHD*), or else $s = 1, r = 0$, and $B_1 \in \text{VarEqAtoms}(\mathcal{D}) \cup \{\perp\}$, in which case

the sentence is called *equality generating* (also called an *EGHD*). If $B_1 = \perp$, the sentence is called a *mutual exclusion*. (\perp may be thought of as an equality which can never be satisfied, such as $(a = b)$.)

As a convenient notation, $\text{GHD}(\mathbf{D})$ will be used to denote the set of all GHDs on \mathbf{D} , with $\text{TGHD}(\mathbf{D})$ (resp. $\text{EGHD}(\mathbf{D})$) the subset consisting of the tuple-generating (resp. equality-generating) sentences. Mutual exclusions will be regarded as special cases of EGHDs in which there is no atom on the right-hand side.

The GHDs are a generalization of the XEIDs of Fagin [Fag82, Sec. 7], and are essentially the source-to-target dependencies of [FKMP05, Def. 2.1]. As such, the GHDs encompass virtually all classes of database constraints which have been studied, including in particular functional and inclusion dependencies. In contrast to XEIDs, the left-hand side need be neither unirelational nor typed. Of course, the more stringent requirement on XEIDs is made for a reason — XEIDs enjoy the property of possessing Armstrong models [Fag82, Thm. 3.1] which GHDs do not. Although they were used as the general class of dependency in [Heg08a], it turns out that this property of possessing Armstrong models is only required of sentences in $\text{WFS}(\mathbf{D}, \exists\wedge+)$ which are used to characterize the information content of database, and not for more general Horn clauses which are used to characterize the underlying constraints. Therefore, there is no need to enforce the additional requirements of XEIDs.

GHDs also generalize XEIDs in a less essential way — there is no requirement that n be greater than zero, although both are not allowed to be zero in the same clause. If $n = 0$, a sentence in $\text{WFS}(\mathbf{D}, \exists\wedge+)$ is obtained; thus, $\text{WFS}(\mathbf{D}, \exists\wedge+) \subseteq \text{GHD}(\mathbf{D})$. As a specific example of such a constraint, consider $(\exists y_1)(\exists y_2)(R(y_1, y_2))$, which states that the relation instance for R is always nonempty. Additionally, constant symbols are allowed in a GHD. For example, a constraint of the form $(\forall x_1)(\forall x_2)(R(x_1, x_2) \Rightarrow S(x_1, x_2, \nu))$ might assert that for every tuple in R , a similar tuple, padded with the constant ν in the third position (with ν representing a null value, for example), is required. Finally, an unsatisfiable right-hand side is allowed. If $B_1 = \perp$, a statement of the form $(\forall x_1)(\forall x_2) \dots (\forall x_m)((A_1 \wedge A_2 \wedge \dots \wedge A_n) \Rightarrow \perp)$ is obtained, with \perp denoting the identically false assertion. As noted in (ghd-5) above, such a sentence is called a mutual exclusion. An example is the antisymmetry constraint $(\forall x_1)(\forall x_2)((R(x_1, x_2) \wedge R(x_2, x_1)) \Rightarrow \perp)$.

The idea of forward chaining on databases by applying Horn-style rules is well known [DG84], and forms one of the cornerstones of logic programming [Llo87]. However, in the current context, rather than reasoning on ground atoms, it is essential to apply forward chaining on information — that is, on sentences in Υ_K . This idea is addressed via notions of information associated with TGHDs, as developed below.

3.16 Information inference for TGHDs In so-called forward chaining in classical propositional logic, given a Horn clause of the form $A_1 \wedge A_2 \wedge \dots \wedge A_n \Rightarrow B$, if all of the A_i 's are known to be true, the rule may “fire” and assert that B is also true. This idea also works for TGHDs operating on ground atoms; a ground substitution is applied to the left-hand side, and if all of the resulting ground atoms are true, the sentence obtained by applying the same substitution to the right-hand side must be true, and further substitutions may then be applied to identify the ground atoms which it implies. More generally, however, when the pool of knowledge consists of sentences in $\text{WFS}(\mathbf{D}, \exists\wedge+)$, the left-hand side and right-hand sides are coupled via variables. In other words, a rule fires for specific bindings of variable common to both sides. It is therefore

not possible to separate the conclusion (*i.e.*, the right-hand side) from the hypotheses (*i.e.*, the left-hand side). The solution is a rather simple one. When a rule of the above form fires, rather than concluding simply B from $A_1 \wedge A_2 \wedge \dots \wedge A_n$, the conjunction $A_1 \wedge A_2 \wedge \dots \wedge A_n \wedge B$ is deduced as a new fact. In this way, any binding of quantified variables which occurred during the inference process are preserved.

Let φ be a GHD of the form (GHD) of 3.15 above, and let s be a substitution on $\{x_1, x_2, \dots, x_n\}$. Call s *GHD-compatible with φ* if $s(x_i) \neq y_j$ for any indices i and j , and assume that s has this property. GHD compatibility ensures only that a substitution does not rename a universally quantified variable to coincide with one which is existentially quantified in the original formula. The *left-hand-side information* of φ with respect to s is the sentence obtained by applying s to the left-hand side of φ . Formally, $\text{LHSinfo}\langle\varphi, s\rangle$ is the sentence in $\text{WFS}(\mathbf{D}, \exists \wedge +)$ obtained from $(\forall x_1)(\forall x_2) \dots (\forall x_m)(A_1 \wedge A_2 \wedge \dots \wedge A_n)$ as follows.

- (lhs-i) For $i \in \{1, 2, \dots, m\}$, if $s(x_i) = v$ and $v \in \text{GenVars}(\mathcal{D})$, replace $(\forall x_i)$ with $(\exists v)$.
- (lhs-ii) For $i \in \{1, 2, \dots, m\}$, if $s(x_i) \in \text{Const}(\mathcal{D})$, delete $(\forall x_i)$.
- (lhs-iii) For $i \in \{1, 2, \dots, n\}$, replace each A_i with the sole element of $\text{Subst}(A_i, s)$.

Note in particular that universally quantified variables become existentially quantified. The existentially quantified versions represent a single but unspecified binding on those formerly universally quantified positions. Now, the *left-plus-right-hand-side information* of φ with respect to s , denoted $\text{LRHSinfo}\langle\varphi, s\rangle$, is the sentence obtained from

$$(\forall x_1)(\forall x_2) \dots (\forall x_m)(\exists y_1)(\exists y_2) \dots (\exists y_r)(A_1 \wedge A_2 \wedge \dots \wedge A_n \wedge B_1 \wedge B_2 \wedge \dots \wedge B_s)$$

by following the steps (i)-(iii) above and, in addition to the following step.

- (rhs) Replace each B_i with $\text{Subst}(B_i, s)$.

For example, if $\varphi = (\forall x_1)(\forall x_2)(R(x_1, x_2) \Rightarrow (\exists y)(S(x_1, x_2, y)))$ and $s = \{a_1/x_1, x/x_2\}$ then $\text{LHSinfo}\langle\varphi, s\rangle = (\exists x)(R(a_1, x))$ and $\text{LRHSinfo}\langle\varphi, s\rangle = (\exists x)(\exists y)(R(a_1, x) \wedge S(a_1, x, y))$.

The following lemma states formally the intuition that if $\text{LHSinfo}\langle\varphi, s\rangle$ is satisfied, and the rule φ holds, then $\text{LRHSinfo}\langle\varphi, s\rangle$ holds as well. It is an immediate consequence of the above definitions.

3.17 Lemma *Let $\varphi \in \text{TGHD}(\mathbf{D})$ and let $M \in \text{DB}(\mathbf{D})$. Then $M \in \text{AtMod}_J(\varphi)$ iff for every GHD-compatible substitution s on $\text{Vars}(\varphi)$, if $\text{LHSinfo}\langle\varphi, s\rangle \in \text{Info}\langle M, \Upsilon_K \rangle$, then $\text{LRHSinfo}\langle\varphi, s\rangle \in \text{Info}\langle M, \Upsilon_K \rangle$ as well. \square*

The case of EGHDs must be handled a bit differently, since the result of a deduction is not a new sentence in Υ_K but rather a restriction on existing sentences. Therefore, the crucial idea is to consider the effect of *reducing* a given sentence $\psi \in \Upsilon_K$ by an EGHD φ , with the latter possibly forcing certain terms of ψ to be equal.

3.18 Information associated with an EGHD Let φ be an EGHD of the form (GHD) of 3.15 above, and let $\psi \in \Upsilon_K$. Here ψ is a sentence to which the φ will be applied. If ψ can be unified with the left-hand side of φ , then the equality defined by the right-hand side of φ must be applied to ψ .

For this to work, a suitable substitution must be applied, so s be any substitution on $\{x_1, x_2, \dots, x_n\}$, the variables of φ . (In an EGHD, there are no existentially quantified variables of the form y_i in (GHD).) Note further that the GHD-compatibility property on a substitution does not apply here, since every EGHD is a universal sentence. However, the result of applying s to the left-hand side of φ must match ψ , so that the rule can fire. Formally, call s φ -compatible for ψ if $\text{AtRep}(\text{LHSinfo}\langle\varphi, s\rangle) \subseteq \text{AtRep}(\psi)$. Thus, s is φ -compatible if, when applied to the left-hand side of φ , the conjuncts which are obtained (after removing any quantifiers) are a subset of those of ψ . Let $\text{RHSinfo}\langle\varphi, s\rangle$ denote the equality atom obtained by applying the substitution s to B_1 , the right-hand side of φ . (If φ is a mutual exclusion, then $\text{RHSinfo}\langle\varphi, s\rangle = \perp$.) The reduction of ψ by $\text{RHSinfo}\langle\varphi, s\rangle$, denoted $\text{Reduction}\langle\psi, \text{RHSinfo}\langle\varphi, s\rangle\rangle$ is defined by cases as follows.

- (red-i) If $\text{RHSinfo}\langle\varphi, s\rangle$ is of the form $(a_i = a_j)$ with $i \neq j$, or $\text{RHSinfo}\langle\varphi, s\rangle = \perp$, then $\text{Reduction}\langle\psi, \text{RHSinfo}\langle\varphi, s\rangle\rangle = \perp$, the identically false assertion.
- (red-ii) If $\text{RHSinfo}\langle\varphi, s\rangle$ is of the form $(a_i = a_i)$ or $(x_i = x_i)$, then $\text{Reduction}\langle\psi, \text{RHSinfo}\langle\varphi, s\rangle\rangle = \psi$.
- (red-iii) If $\text{RHSinfo}\langle\varphi, s\rangle$ is of the form $(x_i = a_j)$ or $(a_j = x_i)$, then $\text{Reduction}\langle\psi, \text{RHSinfo}\langle\varphi, s\rangle\rangle$ is obtained by substituting a_j for x_i in ψ and removing the quantifier $(\exists x_i)$.
- (red-iv) If $\text{RHSinfo}\langle\varphi, s\rangle$ is of the form $(x_i = x_j)$ with $i \neq j$, then $\text{Reduction}\langle\psi, \text{RHSinfo}\langle\varphi, s\rangle\rangle$ is obtained by substituting x_i for x_j in ψ and removing the quantifier $(\exists x_j)$. (Note that the quantifier $(\exists x_i)$ must also be present in this case, and is not removed.)

For example, if $\varphi = (\forall x_1)(\forall x_2)(\forall x_3)((R(x_1, x_2) \wedge R(x_1, x_3)) \Rightarrow (x_2 = x_3))$ and $\psi = (\exists x_1)(\exists x_2)(R(a_1, x_1) \wedge R(a_1, x_2) \wedge S(x_1, x_2))$ then for $s = \{(a_1/x_1, x_1/x_2, x_2/x_3)\}$, $\text{LHSinfo}\langle\varphi, s\rangle = (\exists x_1)(\exists x_2)(R(a_1, x_1) \wedge R(a_1, x_2))$ and so is φ -compatible. $\text{RHSinfo}\langle\varphi, s\rangle = (x_1 = x_2)$, and so $\text{Reduction}\langle\psi, \text{RHSinfo}\langle\varphi, s\rangle\rangle = (\exists x_1)(\exists x_2)(R(a_1, x_1) \wedge R(a_1, x_1) \wedge S(x_1, x_1))$. On the other hand, if $\psi = (\exists x_1)(R(x_1, a_1) \wedge R(x_1, a_2))$ then $s = \{x_1/x_1, a_1/x_2, a_2/x_3\}$ is φ -compatible but $\text{RHSinfo}\langle\varphi, s\rangle = (a_1 = a_2)$. and so $\text{Reduction}\langle\psi, \text{RHSinfo}\langle\varphi, s\rangle\rangle = \perp$.

The following lemma is a routine consequence of the above.

3.19 Lemma *Let $\varphi \in \text{EGHD}(\mathbf{D})$, let $M \in \text{DB}(\mathbf{D})$, and let Ψ be a cover for $\text{Info}\langle M, \Upsilon_K \rangle$ with respect to Υ_K . Then $M \in \text{AtMod}_j(\varphi)$ iff for every $\psi \in \Psi$ and every substitution s which is φ -compatible for ψ , $\text{Reduction}\langle\psi, \text{RHSinfo}\langle\varphi, s\rangle\rangle \in \text{Info}\langle M, \Upsilon_K \rangle$ as well. \square*

Finally, the main result, that schemata constrained by GHDs always admit canonical models conditionally, may be established.

3.20 Theorem — Conditional existence of canonical models *If $\text{Constr}(\mathbf{D}) \subseteq \text{GHD}(\mathbf{D})$, then for any finite $K \subseteq \text{Const}(\mathcal{D})$, \mathbf{D} admits canonical models conditionally with respect to Υ_K .*

PROOF: Let $\Phi \subseteq \text{WFS}(\mathbf{D}, \exists \wedge +, K)$ have the property that it extends finitely to \mathbf{D} for Υ_K . First, let $\varphi \in \text{Constr}(\mathbf{D})$ be a TGHD, and let s be a GHD-compatible substitution for φ . It follows directly from 3.17 and the definition of XInfo (3.11(b)) that whenever $\text{LHSinfo}\langle\varphi, s\rangle \in$

$\mathbf{XInfo}_{\mathbf{D}}\langle\Phi, \Upsilon_K\rangle$, then $\text{LRHSinfo}\langle\varphi, s\rangle \in \mathbf{XInfo}_{\mathbf{D}}\langle\Phi, \Upsilon_K\rangle$ as well. A second application of 3.17 establishes that $M \in \text{MinAtMod}_{\mathcal{J}}(\varphi)$ for every tuple-minimal Armstrong model M of $\mathbf{XInfo}_{\mathbf{D}}\langle\Phi, \Upsilon_K\rangle$ with respect to Υ_K .

The proof for $\varphi \in \text{Constr}(\mathbf{D})$ an EGHD is similar. Choose $\psi \in \Phi$ and let s be a φ -compatible substitution for ψ . In view of 3.19 and the definition of \mathbf{XInfo} , it must be the case that $\text{Reduction}\langle\psi, \text{RHSinfo}\langle\varphi, s\rangle\rangle \in \mathbf{XInfo}_{\mathbf{D}}\langle\Phi, \Upsilon_K\rangle$ as well. Again, a second application of 3.19 establishes that $M \in \text{MinAtMod}_{\mathcal{J}}(\varphi)$ for every tuple-minimal Armstrong model M of $\mathbf{XInfo}_{\mathbf{D}}\langle\Phi, \Upsilon_K\rangle$ with respect to Υ_K . \square

3.21 Weakly acyclic TGHDs To admit canonical models unconditionally, it is necessary to ensure that infinite increasing sequences in $\mathbf{XInfo}_{\mathbf{D}}\langle\Phi, \Upsilon_K\rangle$, as illustrated in 3.13, do not occur. Such infinite models are related to cycles in the tuple-generating dependencies. In [FKMP05, Def. 3.7], the notion of a *weakly acyclic set of TGDs* is developed, and it is shown [FKMP05, Thm. 3.9] that the chase procedure always terminates when the dependencies are limited to an acyclic set of TGDs together with EGDs. The TGDs and EGDs of [FKMP05] differ only in relatively minor ways from the TGHDs and EGHDs of this paper; in particular, the result extends directly to TGDs. The details will not be worked out, but the following result is noted for completeness.

3.22 Corollary — (to 3.20) *If $\text{Constr}(\mathbf{D})$ is finite and consists of a weakly acyclic set of TGHDs, together with any set of EGHDs, then for any finite $K \subseteq \text{Const}(\mathcal{D})$, \mathbf{D} admits canonical models unconditionally with respect to Υ_K .*

PROOF: Combine the result of 3.20 with [FKMP05, Thm. 3.9]. \square

4 Optimal Reflection of Insertions

The focus is now turned to the problem of characterizing optimal reflections of insertions into the view schema. Roughly, the idea is to reflect the information which must be added to the main schema and then construct an Armstrong model (with respect to Υ_K for a suitable K) of that information together with the current instance of the main schema. There are, of course, details to be developed and pitfalls to be avoided, all of which are discussed in this section. First of all, some basic definitions surrounding updates and information are developed.

4.1 Notational convention Throughout the rest of this paper, unless stated specifically to the contrary, take $\Gamma = (\mathbf{V}, \gamma)$ to be a relational view of \mathbf{D} of class $\exists\wedge+$.

4.2 Updates and reflections An *update* on \mathbf{D} is a pair $(M_1, M_2) \in \text{LDB}(\mathbf{D}) \times \text{LDB}(\mathbf{D})$. M_1 is the current instance, and M_2 the new instance. It is an *insertion* if $M_1 \subseteq M_2$, and a *deletion* if $M_2 \subseteq M_1$.

To describe the situation surrounding an update request on Γ , it is sufficient to specify the current instance M_1 of the main schema and the desired new instance N_2 of the view schema \mathbf{V} . The current instance of the view can be computed as $\gamma(M_1)$; it is only the new instance M_2 of the main schema (subject to $N_2 = \gamma(M_2)$) which must be obtained from an update strategy. Formally, an *update request* from Γ to \mathbf{D} is a pair (M_1, N_2) in which $M_1 \in \text{LDB}(\mathbf{D})$ (the old

instance of the main schema) and $N_2 \in \text{LDB}(\mathbf{V})$ (the new instance of the view schema). If $\gamma(M_1) \subseteq N_2$, it is called an *insertion request*, and if $N_2 \subseteq \gamma(M_1)$, it is called a *deletion request*. Collectively, insertion requests and deletion requests are termed *unidirectional update requests*. A *realization* of (M_1, N_2) along Γ is an update (M_1, M_2) on \mathbf{D} with the property that $\gamma(M_2) = N_2$. The update (M_1, M_2) is called a *reflection* (or *translation*) of the view update $(\gamma(M_1), N_2)$. The set of all realizations of (M_1, N_2) along Γ is denoted $\text{UpdRealiz}\langle(M_1, N_2), \Gamma\rangle$. The subset of $\text{UpdRealiz}\langle(M_1, N_2), \Gamma\rangle$ consisting of insertions (resp. deletions) is denoted $\text{InsRealiz}\langle(M_1, N_2), \Gamma\rangle$ (resp. $\text{DelRealiz}\langle(M_1, N_2), \Gamma\rangle$).

4.3 Update difference and optimal reflections The update difference of an update (M_1, M_2) on \mathbf{D} with respect to a set $\Sigma \subseteq \text{WFS}(\mathbf{D})$ is a measure of how much M_2 differs from M_1 in terms of satisfaction of the sentences of Σ . Formally, the *positive* (Δ^+), *negative* (Δ^-), and *total* (Δ) *update differences* of (M_1, M_2) with respect to Σ are defined as follows:

$$\begin{aligned}\Delta^+\langle(M_1, M_2), \Sigma\rangle &= \text{Info}\langle M_2, \Sigma\rangle \setminus \text{Info}\langle M_1, \Sigma\rangle \\ \Delta^-\langle(M_1, M_2), \Sigma\rangle &= \text{Info}\langle M_1, \Sigma\rangle \setminus \text{Info}\langle M_2, \Sigma\rangle \\ \Delta\langle(M_1, M_2), \Sigma\rangle &= \Delta^+\langle(M_1, M_2), \Sigma\rangle \cup \Delta^-\langle(M_1, M_2), \Sigma\rangle\end{aligned}$$

Note that, given $\varphi \in \Delta\langle(M_1, M_2), \Sigma\rangle$, it is always possible to determine whether $\varphi \in \Delta^+\langle(M_1, M_2), \Sigma\rangle$ or $\varphi \in \Delta^-\langle(M_1, M_2), \Sigma\rangle$ by checking whether or not $M_1 \in \text{AtMod}_j(\varphi)$. Given an update request (M_1, N_2) , the quality of a realization (M_1, M_2) is measured by its update difference. Formally, let $\Sigma \subseteq \text{WFS}(\mathbf{D})$, let (M_1, N_2) be an update request from Γ to \mathbf{D} , let $T \subseteq \text{UpdRealiz}\langle(M_1, N_2), \Gamma\rangle$, and let $(M_1, M_2) \in T$.

- (a) (M_1, M_2) is *minimal* in T with respect to Σ if for any $(M_1, M'_2) \in T$, if $\Delta\langle(M_1, M'_2), \Sigma\rangle \subseteq \Delta\langle(M_1, M_2), \Sigma\rangle$, then $\Delta\langle(M_1, M'_2), \Sigma\rangle = \Delta\langle(M_1, M_2), \Sigma\rangle$.
- (b) (M_1, M_2) is *least* in T with respect to Σ if for all $(M_1, M'_2) \in T$, $\Delta\langle(M_1, M_2), \Sigma\rangle \subseteq \Delta\langle(M_1, M'_2), \Sigma\rangle$.

4.4 Update classifiers An *update classifier* for \mathbf{D} is simply a set Σ of information-monotone sentences. In this work, the set Σ will always be taken to be either $\text{GrAtoms}(\mathbf{D})$ or else $\text{WFS}(\mathbf{D}, \exists\wedge+, K) = \Upsilon_K$ for an appropriate set K of constants.

Let (M_1, N_2) be an update request from Γ to \mathbf{D} , let $T \subseteq \text{UpdRealiz}\langle(M_1, N_2), \Gamma\rangle$, and let $(M_1, M_2) \in T$.

- (a) (M_1, M_2) is $\langle\Upsilon_K, T\rangle$ -*admissible* if it is minimal in T with respect to both Υ_K and $\text{GrAtoms}(\mathbf{D})$.
- (b) (M_1, M_2) is $\langle\Upsilon_K, T\rangle$ -*optimal* if it is $\langle\Upsilon_K, T\rangle$ -admissible and least in T with respect to Σ .

Roughly, (M_1, M_2) is admissible if no other realization is better, and it is optimal if it is better than all others, up to the equivalence defined by Σ . Observe that if some update request is $\langle\Upsilon_K, T\rangle$ -optimal, then all $\langle\Upsilon_K, T\rangle$ -admissible update requests are $\langle\Upsilon_K, T\rangle$ -optimal.

As a notational shorthand, if $T = \text{InsRealiz}\langle(M_1, N_2), \Gamma\rangle$ (resp. $T = \text{DelRealiz}\langle(M_1, N_2), \Gamma\rangle$), that is, if T is the set of all possible insertions (resp. deletions) which realize (M_1, N_2) , then $\langle\Upsilon_K, T\rangle$ -admissible and $\langle\Upsilon_K, T\rangle$ -optimal will be abbreviated to $\langle\Upsilon_K, \uparrow\rangle$ -admissible and $\langle\Upsilon_K, \uparrow\rangle$ -optimal (resp. $\langle\Upsilon_K, \downarrow\rangle$ -admissible and $\langle\Upsilon_K, \downarrow\rangle$ -optimal).

For $\Sigma = \text{GrAtoms}(\mathbf{D})$, admissibility reduces to minimality in the sense of symmetric difference of sets as sketched in 1.2. More concretely, given an update request (M_1, N_2) , a realization (M_1, M_2) is $\langle \text{GrAtoms}(\mathbf{D}), T \rangle$ -admissible if for no other realization $(M_1, M'_2) \in T$ is it the case that $\text{SymDiff}\langle M_1, M'_2 \rangle \subseteq \text{SymDiff}\langle M_1, M_2 \rangle$. Similarly, (M_1, M_2) is $\langle \text{GrAtoms}(\mathbf{D}), T \rangle$ -optimal if for every other realization (M_1, M'_2) , $\text{SymDiff}\langle M_1, M_2 \rangle \subseteq \text{SymDiff}\langle M_1, M'_2 \rangle$. Minimality with respect to $\text{GrAtoms}(\mathbf{D})$ is referred to as *tuple minimality*, in harmony with the terminology already introduced for Armstrong models in Section 3.

The major theme of this paper is that tuple minimality, by itself, is not sufficient to characterize optimal updates. Rather, optimality with respect to Υ_K for a suitably chosen set K of constants is also essential. This issue is next addressed in detail.

4.5 The constants associated with an update request In reflecting an update from a view to the main schema, the use of new constants in generic models is crucial. For a constant to be “new”, it is not sufficient that it merely not appear in the current or proposed view instance, or the current instance of the main schema. Rather, it must not appear in any constraint or defining formula for the main schema or view. For example, referring back to the example of the null-value constraint $(\forall x_1)(\forall x_2)(R(x_1, x_2) \Rightarrow S(x_1, x_2, \nu))$ of 3.15, it would be inappropriate to use ν as a generic constant, since it already has another meaning within the global context of all databases. Thus, the pool of generic constants must also exclude any which occur in constraints or in the defining formulas for views. Specifically, the following sequence of definitions leads to the acceptable pool of generic constants.

First of all, define the constant symbols of the schema \mathbf{D} to be those of its constraints.

$$(a) \text{ConstSym}(\mathbf{D}) = \text{ConstSym}(\text{Constr}(\mathbf{D}))$$

Next, define the constant symbols of the view $\Gamma = (\mathbf{V}, \gamma)$ to be those of the main schema \mathbf{D} , together with those of both the view schema \mathbf{V} and the view mapping γ .

$$(b) \text{ConstSym}(\Gamma) = \text{ConstSym}(\text{Constr}(\mathbf{D})) \cup \text{ConstSym}(\text{Constr}(\mathbf{V})) \cup \text{ConstSym}(\gamma)$$

The constant symbols of an update request u are those of its instances.

$$(c) \text{ For } u = (M_1, N_2) \text{ an update request from } \Gamma \text{ to } \mathbf{D}, \text{ConstSym}(u) = \text{ConstSym}(M_1) \cup \text{ConstSym}(N_2).$$

Finally, define the *total constant set of u* , denoted \mathfrak{C}_u , to be the constant symbols of u together with those of Γ .

$$(d) \mathfrak{C}_u = \text{ConstSym}(\Gamma) \cup \text{ConstSym}(u).$$

Note that $\text{ConstSym}(N_2) \subseteq \text{ConstSym}(M_1) \cup \text{ConstSym}(\gamma) \cup \text{ConstSym}(\mathbf{V})$, so the following more compact representation is also valid.

$$(d') \mathfrak{C}_u = \text{ConstSym}(\Gamma) \cup \text{ConstSym}(M_1).$$

Since database schemata, even those of views, are always taken to be constant finite (see 2.5), it follows that \mathfrak{C}_u is always a finite set. Since $\text{Const}_{\mathcal{D}}(A)$ is infinite for every $A \in \mathcal{A}_{\mathcal{D}}$, it follows that there are always infinitely many “available” constant symbols for each attribute A which do not lie in \mathfrak{C}_u . The set \mathfrak{C}_u will thus be taken to be the set of constant symbols which may not be used as generic values in constructing Armstrong models.

4.6 Information lifting Let $N \in \text{DB}(\mathbf{V})$. The *information lifting* of N along Γ , denoted $\text{InfoLift}\langle N, \Gamma \rangle$, is the minimum information in $\text{WFS}(\mathbf{D}, \exists \wedge +)$ which any $M \in \text{DB}(\mathbf{D})$ with $\gamma(M) = N$ must have. Formally, this is recaptured as follows.

$$\text{InfoLift}\langle N, \Gamma \rangle = \{\text{Substf}\langle \gamma, t \rangle \mid t \in N\}$$

Note that $\text{InfoLift}\langle N, \Gamma \rangle \subseteq \Upsilon_K$ with $K = \text{ConstSym}(N_2) \cup \text{ConstSym}(\gamma)$.

Given an insertion request $u = (M_1, N_2)$ from Γ to \mathbf{D} , and let $M_2 \in \text{InsRealiz}\langle (M_1, N_2), \Gamma \rangle$. The least information which M_2 must have is $\text{InfoLift}\langle N_2, \Gamma \rangle \cup M_1$, closed up under the constraints in $\text{Constr}(\mathbf{D})$. This is recaptured formally as

$$\text{LeastRefl}\langle u, \Gamma \rangle = \text{XInfo}_{\mathbf{D}}\langle M_1 \cup \text{InfoLift}\langle N_2, \Gamma \rangle, \Upsilon_{\mathfrak{C}_u} \rangle$$

and is called the *least reflection of u along Γ* .

In the context of the example of 1.3, with $u = (M_{00}, N_{01})$, the information lifting $\text{InfoLift}\langle N_{01}, \Pi_{R'[AB]}^{\mathbf{E}_1} \rangle = \{(\exists x_1)(\exists x_2)(\exists x_3)(R(a_0, b_0, x_1) \wedge R(a_1, b_1, x_2) \wedge R(a_2, b_2, x_3))\}$, while the least reflection $\text{LeastRefl}\langle u, \Pi_{R'[AB]}^{\mathbf{E}_1} \rangle = M_{00} \cup \{(\exists x_1)(\exists x_2)(R(a_2, b_2, x_1) \wedge (S(x_1, x_2)))\}$.

Using the construction of 3.9, a tuple-minimal Armstrong model of $\text{LeastRefl}\langle u, \Pi_{R'[AB]}^{\mathbf{E}_1} \rangle$ is obtained by replacing the variables by distinct and new constants; for example $M_{01} = M_{00} \cup \{R(a_2, b_2, c_2), S(c_2, d_2)\}$.

4.7 Proposition — Characterization of optimal insertions *Assume that \mathbf{D} admits canonical models conditionally, and let $u = (M_1, N_2)$ be an insertion request from Γ to \mathbf{D} . Then (M_1, M_2) is a $\langle \Upsilon_{\mathfrak{C}_u}, \uparrow \rangle$ -optimal realization of u iff the following two conditions hold.*

- (i) M_2 is a tuple-minimal Armstrong model of $\text{LeastRefl}\langle u, \Upsilon_{\mathfrak{C}_u} \rangle$ with respect to $\Upsilon_{\mathfrak{C}_u}$.
- (ii) $\gamma(M_2) = N_2$.

PROOF: Certainly, if (i) and (ii) hold, then M_2 is a $\langle \Upsilon_{\mathfrak{C}_u}, \uparrow \rangle$ -optimal realization of u . On the other hand, if M_2 is a $\langle \Upsilon_{\mathfrak{C}_u}, \uparrow \rangle$ -optimal realization of u , then (ii) holds trivially. Since \mathbf{D} admits canonical models conditionally, if $\text{LeastRefl}\langle u, \Upsilon_{\mathfrak{C}_u} \rangle$ admits a finite cover, then a tuple-minimal Armstrong model of that set is a $\langle \Upsilon_{\mathfrak{C}_u}, \uparrow \rangle$ -optimal realization of u , just by construction. If $\text{LeastRefl}\langle u, \Upsilon_{\mathfrak{C}_u} \rangle$ does not admit a finite cover, then, as illustrated in 3.13, it must contain an infinite increasing sequence $\langle \varphi_0, \varphi_1, \varphi_2, \dots \rangle$ of sentences with $\text{AtMod}_{\mathcal{J}}(\bigcup\{\varphi_i \mid 0 \leq i \leq k\}) \not\subseteq \text{AtMod}_{\mathcal{J}}(\varphi_k)$ for any $k > 0$, and so it cannot have a finite Armstrong model. Thus, any $\langle \Upsilon_{\mathfrak{C}_u}, \uparrow \rangle$ -admissible realization must satisfy some sentence not in $\text{LeastRefl}\langle u, \Upsilon_{\mathfrak{C}_u} \rangle$, and so no $\langle \Upsilon_{\mathfrak{C}_u}, \uparrow \rangle$ -optimal realization can exist. \square

4.8 Example — Orphan tuples Even in the case that $\text{LeastRefl}\langle u, \Upsilon_{\mathfrak{C}_u} \rangle$ satisfies condition (i) of 4.7, it may not satisfy condition (ii); that is, it may not be the case that $\gamma(M_2) = N_2$. The problem lies with so-called *orphan tuples*, which are illustrated via the following example.

Let \mathbf{E}_4 be the schema having the single relation symbol $R[AB]$, constrained by the dependency $(\exists x_1)(\exists x_2)(R(x_1, x_2))$ which simply asserts that the instance of R is nonempty. It is immediate that this dependency is a TGH (with $m = n = 0$ in the pattern (GHD) of 3.15). The view $\Pi_{A+B}^{\mathbf{E}_4} = (\mathbf{W}_4, \pi_{A+B}^{\mathbf{E}_4})$ has two relation symbols $R_A[A]$ and $R_B[B]$, defined by the obvious projections $\pi_A^{\mathbf{E}_4}$ and $\pi_B^{\mathbf{E}_4}$. The only constraints on the view schema are

$(\exists x_1)(R_A(x_1))$ and $(\exists x_1)(R_B(x_1))$, so $\text{Constr}(\mathbf{W}_4)$ consists of GHDs as well. Now with $M_1 = \{R(a_0, b_0), R(a_1, b_1)\}$ the instance of \mathbf{E}_4 , $\pi_{A+B}^{\mathbf{E}_4}(M_1) = N_1 = \{R_A(a_0), R_A(a_1), R_B(b_0), R_B(b_1)\}$. Let N_2 be the view instance obtained by inserting $R_A(a_2)$ into N_1 , and define $u = (M_1, N_2)$. Then $\text{LeastRefl}\langle u, \Upsilon_{\mathfrak{C}_u} \rangle = M_1 \cup (\exists x_1)(R(a_2, x_1))$. A canonical model M_2 of this least reflection is of the form $M_1 \cup \{R(a_2, b_2)\}$, with $b_2 \notin \mathfrak{C}_u$, so $\pi_{A+B}^{\mathbf{E}_4}(M_2) = N_2 \cup \{R_B(b_2)\}$. Here $R_B(b_2)$ is termed an *orphan tuple*; it represents newly inserted information which has made its way back to the view. The tuple can be made to “disappear” by replacing b_2 with an existing value for attribute B , say b_1 . However, in that case, while the resulting instance $M'_2 = M_1 \cup \{R(a_2, b_1)\}$ does map to N_2 under $\pi_{A+B}^{\mathbf{E}_4}$, $\text{Info}\langle M_2, \Upsilon_{\mathfrak{C}_u} \rangle$ is a proper subset of $\text{Info}\langle M'_2, \Upsilon_{\mathfrak{C}_u} \rangle$, and so M'_2 is not $\langle \Upsilon_{\mathfrak{C}_u}, \uparrow \rangle$ -optimal. Thus, (M_1, M'_2) is a $\langle \Upsilon_{\mathfrak{C}_u}, \uparrow \rangle$ -admissible solution to the update request (M_1, N_2) , but it is not optimal. This problem cannot be made to disappear via clever formulation; in this example, there are no optimal solutions. Fortunately, orphan tuples can be ruled out by requiring that Γ reflect deletions, as described below.

4.9 Reflection of deletions The view $\Gamma = (\mathbf{V}, \gamma)$ *reflects deletions* if every deletion request (M_1, N_2) from Γ to \mathbf{D} admits a realization which is itself a deletion.

In the example of 4.8 above, $\Pi_{A+B}^{\mathbf{E}_4}$ does not reflect deletions. There is no realization of the update request (M'_1, N_2) which is a deletion, since with $M_1 = \{R(a_0, b_0), R(a_1, b_1)\}$ the instance of the main schema \mathbf{E}_4 , there is no way to realize the deletion of $R_2(b_1)$ from the view instance N_1 as a deletion from M_1 ; $R(a_1, b_1)$ must be deleted, which would also remove $R_1(a_1)$.

4.10 Lemma — Reflection of deletions implies no orphan tuples *Assume that \mathbf{D} admits canonical models conditionally, and let $u = (M_1, N_2)$ be an insertion request from Γ to \mathbf{D} for which $\text{LeastRefl}\langle u, \Upsilon_{\mathfrak{C}_u} \rangle$ admits a tuple-minimal Armstrong model M_2 with respect to $\Upsilon_{\mathfrak{C}_u}$. Then, if Γ reflects deletions, $\gamma(M_2) = N_2$.*

PROOF: It is clear that $N_2 \subseteq \gamma(M_2)$. By the definition of reflection of deletions, there is an $M'_2 \in \text{LDB}(\mathbf{D})$ with $M'_2 \subseteq M_2$ and $\gamma(M'_2) = N_2$. However, since M_2 is already tuple minimal, it follows that $M'_2 = M_2$. Thus, $\gamma(M_2) \setminus N_2 = \emptyset$, and so $\gamma(M_2) = N_2$, as required. \square

Finally, it can be established that for views with reflect deletions, condition (ii) of 4.7 is superfluous.

4.11 Proposition *Assume that Γ reflects deletions, that \mathbf{D} admits canonical models conditionally, and let $u = (M_1, N_2)$ be an insertion request from Γ to \mathbf{D} . Then (M_1, M_2) is a $\langle \Upsilon_{\mathfrak{C}_u}, \uparrow \rangle$ -optimal realization of u iff M_2 is a tuple-minimal Armstrong model of $\text{LeastRefl}\langle u, \Upsilon_{\mathfrak{C}_u} \rangle$ with respect to $\Upsilon_{\mathfrak{C}_u}$.*

PROOF: The proof follows immediately from 4.7 and 4.10. \square

4.12 Example — Dependence upon the instance of the main schema There is another issue which arises in applying 4.7 and 4.11; namely, that whether or not an optimal insertion exists may depend upon M_1 and not simply the view instance $\gamma(M_1)$. To illustrate this phenomenon, consider the schema \mathbf{E}_5 containing the single relation symbol $R[ABCDE]$ governed by the FDs in $\mathcal{F} = \{A \rightarrow D, B \rightarrow E, DE \rightarrow C\}$. The view to be updated is the projection onto ABC ; it contains the single relation symbol $R_{ABC}[ABC]$, and

is represented more formally as $\Pi_{ABC}^{\mathbf{E}_5} = (\mathbf{W}_5, \pi_{ABC}^{\mathbf{E}_5})$. It is easy to see that $\text{Constr}(\mathbf{W}_5) = \{AB \rightarrow C\}$, so that the view is constrained by FDs alone. Each of the two instances $M_{20} = \{R(a_0, b_0, c_0, d_0, e_0), R(a_1, b_1, c_1, d_1, e_1)\}$, and $M_{21} = \{R(a_0, b_0, c_0, d_0, e_0), R(a_1, b_1, c_1, d_1, e_0)\}$ of \mathbf{E}_5 maps to the view instance $N_1 = \{R_{ABC}(a_0, b_0, c_0), R_{ABC}(a_1, b_1, c_1)\}$ under $\pi_{ABC}^{\mathbf{E}_5}$. Consider the view update which inserts the tuple $R_{ABC}(a_0, b_1, c_2)$, so that the desired new view instance is $N_2 = N_1 \cup \{R_{ABC}(a_0, b_1, c_2)\}$. For convenience, write $u_{20} = (M_{20}, N_2)$ and $u_{21} = (M_{21}, N_2)$. Then $\text{LeastRefl}\langle u_{20}, \Upsilon_{\mathbf{e}_u} \rangle$ is given by $M'_{20} = M_{20} \cup R(a_0, b_1, c_2, d_0, e_1)$, and it is easy to see that $\gamma(M'_{20}) = N_2$, as desired. On the other hand, $\text{LeastRefl}\langle u_{21}, \Upsilon_{\mathbf{e}_u} \rangle$ does not exist. Indeed, the FDs stipulate that it would need to be $M'_{21} = M_{21} \cup R(a_0, b_1, c_2, d_0, e_0)$, but the presence of both $R(a_0, b_0, c_0, d_0, e_0)$ and $R(a_0, b_1, c_2, d_0, e_0)$ in the same instance violates the FD $DE \rightarrow C$. In the case that the instance of \mathbf{E}_5 is M'_1 , there is no insertion which will realize the insertion of $R_{ABC}(a_0, b_1, c_2)$ into $\gamma(M'_1) = N_1$. To realize this update, tuples of M'_1 must either be deleted or else altered.

This example thus shows that the principles of *reflection of monotonicity* and *invariance of admissibility*, as defined in Sec. 1, cannot always be realized simultaneously, even when the reflections for certain databases of the main schema are very well behaved. Fortunately, this phenomenon can be ruled out by requiring that Γ reflect insertions, as described below.

4.13 Reflection of insertions The view $\Gamma = (\mathbf{V}, \gamma)$ *reflects insertions* if every insertion request (M_1, N_2) from Γ to \mathbf{D} which admits a realization admits one which is itself a insertion.

In the example of 4.12 above, $\Pi_{ABC}^{\mathbf{E}_5}$ does not reflect insertions, since there is no reflection of the insertion request (M'_1, N_2) which is itself an insertion.

4.14 Observation *If Γ reflects insertions, then for every insertion request $u = (M_1, N_2)$ from Γ to \mathbf{D} , $\text{LDB}(\mathbf{D}) \cap \text{AtMod}_j(\text{LeastRefl}\langle u, \Upsilon_{\mathbf{e}_u} \rangle)$ is nonempty.*

PROOF: Since $\text{LeastRefl}\langle u, \Upsilon_{\mathbf{e}_u} \rangle$ is the least information which any realization of u which is an insertion must contain, it must be consistent if any insertion has that property. \square

4.15 Strong monotonicity of views In view of 4.11 and 4.14, it is clear that for the reflection of updates to be well behaved, a view should reflect both deletions and insertions. Because of the importance of this property, it is given a special name. Call Γ *strongly monotonic* if it reflects both insertions and deletions.

Finally, conditions which guarantee the existence of optimal insertions may be established.

4.16 Theorem *If \mathbf{D} admits canonical models unconditionally with respect to Υ_K and Γ is strongly monotonic, then every insertion request u from Γ to \mathbf{D} admits a $\langle \Upsilon_{\mathbf{e}_u}, \uparrow \rangle$ -optimal realization.*

PROOF: The proof follows from 4.11 and 4.14. \square

4.17 Corollary *If $\text{Constr}(\mathbf{D})$ is finite and consists of a weakly acyclic set of TGHDS, together with any set of EGHDs, and Γ is strongly monotonic, then every insertion request u from Γ to \mathbf{D} has a $\langle \Upsilon_{\mathbf{e}_u}, \uparrow \rangle$ -optimal realization.*

PROOF: The proof follows from 4.16 and 3.22. \square

5 Characterization of Strongly Monotonic Views

The characterizations 4.16 and 4.17 lead to the further problem of identifying conditions under which a view is strongly monotonic. In general, this does not appear to be an easy question to answer. However, for schemata constrained by FDs and unary inclusion dependencies (UINDs), and for views defined by projections, it is possible to identify some sufficient conditions which are easily verified in practice.

5.1 Partial dependence and complete sets In the example of 4.12, the problem is that there is a sort of weak dependence of A upon C via $A \rightarrow D; D \subseteq DE; DE \rightarrow C$, while the Fd $A \rightarrow C$ itself does not hold. To obtain strong monotonicity in the context of projections of views constrained by FDs, it is precisely this sort of weak dependence which must not be present without the associated FD also holding. To formalize this idea for a general schema \mathbf{D} which is constrained by FDs, let $R \in \text{Rels}(\mathbf{D})$, let \mathcal{F}_R be a set of FDs on R , and let $A, B \in \text{Ar}_{\mathbf{D}}(R)$.

- (a) A *functionally influences* B , denoted $A \dashrightarrow B$, if there is a sequence $\langle A_0, A_1, A_2, \dots, A_n \rangle$ of elements of $\text{Ar}_{\mathbf{D}}(R)$ with $A = A_0$ and $B = A_k$, and a sequence $X_1 \rightarrow A_1, X_2 \rightarrow A_2, \dots, X_k \rightarrow A_k$ of FDs in the closure of \mathcal{F}_R with the property that $A_i \in X_{i+1}$ for $i \in \{0, \dots, k\}$. This may be visualized as follows.

$$A = A_0 \in X_1; X_1 \rightarrow A_1; A_1 \in X_2; X_2 \rightarrow A_2; \dots A_{k-1} \in X_k; X_k \rightarrow A_k = B$$

Functional influence is weaker than functional dependence. $A \dashrightarrow B$ simply means that the value of A could influence the value of B , subject to information about the values of other attributes. Put another way, if $A \dashrightarrow B$ does not hold, then the value of A cannot influence the value of B via the FDs which hold on the schema.

- (b) Call a subset $Y \subseteq \text{Ar}_{\mathbf{D}}(R)$ *complete for* \mathcal{F}_R if whenever $A, B \in Y$ with $A \dashrightarrow B$, then there is a $Z \subseteq Y$ with $A \in Z$ and $Z \rightarrow B \in \text{Closure}(\mathcal{F}_R, \text{WFS}(\mathbf{D}, \exists \wedge +))$.

In other words, completeness states that if $A \dashrightarrow B$ holds in Y , then an FD whose left hand side contains A and whose right-hand side is B also embeds into Y .

5.2 Simple projective views and complete views Informally, $\Gamma = (\mathbf{V}, \gamma)$ is a *simple projective view* of \mathbf{D} if the schema \mathbf{V} consists of at most one projection of each relation symbol of \mathbf{D} . Formally, a simple projective view $\Gamma = (\mathbf{V}, \gamma)$ of \mathbf{D} is defined by an injective function $\text{SP}_{\Gamma} : \text{Rels}(\mathbf{V}) \rightarrow \text{Rels}(\mathbf{D})$ with $\text{Ar}_{\mathbf{V}}(R) \subseteq \text{Ar}_{\mathbf{D}}\text{SP}_{\Gamma}(R)$ for each $R \in \text{Rels}(\mathbf{V})$. $\text{SP}_{\Gamma}(R)$ is the relation of which R is a projection. The property that SP_{Γ} be injective; that is, that each relation of \mathbf{V} be the projection of a distinct relation in \mathbf{D} , is critical.

Now assume that Γ is a simple projective view and that each $S \in \text{Rels}(\mathbf{D})$ is constrained by a set \mathcal{F}_S of FDs. Call Γ *FD-complete* if for each $R \in \text{Rels}(\mathbf{V})$, the set $\text{Ar}_{\mathbf{V}}(R)$ is complete for $\mathcal{F}_{\text{SP}_{\Gamma}(R)}$. In other words, for each relation symbol of \mathbf{V} , the projected attributes must be complete in the relation of \mathbf{D} from which they originate.

5.3 Proposition *Suppose that \mathbf{D} is constrained solely by FDs, and that Γ is a simple projective view which is FD-complete. Then Γ is strongly monotonic.*

PROOF: Since the relation symbols of \mathbf{D} are independent of one another, it suffices to consider the situation in which $\mathbf{Rels}(\mathbf{D})$ consists of single relation symbol $R[X]$, constrained by FDs \mathcal{F}_R , with $\mathbf{Rels}(\mathbf{V})$ consisting of a single relation symbol $R'[Y]$ with $Y \subseteq X$, and $\gamma = \pi_Y^{\mathbf{D}}$ defining the projection of $R[X]$ onto $R'[Y]$.

First of all, since FDs are EGDS, deletions are reflected trivially. Now let $u = (M_1, N_2)$ be an insertion request, and let $P = N_2 \setminus \gamma(M_1)$. Each tuple $t' = R'(a_1, a_2, \dots, a_m) \in P$ must be lifted to a tuple $t = R(a_1, a_2, \dots, a_m, b_1, b_2, \dots, b_k)$ in \mathbf{D} , with the b_i 's values for the attributes in $X \setminus Y$. In view of the FD-completeness property of Γ , the values of the b_i 's cannot be influenced by the values of the a_i 's. (Put another way, the classical chase procedure [BV84] will not force the values of the b_i 's to match those of any existing tuples.) Thus the insertion will not have the problems which are illustrated in 4.12, and so the update may be realized as an insertion. \square

5.4 Examples — FDs and strong monotonicity As noted above, the example of 4.12 is not strongly monotonic. However, if the set $\mathcal{F} = \{A \rightarrow D, B \rightarrow E, DE \rightarrow C\}$ of FDs is replaced by $\mathcal{F}' = \{A \rightarrow D, B \rightarrow E, D \rightarrow C\}$, then $A \rightarrow C$ is in the closure of \mathcal{F}' , ABC is FD-complete, and the associated projection is strongly monotonic.

5.5 Examples — UINDs and strong monotonicity It is possible to obtain conditions under which simple projective views governed by UINDs are strongly monotonic, although some care is necessary. A few examples will illustrate the central issues.

First of all, let \mathbf{E}_6 denote the schema with two binary relation symbols $R_1[AB]$ and $R_2[AB]$, governed by the IND $R_1[AB] \sqsubseteq R_2[AB]$. $R_2[AB]$ is also governed by the FD $A \rightarrow B$, and this FD is inherited by $R_1[AB]$ via the IND. The view $\Pi_{R_1}^{\mathbf{E}_6}$ preserves $R_1[AB]$ but discards $R_2[AB]$. Let $M_1 = \{R_2(a_0, b_0)\}$ be the current instance of \mathbf{E}_6 , so that the instance of the view is \emptyset . Consider the update request (M_1, N_2) with $N_2 = \{R_1(a_0, b_0)\}$. This update can be realized as the insertion of $\{R_1(a_0, b_0)\}$ to R in M_1 . However, if the instance of \mathbf{E}_6 were $M'_1 = \{R_2(a_0, b_1)\}$ instead, this view insertion would not be realizable as an insertion to the main schema, since the insertion of $R_1(a_0, b_0)$ requires the insertion of $R_2(a_0, b_0)$, and $\{R_2(a_0, b_0), R_2(a_0, b_1)\}$ together violate the FD $A \rightarrow B$ on R_2 . Thus, any reflection of (M'_1, N_2) must delete $R_2(a_0, b_1)$. Hence $\Pi_{R_1}^{\mathbf{E}_6}$ does not reflect insertions. In general, non-unary INDs are very problematic with respect to strong monotonicity, and so in this paper attention will be restricted to UINDs. Note that if the IND above is changed to $R_1[A] \sqsubseteq R_2[A]$, then the problem disappears.

Next, let \mathbf{E}_7 denote the schema with the single relation symbol $R[ABC]$, constrained by $R[B] \sqsubseteq R[C]$, and let $\Pi_A^{\mathbf{E}_7} = (\mathbf{W}_7, \pi_A^{\mathbf{E}_7})$ be the view which projects $R[ABC]$ onto $R'[A]$. Let $M_1 = \{R(a_0, b_0, b_0)R(a_1, b_1, b_1)\}$, so that $\pi_A^{\mathbf{E}_7}(M_1) = N_1 = \{R'(a_0), R'(a_1)\}$. The view instance $N_2 = N_1 \setminus R'(a_1) = \{R'(a_0)\}$ may be realized via the deletion of $R(a_1, b_1, b_1)$ from M_1 . On the other hand, for $M'_1 = \{R(a_0, b_0, b_1)R(a_1, b_1, b_0)\}$, $\pi_A^{\mathbf{E}_7}(M'_1) = N_1$ as well, yet there is no subset of M'_1 which maps to N_2 under $\pi_A^{\mathbf{E}_7}$. Hence $\Pi_A^{\mathbf{E}_7}$ does not reflect deletions. A similar problem occurs if the view is taken to be the projection $\Pi_{AB}^{\mathbf{E}_7}$ onto AB or the projection $\Pi_{AC}^{\mathbf{E}_7}$ onto AC . Thus, the second principle to enforce is that for every intrarelatational UIND $R[A_1] \sqsubseteq R[A_2]$, if R is projected at all to the view, then the projection must contain $R[A_1A_2]$. Less formally, all intrarelatational UINDs must be completely visible in the projection.

Finally, let \mathbf{E}_8 be the schema with three relation symbols $R_1[AB]$, $R_2[AB]$, and $R_3[AB]$, constrained by $R_1[A] \sqsubseteq R_2[A]$ and $R_2[B] \sqsubseteq R_3[B]$. The simple projective view Ω_7 contains the pro-

jections $R'_1[A]$ and $R'_3[A]$ of R_1 and R_3 , respectively. Let $M_1 = \{R_1(a_0, b_0), R_2(a_0, b_1), R_3(a_0, b_1), R_3(a_1, b_0)\}$, so that the view instance is $N_1 = \{R'_1(a_0), R'_3(a_0), R'_3(a_1)\}$. Let $N_2 = \{R'_1(a_0), R'_3(a_0)\}$. It is easy to see that the deletion request (M_1, N_2) is realizable by deleting $R_3(a_1, b_0)$ from M_1 . However, for $M'_1 = \{R_1(a_0, b_0), R_2(a_0, b_1), R_3(a_1, b_1), R_3(a_0, b_0)\}$, the deletion request (M'_1, N_2) is not realizable as a deletion, since $R_3(a_1, b_1)$ must be deleted, which would violate $R_2[B] \sqsubseteq R_3[B]$. Note that parts of the “chain” $R_1[A] \sqsubseteq R_2[A] \leftrightarrow R_2[B] \sqsubseteq R_3[B] \leftrightarrow R_3[A]$ are not visible in the view, with \leftrightarrow meaning “occurs in the same relation. The final condition to be enforced is that all intermediate entries in such a chain must appear in the view whenever the end points do.

The next task is to formalize all of this.

5.6 UINDs and projective views The set of *unary projections* of \mathbf{D} consists of all expressions of the form $R[A]$ with $R \in \text{Rels}(\mathbf{D})$ and $A \in \text{Ar}_{\mathbf{D}}(R)$. Thus, the unary projections are precisely those which can occur as the left-hand or right-hand side of a UIND. Formally, $\text{UProj}(\mathbf{D}) = \{R[A] \mid R \in \text{Rels}(\mathbf{D}) \text{ and } A \in \text{Ar}_{\mathbf{D}}(R)\}$.

Define $\text{UIND}(\mathbf{D})$ to be the set of all UINDs which are implied by $\text{Constr}(\mathbf{D})$. Say that $R[A]$ *participates* in $\text{UIND}(\mathbf{D})$ if it appears as either the left-hand side or else the right-hand side of some nontrivial $\varphi \in \text{UIND}(\mathbf{D})$. Here a *nontrivial* UIND is one which is not true in $M \in \text{DB}(\mathbf{D})$; *i.e.*, one which is not of the form $R[A] \sqsubseteq R[A]$. For $R[A_1], R[A_2] \in \text{UProj}(\mathbf{D})$ both participants in $\text{UIND}(\mathbf{D})$ and over the same relation R , write $R[A_1] \leftrightarrow R[A_2]$. Note that $R[A_1] \leftrightarrow R[A_2]$ does not necessarily imply that one of $R[A_1] \sqsubseteq R[A_2]$ or $R[A_2] \sqsubseteq R[A_1]$ holds; $R[A_1]$ and $R[A_2]$ may well participate in distinct UINDs of \mathbf{D} .

Define the *UIND-graph* of \mathbf{D} , denoted $\text{UGraph}(\mathbf{D})$, to be the directed graph whose vertices are the members of $\text{UProj}(\mathbf{D})$, with an edge from $R_1[A_1]$ to $R_2[A_2]$ iff $R_1[A_1] \neq R_2[A_2]$ and either $R_1[A_1] \sqsubseteq R_2[A_2] \in \text{UIND}(\mathbf{D})$ or else $R_1 = R_2$ and $R_1[A_1] \leftrightarrow R_2[A_2]$. Thus, the UIND-graph recaptures “chains of influence for UINDs, much as functional influence does for FDs.

Finally, assume that $\Gamma = (\mathbf{V}, \gamma)$ is a simple projective view of \mathbf{D} . Say that $R[A] \in \text{UProj}(\mathbf{D})$ is *visible in* Γ if there is some $R' \in \text{Rels}(\mathbf{V})$ which is a projection of R under γ . Call Γ *UIND-complete* if for every directed path $\rho = \langle R_1[A_1], R_2[A_2], \dots, R_k[A_k] \rangle$ in $\text{UGraph}(\mathbf{D})$, if $R_1[A_1]$ and $R_k[A_k]$ are visible in Γ , so too are all intermediate elements in ρ .

The following result, analogous to 5.3, may now be established.

5.7 Proposition *Suppose that \mathbf{D} is constrained by UINDs, and that Γ is a simple projective view which is UIND-complete. Then Γ is strongly monotonic.*

PROOF: The proof is a similar to that of 5.3. The central idea is to observe that all connections between values which are forced by the UINDs are already visible in the view, in the sense that they are deducible from the constraints of the view instance, and so it is completely decidable within the view whether or not an insertion or deletion will violate a UIND. The details are left to the reader. \square

5.8 Proposition *Suppose that \mathbf{D} is constrained by FDs and UINDs, and that Γ is a simple projective view which is both FD-complete and UIND-complete. Then Γ is strongly monotonic.*

PROOF: The proof follows from 5.3 and 5.7, together with the classical result that FDs and UINDs have trivial interaction [CKV90]. \square

5.9 Theorem *Suppose that \mathbf{D} is constrained by FDs and weakly acyclic UINDs, and that Γ is a simple projective view which is both FD-complete and UIND-complete. Then every insertion request u from Γ to \mathbf{D} admits a $\langle \Upsilon_{\mathfrak{c}_u}, \uparrow \rangle$ -optimal realization.*

PROOF: Combine 3.22, 4.17, and 5.8. \square

6 Optimal Reflection of Deletions

At first glance, the information-based modelling of deletions to views would appear to be much simpler than that for insertions. Indeed, in large part, the only relevant information-monotone family is $\text{GrAtoms}(\mathbf{D})$. Thus, the following is immediate.

6.1 Observation *Let $u = (M_1, N_2)$ be a deletion request from Γ to \mathbf{D} , $(M_1, M_2) \in \text{DelRealiz}\langle (M_1, N_2), \Gamma \rangle$, and let $K = \text{ConstSym}(u)$.*

- (a) *(M_1, M_2) is $\langle \Upsilon_K, \downarrow \rangle$ -admissible iff for all $(M_1, M'_2) \in \text{DelRealiz}\langle (M_1, N_2), \Gamma \rangle$ with $M_2 \subseteq M'_2$, it must be that $M_2 = M'_2$.*
- (b) *(M_1, M_2) is $\langle \Upsilon_K, \downarrow \rangle$ -optimal iff for all $(M_1, M'_2) \in \text{DelRealiz}\langle (M_1, N_2), \Gamma \rangle$, $M'_2 \subseteq M_2$. \square*

6.2 Admissibility and optimality for deletions As the admissibility and optimality of deletions do not depend upon the set of constant symbols in the instances defining the update, a more concise notation may be employed. For $u = (M_1, N_2)$ be a deletion request from Γ to \mathbf{D} and $(M_1, M_2) \in \text{DelRealiz}\langle (M_1, N_2), \Gamma \rangle$, write \downarrow -admissible as an abbreviation for $\langle \Upsilon_{\mathfrak{c}_u}, \downarrow \rangle$ -admissible, and \downarrow -optimal as an abbreviation for $\langle \Upsilon_{\mathfrak{c}_u}, \downarrow \rangle$ -optimal.

It is also worth observing that equality-generating dependencies and mutual-exclusion dependencies are always preserved under deletion, so no special handling of them is required.

6.3 Observation — EGHDs and mutual-exclusion TGHDs preserved under deletion *Let φ be either an EGHD or a mutual-exclusion TGHD on \mathbf{D} . If $M \in \text{AtMod}_\mathcal{J}(\varphi)$, then $M' \in \text{AtMod}_\mathcal{J}(\varphi)$ for every $M' \subseteq M$. \square*

Despite these simplifications, the optimal support of deletions is far from trivial. The root of the problem is that while TGDs are well suited for insertions, they display inherent disjunction in the context of deletion. A simple example will help illustrate.

6.4 Example — Strong monotonicity does not ensure \downarrow -optimal realizations Let \mathbf{E}_9 be the relational schema with three relation symbols $R[A]$, $S[A]$, and $T[A]$, constrained by the single TGHD $(\forall x)((R(x) \wedge S(x)) \Rightarrow T(x))$, let \mathbf{G}_9 be the schema whose single relation symbol is $T'[A]$, and let $\Omega_9 = (\mathbf{G}_9, \omega_9)$ be the view of \mathbf{G}_9 with $\omega_9^{T'} = T(x_A)$. In words, the view Ω_9 preserves $T[A]$ (as $T'[A]$) but discards $R[A]$ and $S[A]$ completely.

Clearly Ω_9 is strongly monotonic. However, it does not always admit reflections which are \downarrow -optimal. Indeed, let $M_1 = \{R(a_0), R(a_1), S(a_0), T(a_0)\}$ be the current instance of \mathbf{E}_9 , so

that $N_1 = \{T'(a_0)\}$ is the instance of \mathbf{G}_9 . Let $N_2 = \emptyset$. For $M'_2 = \{R(a_0), R(a_1)\}$ and $M''_2 = \{S(a_0), R(a_1)\}$, each of (M_1, M'_2) and (M_1, M''_2) are \downarrow -admissible realizations of $u = (M_1, N_2)$ with respect to $\mathbf{WFS}(\mathbf{E}_9, \exists\wedge+)$, and so neither is \downarrow -optimal.

Observe that for $M_3 = M'_2 \cap M''_2 = \{R(a_1)\}$, $(M_1, M_3) \in \text{DelRealiz}\langle(M_1, N_2), \Omega_9\rangle$ as well, although it is not minimal. This leads to a weaker “minimax”-style of optimality, in which every tuple which is deleted in some \downarrow -admissible realization is deleted. This is formalized as follows.

6.5 Weak \downarrow -optimality Let $u = (M_1, N_2)$ be a deletion request from Γ to \mathbf{D} .

- (a) Define $\text{WeakOpt}_{\downarrow}\langle u, \Gamma \rangle = \bigcap \{M_3 \mid (M_1, M_3) \in \text{DelRealiz}\langle(M_1, N_2), \Gamma\rangle \text{ and } \downarrow\text{-admissible}\}$.
- (b) Call $(M_1, M_2) \in \text{DelRealiz}\langle(M_1, N_2), \Gamma\rangle$ *weakly \downarrow -optimal* if $M_2 = \text{WeakOpt}_{\downarrow}\langle u, \Gamma \rangle$.

In 6.4 above, M_3 is weakly \downarrow -optimal but not \downarrow -optimal in the sense of 6.2.

To identify conditions under which a weakly \downarrow -optimal reflections are admitted, it useful to introduce a new way of viewing the combination of a schema and a view.

6.6 The combined schema induced by a view Rather than regarding the main schema \mathbf{D} and the view schema \mathbf{V} as distinct, it is quite possible to combine them into a single schema, with the view mappings regarded as additional constraints. The alternative representation turns out to be very useful, since all constraints, both those of the schemata and those induced by the view mappings, may be considered at once.

To formalize this idea, it is necessary to assume that $\text{Rels}(\mathbf{D}) \cap \text{Rels}(\mathbf{V}) = \emptyset$. This is not a problem since relations may always be renamed as necessary.

The *combined schema* $\text{CombSch}\langle \mathbf{D}, \Gamma \rangle$ has as its relational symbols $\text{Rels}(\mathbf{D}) \cup \text{Rels}(\mathbf{V})$. The constraints of $\text{CombSch}\langle \mathbf{D}, \Gamma \rangle$ are those in $\text{Constr}(\mathbf{D})$, together with, for each $R \in \text{Rels}(\mathbf{V})$, the *definitional constraint*

$$\text{(DefC)} \quad (\forall x_{A_1})(\forall x_{A_2}) \dots (\forall x_{A_m}) ((R(x_{A_1}, x_{A_2}, \dots, x_{A_n})) \Leftrightarrow \gamma^R)$$

In the above, $\{x_{A_1}, x_{A_2}, \dots, x_{A_m}\}$ are precisely the attribute variables which occur in the interpretation formula γ^R . It is easier to see the full nature of this constraint when γ^R is expanded into its full form $(\exists x_1)(\exists x_2) \dots (\exists x_n)(\bar{\gamma}^R(x_{A_1}, x_{A_2}, \dots, x_{A_m}, x_1, x_2, \dots, x_n))$. (See 2.6 for a clarification of the notation $\bar{\gamma}^R$.) The complete expansion then becomes

$$\text{(DefC')} \quad (\forall x_{A_1})(\forall x_{A_2}) \dots (\forall x_{A_m}) ((R(x_{A_1}, x_{A_2}, \dots, x_{A_n})) \Leftrightarrow (\exists x_1)(\exists x_2) \dots (\exists x_n)(\bar{\gamma}^R)(x_{A_1}, x_{A_2}, \dots, x_{A_m}, x_1, x_2, \dots, x_n))$$

The definitional constraint (DefC') for R may be broken into the *forward component*

$$\text{(DefC-Fwd)} \quad (\forall x_{A_1})(\forall x_{A_2}) \dots (\forall x_{A_m})(\forall x_1)(\forall x_2) \dots (\forall x_n) (\bar{\gamma}^R(x_{A_1}, x_{A_2}, \dots, x_{A_m}, x_1, x_2, \dots, x_n) \Rightarrow (R(x_1, x_2, \dots, x_n)))$$

and the *reverse component*

$$\text{(DefC-rev)} \quad (\forall x_{A_1})(\forall x_{A_2}) \dots (\forall x_{A_m}) ((R(x_{A_1}, x_{A_2}, \dots, x_{A_n})) \Rightarrow (\exists x_1)(\exists x_2) \dots (\exists x_n)(\bar{\gamma}^R)(x_{A_1}, x_{A_2}, \dots, x_{A_m}, x_1, x_2, \dots, x_n))$$

It is easy to see that both the forward and the reverse components are TGHs. Thus, if $\text{Constr}(\mathbf{D})$ has a cover consisting of GHs, so too does $\text{Constr}(\text{CombSch}\langle\mathbf{D}, \Gamma\rangle)$.

This idea has already been illustrated in 1.2 with \mathbf{E}'_0 the combined schema associated with $\text{CombSch}\langle\mathbf{E}_0, \Pi_{AB}^{\mathbf{E}_0}\rangle$. The single definitional constraint is $(\forall x)(\forall y)(R_{AB}(x, y) \Leftrightarrow (\exists z)(R(x, y, z)))$, and this decomposes into the forward constraint $(\forall x)(\forall y)(\forall z)(R(x, y, z) \Rightarrow R_{AB}(x, y))$ and the reverse constraint $(\forall x)(\forall y)(R_{AB}(x, y) \Rightarrow (\exists z)(R(x, y, z)))$. Note that the forward constraint is always universal.

6.7 Universal pairs Call the pair $\langle\mathbf{D}, \Gamma\rangle$ *universal* if $\text{Constr}(\text{CombSch}\langle\mathbf{D}, \Gamma\rangle)$ consists entirely of universal GHs. In other words, this means that both the constraints of \mathbf{D} and the view interpretation mappings consist of total dependencies, without any existential quantification. For example, in 6.4, the pair $\langle\mathbf{E}_9, \Omega_9\rangle$ is universal.

Under this assumption of universality, weakly \downarrow -optimal solutions exist whenever a solution which is a deletion is possible.

6.8 Proposition *Let $\langle\mathbf{D}, \Gamma\rangle$ be a universal pair. Then for every deletion request $u = (M_1, N_2)$ from Γ to \mathbf{D} with $\text{DelRealiz}\langle u, \Gamma\rangle \neq \emptyset$, $\text{WeakOpt}_{\downarrow}\langle u, \Gamma\rangle$ is a weak \downarrow -optimal realization of u .*

PROOF: Let $(M_1, M_2) \in \text{DelRealiz}\langle (M_1, N_2), \Gamma\rangle$ and let $t \in N_2$. Since $\langle\mathbf{D}, \Gamma\rangle$ is a universal pair, $\text{Subst}(\gamma, t \rightarrow)$ must be a conjunction of ground atoms, and these ground atoms must be in every M for which $t \in \gamma(M)$. In particular, each such conjunct must be in $\bigcap\{M_3 \mid (M_1, M_3) \in \text{DelRealiz}\langle (M_1, N_2), \Gamma\rangle \text{ and } \downarrow\text{-admissible}\}$. Thus $\gamma(\text{WeakOpt}_{\downarrow}\langle u, \Gamma\rangle) = N_2$.

It remains to verify that $\text{WeakOpt}_{\downarrow}\langle u, \Gamma\rangle \in \text{LDB}(\mathbf{D})$. However, it is a very easy exercise to show that all universal Horn sentences (and hence all total TGDs) are preserved under intersection, whence the result. \square

6.9 Example — Lack of weak \downarrow -optimal realizations It is natural to conjecture that the result of 6.8 extends to situations involving existential quantification. Unfortunately, this is not the case. If non-total TGDs are allowed, relatively simple examples of deletions requests exist which admit \downarrow -admissible realizations but no weak \downarrow -optimal realization.

Let \mathbf{E}_{10} be the schema with three relational symbols $R[AB]$, $S[BC]$, and $T[BC]$, with the following three constraints.

$$\begin{aligned} &(\forall x_1)(\forall x_2)(\forall x_3)((R(x_1, x_2) \wedge R(x_1, x_3)) \Rightarrow S(x_2, x_3)) \\ &(\forall x_1)(\forall x_2)(R(x_1, x_2) \Rightarrow T(x_2, x_2)) \\ &(\forall x_1)(T(x_1, x_1) \Rightarrow (\exists y_1)(R(y_1, x_1))) \end{aligned}$$

Let $\Omega_{10} = (\mathbf{G}_{10}, \omega_{10})$ be the view which retains the relations of S and T , but discards R . Let $M_1 = \{R(a_0, b_0), R(a_0, b_1), R(a_1, b_0), R(a_1, b_1), S(b_0, b_0), S(b_0, b_1), S(b_1, b_0), S(b_1, b_1), T(b_0, b_0), T(b_1, b_1)\}$. It is easy to see that $M_1 \in \text{LDB}(\mathbf{E}_{10})$. The corresponding view instance $\omega_{10}(M_1) = \{S(b_0, b_0), S(b_0, b_1), S(b_1, b_0), S(b_1, b_1), T(b_0, b_0), T(b_1, b_1)\}$. Let the desired new view instance be $N_2 = \{S(b_0, b_0), S(b_1, b_1), T(b_0, b_0), T(b_1, b_1)\}$. Thus, the tuples in $P = \{S(b_1, b_0), S(b_0, b_1)\}$ are to be deleted. It is easy to see that this update admits two \downarrow -admissible realizations, one which deletes $P \cup \{R(a_0, b_0), R(a_1, b_1)\}$ and the other which deletes

$P \cup \{R(a_0, b_1), R(a_1, b_0)\}$. Thus, the two possibilities for the new instance of \mathbf{E}_{10} are $M'_2 = \{R(a_0, b_1), R(a_1, b_0), S(b_0, b_0), S(b_1, b_1), T(b_0, b_0), T(b_1, b_1)\}$ and $M''_2 = \{R(a_0, b_0), R(a_1, b_1), S(b_0, b_0), S(b_1, b_1), T(b_0, b_0), T(b_1, b_1)\}$. However, $M'_2 \cap M''_2 = \{S(b_0, b_0), S(b_1, b_1), T(b_0, b_0), T(b_1, b_1)\} \notin \text{LDB}(\mathbf{E}_{10})$. Indeed, the existential quantification has introduced an exclusive-or requirement in which $(R(a_0, b_0) \wedge R(a_1, b_1)) \vee (R(a_0, b_1) \wedge R(a_1, b_0))$ must hold, but both disjuncts must not hold simultaneously. Furthermore, it is not difficult to see that Ω_{10} reflects deletions. Therefore, a general approach which addresses the disjunction problem for deletions seems impossible.

It should be noted that $\text{Constr}(\mathbf{E}_{10})$ is not typed in the sense of [Fag82, p. 955]. For this paper, whether such an example is possible with such typed constraints is left as an open question. Rather, attention is turned to a more restrictive but nevertheless useful class which does admit full \downarrow -optimal realizations.

6.10 Unit-head TGDs, schemata, and views A TGHD of the form (GHD) of 3.15 is called a *unit-head dependency* if $n = 1$; that is, if there is only one atom on the left-hand side of the rule. The most important example of a unit-head dependency is the inclusion dependency. The schema \mathbf{D} is called *unit-head* if $\text{Constr}(\mathbf{D})$ has a cover consisting of EGDs, mutual-exclusion TGHDs, and unit-head TGHDs.

This idea extends to combined schemata as well. Call the pair $\langle \mathbf{D}, \Gamma \rangle$ *unit head* if $\text{CombSch}\langle \mathbf{D}, \Gamma \rangle$ is unit head. Clearly, $\text{CombSch}\langle \mathbf{D}, \Gamma \rangle$ is unit head iff \mathbf{D} has that property and, for each $R \in \text{Rels}(\mathbf{V})$, the interpretation formula $\bar{\gamma}^R$ consists of a single (not necessarily ground) atom.

All INDs, are unit-head. Furthermore, views which are defined via projection and selection (but not join), are also unit head. Therefore, many practical examples are represented under this class.

The following lemma identifies a critical property of unit-head schemata.

6.11 Lemma *Let \mathbf{D} be a unit-head schema, and let $M, M_1, M_2 \in \text{LDB}(\mathbf{D})$ with $M_i \subseteq M$ for $i \in \{1, 2\}$. Then $M_1 \cup M_2 \in \text{LDB}(\mathbf{D})$ as well.*

PROOF: Without loss of generality, assume that $\text{Constr}(\mathbf{D})$ itself consists of EGDs and unit-head TGHDs. Let $\varphi \in \text{Constr}(\mathbf{D})$ be a TGHD and let s be a constant substitution into the universal variables of φ such that $\text{LHSinfo}\langle \varphi, s \rangle \in M_1 \cup M_2$. Let $i \in \{1, 2\}$ for which $\text{LHSinfo}\langle \varphi, s \rangle \in M_i$. There must then be a tuple $t \in M_i$ which satisfies $\text{RHSinfo}\langle \varphi, s \rangle$, since $M_i \in \text{LDB}(\mathbf{D})$. Thus $M_1 \cup M_2 \in \text{AtMod}_j(\varphi)$, and so $M_1 \cup M_2$ satisfies every TGHD in $\text{Constr}(\mathbf{D})$. Since $M_1 \cup M_2 \subseteq M$, it follows from 6.3 that $M_1 \cup M_2 \in \text{AtMod}_j(\psi)$ for every EGD and mutual-exclusion TGHD in $\text{Constr}(\mathbf{D})$. Thus, $M_1 \cup M_2 \in \text{LDB}(\mathbf{D})$. \square

6.12 Theorem *Let $\langle \mathbf{D}, \Gamma \rangle$ be a unit-head pair. Then every deletion request $u = (M_1, N_2)$ from Γ to \mathbf{D} for which $\text{DelRealiz}\langle (M_1, N_2), \Gamma \rangle \neq \emptyset$ admits a unique \downarrow -optimal realization.*

PROOF: Given a deletion request $u = (M', N_2)$ from Γ to \mathbf{D} , let M'_1 and M'_2 be \downarrow -admissible realizations of u , and let $M_i = M'_i \cup N_2 \in \text{LDB}(\text{CombSch}\langle \mathbf{D}, \Gamma \rangle)$ for $i \in \{1, 2\}$. Now just apply 6.11 to establish that $M_1 \cup M_2 \in \text{LDB}(\mathbf{D})$. It thus follows that $M_1 = M_2$, else one of M_1 and

M_2 would not be tuple minimal. Hence the part of M_1 which corresponds to the relations of \mathbf{D} must be the \downarrow -optimal solution. \square

In the context of the traditional dependencies and view constructions of the relational model, this theorem leads to the following corollary.

6.13 Corollary *Let \mathbf{D} be constrained by FDs and INDs, and let Γ be defined by projections and selections on the relations of \mathbf{D} . Then, if Γ reflects deletions, every deletion request from Γ to \mathbf{D} admits a unique \downarrow -optimal realization. \square*

Invoking 5.8 yields the following more focused characterization.

6.14 Corollary *Let \mathbf{D} be constrained by FDs and weakly acyclic INDs, and let Γ be a simple-projective view which is strongly monotonic. Then every unidirectional update request (i.e., insertion request or deletion request) u from Γ to \mathbf{D} admits a unique optimal realization ($\langle \Upsilon_{\mathbf{e}_u}, \uparrow \rangle$ -optimal or \downarrow -optimal, as the case may be). \square*

7 Conclusions and Further Directions

A strategy for the optimal reflection of view updates has been developed, based upon the concept of least information change. The property of strong monotonicity — that view insertions may always be reflected as main-schema insertions and view deletions may always be reflected as main-schema deletions, has been shown to be critical. Under this assumption, and in the context of generalized Horn dependencies, it has been shown that optimal insertions are supported in a reasonable fashion — they are unique up to a renaming of the newly-inserted constants. It has furthermore been shown that optimal deletions are supported under unit-head conditions. Nonetheless, a number of issues remain for future investigation. Among the most important are the following.

Deletion beyond the unit-head context The theory for deletions developed in Section 6 is largely restricted to unit-head pairs. It would be useful to extend these results to a wider class of schemata. As noted at the end of 6.9, it is not known (at least to the author) whether weak \downarrow -optimality may be obtained for typed GHDs. This topic warrants further investigation.

Characterization of strong monotonicity for wider classes of schema constraints The characterization in Section 5 of views which are strongly monotonic is limited to simple projections constrained by FDs and UINDs. Since strong monotonicity is central to the support of optimal updates, an investigation into broader characterization would certainly be worthwhile.

Optimization of tuple modification Although the general formulation applies to all types of updates, the results focus almost entirely upon insertions and deletions. Modification of single tuples (“updates” in SQL), on the other hand, are of fundamental importance. With the standard update classification pair introduced in 4.4 and used throughout the paper, only very special cases admit optimal solutions. The difficulty arises from the fact that the framework, which is based entirely upon information content, cannot distinguish between the

process of modifying a tuple and that of deleting it and then inserting a new one. Consequently, both appear as admissible updates, but neither is optimal relative to the other. Further work must therefore look for a way to recapture the distinction between tuple modification and a delete-insert pair.

Application to database components This investigation began as an effort to understand better how updates are propagated between database components, as forwarded in [Heg08b, Sec. 4], but then took on a life of its own as it was discovered that the component-based problems were in turn dependent upon more fundamental issues. Nevertheless, it is important to return to the roots of this investigation — database components. This includes not only the purely autonomous case, as sketched in [Heg08b, Sec. 4], but also the situation in which users cooperate to achieve a suitable reflection, as introduced in [HS07]

Relationship to work in logic programming As already noted in the introduction, the problem of view update has also been studied extensively in the context of deductive databases. The connection between update preference based upon distance measures, as identified in Section 1 and the current approach beg a rapprochement. In addition, the connection between the current work and that of identifying algorithms for finding all possible reflections [BM04] is of interest. Furthermore, some recent work has introduced the idea of using active constraints to establish a preference order on admissible updates [GSTZ03]. Thus, rather than employing a preference based upon information content, one based upon explicit rules is employed. The relationship between such approaches and that of this paper warrants further investigation. Also, there has been a substantial body of work on updates to disjunctive deductive databases [FGM96], in which the extensional database itself consists of a collection of alternatives. The approach of minimizing information change in the disjunctive context deserves further attention as well.

Acknowledgment Much of this research was carried out while the author was a visitor at the Information Systems Engineering Group at the University of Kiel. He is indebted to Bernhard Thalheim and the members of the group for the invitation and for the many discussions which led to this work. Also, the anonymous reviewers made numerous suggestions which (hopefully) have led to a more readable presentation. Thanks are also due to Peggy Schmidt who read the manuscript and made many valuable suggestions.

References

- [AHV95] Serge Abiteboul, Richard Hull, and Victor Vianu, *Foundations of Databases*, Addison-Wesley, 1995.
- [ADB07] Ofer Arieli, Marc Denecker, and Maurice Bruynooghe, “Distance semantics for database repair,” *Ann. Math. Artif. Intell.*, **50**(2007), pp. 389–415.
- [ADNB06] Ofer Arieli, Marc Denecker, Bert Van Nuffelen, and Maurice Bruynooghe, “Computational methods for database repair by signed formulae,” *Ann. Math. Artif. Intell.*, **46**(2006), pp. 4–37.

- [BS81] François Bancilhon and Nicolas Spyratos, “Update semantics of relational views,” *ACM Trans. Database Systems*, **6**(1981), pp. 557–575.
- [BV84] Catriel Beeri and Moshe Y. Vardi, “A proof procedure for data dependencies,” *J. Assoc. Comp. Mach.*, **31**(1984), pp. 718–741.
- [BM04] Andreas Behrend and Rainer Manthey, “Update propagation in deductive databases using soft stratification,” in: Georg Gottlob, András A. Benczúr, and János Demetrovics, eds., *Advances in Databases and Information Systems, 8th East European Conference, ADBIS 2004, Budapest, Hungary, September 22-25, 2004, Proceedings*, pp. 22–36, Volume 3255 of Lecture Notes in Computer Science, Springer-Verlag, 2004.
- [BL97] Fadila Bentayeb and Dominique Laurent, “Inversion de l’algèbre relationnelle et mises à jour,” Technical Report 97-9, Université d’Orléans, LIFO, 1997.
- [BL98] Fadila Bentayeb and Dominique Laurent, “View updates translations in relational databases,” in: *Proc. DEXA ’98, Vienna, Sept. 24-28, 1998*, pp. 322–331, 1998.
- [CGT90] Stefano Ceri, Georg Gottlob, and Letizia Tanca, *Logic Programming and Databases*, Springer-Verlag, 1990.
- [CL73] Chin-Liang Chang and Richard Char-Tung Lee, *Symbolic Logic and Mechanical Theorem Proving*, Academic Press, 1973.
- [CKV90] Stavros Cosmadakis, Paris C. Kannelakis, and Moshe Y. Vardi, “Polynomial-time implication problems for unary inclusion dependencies,” *J. Assoc. Comp. Mach.*, **37**(1990), pp. 15–46.
- [DB82] Umeshwar Dayal and Philip A. Bernstein, “On the correct translation of update operations on relational views,” *ACM Trans. Database Systems*, **8**(1982), pp. 381–416.
- [DG84] William F. Dowling and Jean H. Gallier, “Linear-time algorithms for testing the satisfiability of propositional Horn clauses,” *J. Logic Programming*, **3**(1984), pp. 267–284.
- [EM97] Thomas Eiter and Heikki Mannila, “Distance measures for point sets and their computation,” *Acta Inf.*, **34**(1997), pp. 109–133.
- [Fag82] Ronald Fagin, “Horn clauses and database dependencies,” *J. Assoc. Comp. Mach.*, **29**(1982), pp. 952–985.
- [FKMP05] Ronald Fagin, Phokion G. Kolaitis, Renée J. Miller, and Lucian Popa, “Data exchange: Semantics and query answering,” *Theoret. Comput. Sci.*, **336**(2005), pp. 89–124.
- [FV83] Ronald Fagin and Moshe Y. Vardi, “Armstrong databases for functional and inclusion dependencies,” *Info. Process. Lett.*, **16**(1983), pp. 13–19.

- [FGM96] José Alberto Fernández, John Grant, and Jack Minker, “Model theoretic approach to view updates in deductive databases,” *J. Automated Reasoning*, **17**(1996), pp. 171–197.
- [GN87] Michael R. Genesereth and Nils J. Nilsson, *Logical Foundations of Artificial Intelligence*, Morgan-Kaufmann, 1987.
- [GSTZ03] Sergio Greco, Cristina Sirangelo, Irina Trubitsyna, and Ester Zumpano, “Preferred repairs for inconsistent databases,” in: *7th International Database Engineering and Applications Symposium (IDEAS 2003), 16-18 July 2003, Hong Kong, China*, pp. 202–211, IEEE Computer Society, 2003.
- [Heg04] Stephen J. Hegner, “An order-based theory of updates for closed database views,” *Ann. Math. Art. Intell.*, **40**(2004), pp. 63–125.
- [Heg08a] Stephen J. Hegner, “Information-optimal reflections of view updates on relational database schemata,” in: Sven Hartmann and Gabriele Kern-Isberner, eds., *Foundations of Information and Knowledge Systems: Fifth International Symposium, FoIKS 2008, Pisa, Italy, February 11-15, 2008, Proceedings*, pp. 112–131, Volume 4932 of Lecture Notes in Computer Science, Springer-Verlag, 2008.
- [Heg08b] Stephen J. Hegner, “A model of database components and their interconnection based upon communicating views,” in: Hannu Jakkola, Yashui Kiyoki, and Takehiro Tokuda, eds., *Information Modelling and Knowledge Systems XIX*, pp. 79–100, Frontiers in Artificial Intelligence and Applications, IOS Press, 2008.
- [Heg08c] Stephen J. Hegner, “Semantic bijectivity and the uniqueness of constant-complement updates in the relational context,” in: Klaus-Dieter Schewe and Bernhard Thalheim, eds., *International Workshop on Semantics in Data and Knowledge Bases, SDKB 2008, Nantes, France, March 29, 2008, Proceedings*, pp. 172–191, Volume 4925 of Lecture Notes in Computer Science, Springer-Verlag, 2008.
- [HS07] Stephen J. Hegner and Peggy Schmidt, “Update support for database views via cooperation,” in: Yannis Ioannis, Boris Novikov, and Boris Rachev, eds., *Advances in Databases and Information Systems, 11th East European Conference, ADBIS 2007, Varna, Bulgaria, September 29 - October 3, 2007, Proceedings*, pp. 98–113, Volume 4690 of Lecture Notes in Computer Science, Springer-Verlag, 2007.
- [HS73] Horst Herrlich and George. E. Strecker, *Category Theory*, Allyn and Bacon, 1973.
- [Hor51] Alfred Horn, “On sentences which are true of direct unions of algebras,” *J. Symbolic Logic*, **16**(1951), pp. 14–21.
- [Hut97] Alan Hutchinson, “Metrics on terms and clauses,” in: Maarten van Someren and Gerhard Widmer, eds., *Machine Learning: ECML-97, 9th European Conference on Machine Learning, Prague, Czech Republic, April 23-25, 1997, Proceedings*, pp. 138–145, Volume 1224 of Lecture Notes in Computer Science, 1997.

- [JAK82] Barry E. Jacobs, Alan R. Aronson, and Anthony C. Klug, “On interpretations of relational languages and solutions to the implied constraint problem,” *ACM Trans. Database Systems*, **7**(1982), pp. 291–315.
- [Kel85] Arthur M. Keller, “Updating relational databases through views,” PhD thesis, Stanford University, 1985.
- [Lan90] Rom Langerak, “View updates in relational databases with an independent scheme,” *ACM Trans. Database Systems*, **15**(1990), pp. 40–66.
- [Llo87] John W. Lloyd, *Foundations of Logic Programming, Second Extended Edition*, Springer-Verlag, 1987.
- [Mak87] Johann A. Makowsky, “Why Horn formulas matter in computer science: Initial structures and generic examples,” *J. Comput. System Sci.*, **34**(1987), pp. 266–292.
- [Mon76] J. Donald Monk, *Mathematical Logic*, Springer-Verlag, 1976.
- [NC97] Shan-Hwei Nienhuys-Cheng, “Distance between Herbrand interpretations: A measure for approximations to a target concept,” in: *Inductive Logic Programming, 7th International Workshop, ILP-97, Prague, Czech Republic, September 17-20, 1997, Proceedings*, pp. 213–226, Volume 1297 of Lecture Notes in Computer Science, Springer, 1997.
- [PDGV89] Jan Paredaens, Paul De Bra, Marc Gyssens, and Dirk Van Gucht, *The Structure of the Relational Database Model*, Springer-Verlag, 1989.