

# Accuracy and Efficiency of Fixpoint Methods for Approximate Query Answering in Locally Complete Databases

Álvaro Cortés-Calabuig<sup>†</sup>, Marc Denecker<sup>†</sup>, Ofer Arieli<sup>‡</sup>, Maurice Bruynooghe<sup>†</sup>

<sup>†</sup> Department of Computer Science, Katholieke Universiteit Leuven, Belgium.  
{alvaro.cortes,marc.denecker,maurice.bruynooghe}@cs.kuleuven.be

<sup>‡</sup> Department of Computer Science, The Academic College of Tel-Aviv, Israel.  
oarieli@mta.ac.il

## Abstract

Standard databases convey Reiter’s closed-world assumption that an atom not in the database is false. This assumption is relaxed in locally closed databases that are sound but only *partially complete* about their domain. One of the consequences of the weakening of the closed-world assumption is that query answering in locally closed databases is undecidable. In this paper, we develop efficient *approximate* methods for query answering, based on fixpoint computations, and investigate conditions that assure the optimality of these methods. Our approach of approximative reasoning may be incorporated in different contexts where incompleteness plays a major role and efficient reasoning is imperative.

## Introduction

The *Closed-World Assumption* (CWA) on databases (Reiter 1982) expresses that an atom not in the database is false. However, as the following example shows, in many cases database information is only *partially* complete, and so applying the CWA is not correct, and may lead to wrong conclusions.

**Example 1** Consider the following database of a computer science department. This database stores information about the telephone numbers of the department’s members and collaborators.

Telephone		Department	
Name	Telephone	Name	Department
Leen Desmet	6531421	Bart Delvaux	Computer Science
Leen Desmet	09-23314	Leen Desmet	Philosophy
Bart Delvaux	5985625	Tom Demans	Computer Science
Tom Demans	5845213	David Finner	Biology

Assume that the database is complete with respect to the telephone number of all department members, but it is not complete regarding external collaborators. Thus, appropriate answers for the queries  $\text{Tel}(\text{Bart Delvaux}, 3962836)$  and  $\text{Tel}(\text{Leen Desmet}, 3212445)$  should be ‘no’ and ‘unknown’, respectively. If completeness of the database is taken for granted, the answer for both of these queries is ‘no’. Similarly, the answer under the CWA for the query  $\exists x \text{Tel}(\text{David Finner}, x)$  is ‘no’, but as the database is complete only with respect to the computer science department,

one cannot exclude the possibility that David Finner has a phone number, so the correct answer in this case should be ‘unknown’.

For dealing with situations like that of Example 1, the databases considered here are *locally closed* in the sense that in addition to a standard database instance they also contain a collection of *local closed world assumptions*, expressing conditions (the *windows of expertise*) that guarantee complete knowledge about specific database predicates.

**Example 2** In the context of the previous example, an expression of the form  $\text{LCWA}(\text{Tel}(x, y), \text{Dept}(x, \text{CS}))$  is intuitively understood as expressing that in the database instance, the relation  $\text{Tel}$  contains all the telephone numbers of all the members of the computer science department. Here, the window of expertise is  $\text{Dept}(x, \text{CS})$  and it is referring to the database relation  $\text{Tel}(x, y)$ .

Not surprisingly, query answering in locally closed databases (also referred here as locally complete databases) turns out to be intractable in general. This provides the motivation for developing efficient *approximate methods* for query answering. This approach is also motivated by the fact that in many applications there is no need to have *all* the answers to a query and often it is enough to have a sufficiently large subset of them. For instance, a company searching in an (incomplete) database for a provider of some urgently required service will be satisfied by finding *some* candidate providers, so an exhaustive search for all the providers is not always needed.

The current work builds upon (Cortés-Calabuig et al. 2007), which in turn relies on an extension of the formalism of Levy (1996) for representing partially complete information in database systems. In the former, a polynomial algorithm for querying a subclass of locally closed databases, called hierarchically closed, is introduced. This algorithm is based on implicit approximations of all the models of a locally closed database using three-valued structures. In order to guarantee completeness, strong syntactical conditions are imposed on the hierarchically closed databases and on the queries. Optimality results are provided for so-called conjunctive and disjunctive queries.

The approximation techniques mentioned above are extended and improved in this paper. In particular, we introduce a novel approach, based on query rewriting and fixpoint

computations, whose main benefits are the following:

- It generalizes the class of queries and databases considered in (Cortés-Calabuig et al. 2007). More specifically, in the new setting the query answering mechanism always terminates, even in case of arbitrary circular dependencies between predicates of LCWAs and the predicate in its windows of expertise.
- Certain types of integrity constraints are integrated in the query answering process. In particular, whenever the database is consistent, primary key constraints can be incorporated into the framework. Thus, in such cases, our framework may be used for approximative query answering in standard, constraint databases.
- A thorough analysis of the optimality of approximative query answering is developed. It shows that our approach retrieves correct and complete answers for a considerably larger class of queries and locally closed databases than those that are covered in (Cortés-Calabuig et al. 2007).

The rest of this paper is organized as follows. In the next section we recall the basic concepts, properties and semantics of locally closed databases as introduced in (Levy 1996; Cortés-Calabuig et al. 2005; 2006; 2007). Then we survey some previous results regarding query answering in locally closed databases and show an undecidability result concerning the closed-world information, which motivates the use of approximative query answering techniques. The main contribution of this paper is given in the next three sections, where these techniques are defined and analyzed. Finally, we discuss some related works and conclude.

## Representing Incompleteness in Databases

First, we recall the basic concepts concerning the local closed world assumption. The definitions in this section are taken from (Cortés-Calabuig et al. 2005).

We denote by  $\Sigma$  a finite first-order vocabulary, consisting of sets  $\mathcal{R}(\Sigma)$  of predicate symbols and  $\mathcal{C}(\Sigma)$  of constants. In addition,  $\Sigma$  contains equality  $=$  and the propositional constants  $\mathbf{t}$  (*true*) and  $\mathbf{f}$  (*false*), interpreted in the standard way. First-order formulas over  $\Sigma$  are constructed as usual.  $\Psi[\bar{x}]$  denotes a formula with free variables that are a subset of  $\bar{x}$ . An interpretation  $I$  for  $\Sigma$  ( $\Sigma$ -structures) is also defined as usual. In particular, a Herbrand interpretation has a domain  $\mathcal{C}(\Sigma)$ , such that each element of  $\mathcal{C}(\Sigma)$  interprets itself.

For a domain element  $a \in \text{Dom}(I)$ , let us define  $a^I = a$ . For a tuple  $\bar{t}$  of terms, we define  $\bar{t}^I = (t_1^I, \dots, t_n^I)$ . Setting  $\mathbf{f} \leq \mathbf{t}$  and  $\mathbf{f}^{-1} = \mathbf{t}$ ,  $\mathbf{t}^{-1} = \mathbf{f}$ , the truth value of a variable free formula  $\varphi$  in  $I$ , denoted  $\varphi^I$ , is defined recursively by

$$\begin{aligned} P(\bar{t})^I &= \mathbf{t} \text{ if } \bar{t}^I \in P^I, \text{ otherwise } P(\bar{t})^I = \mathbf{f}; \\ (\psi \wedge \phi)^I &= \min_{\leq}(\psi^I, \phi^I); \quad (\psi \vee \phi)^I = \max_{\leq}(\psi^I, \phi^I); \\ (\neg\psi)^I &= (\psi^I)^{-1}; \\ (\forall x \psi[x])^I &= \min_{\leq}\{(\psi[a])^I \mid a \in \text{Dom}(I)\}; \\ (\exists x \psi[x])^I &= \max_{\leq}\{(\psi[a])^I \mid a \in \text{Dom}(I)\}. \end{aligned}$$

We denote by  $I \models \varphi$  that  $\varphi^I = \mathbf{t}$ .

A *database instance*  $D$  is a Herbrand interpretation with a finite domain  $\text{Dom}^D \subseteq \mathcal{C}(\Sigma)$ . The set  $\text{Dom}^D$  is sometimes

called the *active domain* of the database instance and contains at least all the constants in the tables of  $D$  (and often only those). For some variable-free atomic formula  $A$ , we write  $A \in D$  to denote that  $A = P(\bar{d})$ , where  $\bar{d} \in P^D$ . In what follows, we will often specify a database instance  $D$  by a set of atoms. Unless the domain of  $D$  is explicitly mentioned, it consists of the set of constants that appear in these atoms. In our setting, databases represent *partial knowledge* on the domain of discourse, and as such their instances cannot be viewed as the (unique) possible state of the world.

**Definition 3** A local closed-world assumption (LCWA) is an expression  $\mathcal{LCWA}(P(\bar{x}), \Psi[\bar{x}])$ , where  $P \in \mathcal{R}(\Sigma)$  is called the LCWA's object and  $\Psi[\bar{x}]$ , called the LCWA's window of expertise, is a first-order formula over  $\Sigma$ .<sup>1</sup>

The intuitive reading of the expression in Definition 3 is the following: “for all objects  $\bar{x}$  such that  $\Psi(\bar{x})$  holds in the *real world*, if an atom of the form  $P(\bar{x})$  is true in the real world, then  $P(\bar{x})$  occurs in the *database*”. Note that in  $P(\bar{x})$  the value of the variables  $\bar{x}$  are constrained by the windows of expertise,  $\Psi$ .

**Definition 4** A locally closed database  $\mathfrak{D}$  over  $\Sigma$  is a pair  $(D, \mathcal{L})$  of a database instance  $D$  over  $\Sigma$  and a finite set  $\mathcal{L}$  of local closed-world assumptions over  $\Sigma$ .

Given a locally closed database  $\mathfrak{D} = (D, \mathcal{L})$ , we denote by  $\text{dom}(\mathfrak{D})$  its active domain. That is,  $\text{dom}(\mathfrak{D})$  is the union of active domain  $\text{Dom}^D$  of  $D$  and the set of constants in  $\mathcal{L}$ . A signature where  $\mathcal{R}(\Sigma_{\mathfrak{D}}) = \mathcal{R}(\Sigma)$  and  $\mathcal{C}(\Sigma_{\mathfrak{D}}) = \text{dom}(\mathfrak{D})$  is denoted by  $\Sigma_{\mathfrak{D}}$ . The *base predicates* of  $\mathfrak{D}$  are the propositional constants  $\mathbf{t}$ ,  $\mathbf{f}$  and all the predicates  $P$  such that the local closed-world assumption  $\mathcal{LCWA}(P(\bar{x}), \mathbf{t})$  appears in  $\mathcal{L}$ .

**Example 5** Abbreviate the database of Example 1 as follows:

$$D = \left\{ \begin{array}{l} \text{Tel}(\text{LD}, 6531421), \text{Tel}(\text{BD}, 5985625), \\ \text{Tel}(\text{TD}, 5845213), \text{Tel}(\text{LD}, 09-23314), \\ \text{Dept}(\text{BD}, \text{CS}), \text{Dept}(\text{LD}, \text{Phil}), \\ \text{Dept}(\text{TD}, \text{CS}), \text{Dept}(\text{DF}, \text{Bio}) \end{array} \right\}$$

As noted in Example 2, the local closed world assumption  $\mathcal{LCWA}(\text{Tel}(x, y), \text{Dept}(x, \text{CS}))$  states that all the telephone numbers of the computer science department members are known and occur in the database instance. That is, for every  $x_0$  in  $\{x \mid \text{Dept}(x, \text{CS})\}$  (the window of expertise for Tel), all true atoms of the form  $\text{Tel}(x_0, y)$  are in the database. Similarly,  $\mathcal{LCWA}(\text{Tel}(x, y), x = \text{LD})$  expresses that  $D$  contains all telephone numbers of Leen Desmet.

## Semantics of locally closed databases

The intuitive meaning behind the LCWA expressions of Definition 3 can be formally captured using first-order formulas.

<sup>1</sup>Definition 3 is based on the notion of Levy's local completeness statements (Levy 1996). Apart for the (innocent) difference that we use logical notation rather than Levy's database notation, LCWAs are more expressive by allowing arbitrary first-order formulas – instead of conjunction of atoms – in the window of expertise.

**Definition 6** Let  $D$  be a database based over a first-order vocabulary  $\Sigma$  and let  $P$  be a predicate in  $D$ . Denote by  $P^D$  the  $P$ -tuples in  $D$ . Given a tuple  $\bar{x}$  of terms, we denote by  $P(\bar{x}) \in D$  the formula  $\bigvee_{\bar{a} \in P^D} (\bar{x} = \bar{a})$ .

**Definition 7** Let  $D$  be a database over a vocabulary  $\Sigma$  and let  $\theta = \mathcal{LCWA}(P(\bar{x}), \Psi[\bar{x}])$  be an LCWA over  $\Sigma$ . The meaning of  $\theta$  in  $D$  is given by the formula

$$\mathcal{M}_D(\theta) = \forall \bar{x} (\Psi[\bar{x}] \supset (P(\bar{x}) \supset (P(\bar{x}) \in D))).$$

Observe that the meaning of a local closed-world assumption is induced by the database instance under consideration, as it contains the subformula  $P(\bar{x}) \in D$ . As such, a local closed-world assumption is a *non-monotonic* construct.

**Example 8** Consider again Example 5. The meaning of  $\theta = \mathcal{LCWA}(\text{Tel}(x, y), \text{Dept}(x, \text{CS}))$  is given by

$$\mathcal{M}_D(\theta) = \forall x \forall y (\text{Dept}(x, \text{CS}) \supset (\text{Tel}(x, y) \supset ((x = \text{LD} \wedge y = 6531421) \vee (x = \text{LD} \wedge y = 09-23314) \vee (x = \text{BD} \wedge y = 5985625) \vee (x = \text{TD} \wedge y = 5845213)))).$$

As shown in (Cortés-Calabuig et al. 2006), any collection of LCWAs on the same predicate may be combined into one (disjunctive) LCWA. We therefore assume that each predicate  $P$  in  $\mathcal{R}(\Sigma)$  is the object of exactly *one* LCWA expression, whose window of expertise is denoted  $\Psi_P$ .

The meaning of a locally closed database  $\mathfrak{D} = (D, \mathcal{L})$  is expressed by a first-order formula consisting of the conjunction of the database atoms, the meaning of the given local closed-world assumptions, and the following two axioms:

- *Domain Closure*:  $\text{DCA}(\text{dom}(\mathfrak{D})) = \forall x (\bigvee_{i=1}^n x = c_i)$
- *Unique Names*:  $\text{UNA}(\text{dom}(\mathfrak{D})) = \bigwedge_{1 \leq i < j \leq n} c_i \neq c_j$

where  $\text{dom}(\mathfrak{D}) = \{c_1, \dots, c_n\}$ .

**Definition 9** Let  $\mathfrak{D} = (D, \mathcal{L})$  be a locally closed database over  $\Sigma$ . The meaning  $\mathcal{M}(\mathfrak{D})$  of  $\mathfrak{D}$  is the first-order sentence

$$\text{UNA}(\text{dom}(\mathfrak{D})) \wedge \text{DCA}(\text{dom}(\mathfrak{D})) \wedge \bigwedge_{A \in D} A \wedge \bigwedge_{\theta \in \mathcal{L}} \mathcal{M}_D(\theta).$$

The formula  $\mathcal{M}(\mathfrak{D})$  expresses incomplete knowledge about the real world. Thus, in general, it has several (Herbrand) models. A  $\Sigma_{\mathfrak{D}}$ -model  $M$  of  $\mathcal{M}(\mathfrak{D})$  is called a model of  $\mathfrak{D}$ , and this is denoted by  $M \models \mathfrak{D}$ . If every model of  $\mathfrak{D}$  is also a model of a formula  $\varphi$  over  $\Sigma_{\mathfrak{D}}$  we say that  $\mathfrak{D}$  *entails*  $\varphi$  (or  $\varphi$  *follows* from  $\mathfrak{D}$ ), and denote this by  $\mathfrak{D} \models \varphi$ .

## Query Answering in Locally Closed Databases

In this section, we provide the basic tools for reasoning with locally closed databases. Query answering in such databases may be represented as follows:

**Definition 10** Given a locally closed database  $\mathfrak{D}$  over  $\Sigma$ , a first-order query  $Q[\bar{x}]$  over  $\Sigma$  (whose free variables are in  $\bar{x}$ ), and a tuple  $\bar{t}$  of constants in  $\text{dom}(\mathfrak{D})$ , we say that

- $\bar{t}$  is a certain answer in  $\mathfrak{D}$  for  $Q[\bar{x}]$ , if  $\mathfrak{D} \models Q[\bar{t}/\bar{x}]$
- $\bar{t}$  is a possible answer in  $\mathfrak{D}$  for  $Q[\bar{x}]$ , if  $\mathfrak{D} \cup Q[\bar{t}/\bar{x}]$  is satisfiable (or, equivalently, if  $\mathfrak{D} \not\models \neg Q[\bar{t}/\bar{x}]$ ).

We denote by  $\text{Cert}_{\mathfrak{D}}(Q[\bar{x}])$  the set of certain answers of  $Q[\bar{x}]$  in  $\mathfrak{D}$  and by  $\text{Poss}_{\mathfrak{D}}(Q[\bar{x}])$  the set of possible answers of  $Q[\bar{x}]$  in  $\mathfrak{D}$ .

An interesting property of a query  $Q[\bar{x}]$  in a locally closed database  $\mathfrak{D}$  is whether  $\mathfrak{D}$  has complete knowledge on  $Q[\bar{x}]$ . This has been defined as follows:

**Definition 11** (Levy 1996) A locally closed database  $\mathfrak{D}$  over  $\Sigma$  has complete world information (CWI) on a query  $Q[\bar{x}]$ , if for each tuple  $\bar{t}$  of constants in  $\text{dom}(\mathfrak{D})$ , either  $\mathfrak{D} \models Q[\bar{t}]$  or  $\mathfrak{D} \models \neg Q[\bar{t}]$ .

Next, we focus on the computational complexity of query answering. Following the usual measure of complexity in database systems, the results below are specified in terms of data complexity, that is, in terms of the size  $|D|$  of the database instance (assuming that all the rest is fixed). Accordingly, we consider the following decision problems:

$$\begin{aligned} \text{Poss}_{\mathcal{L}}(Q[\bar{x}]) &= \{(D, \bar{t}) \mid \bar{t} \in \text{Poss}_{(D, \mathcal{L})}(Q[\bar{x}])\}, \\ \text{Cert}_{\mathcal{L}}(Q[\bar{x}]) &= \{(D, \bar{t}) \mid \bar{t} \in \text{Cert}_{(D, \mathcal{L})}(Q[\bar{x}])\}, \\ \text{CWI}_{\mathcal{L}}(Q[\bar{x}]) &= \{D \mid (D, \mathcal{L}) \text{ has CWI on } Q[\bar{x}]\}. \end{aligned}$$

As proven in (Cortés-Calabuig et al. 2007), the decision problem  $\text{Poss}_{\mathcal{L}}(Q[\bar{x}])$  is in NP for all  $\mathcal{L}$  and  $Q[\bar{x}]$ , and is NP-hard for some of them.  $\text{Cert}_{\mathcal{L}}(Q[\bar{x}])$  is in coNP for each  $\mathcal{L}$  and  $Q[\bar{x}]$ , and is coNP-hard for some of them.

When  $\mathfrak{D}$  has complete information about a query, there is no uncertainty about its answers, so such queries are of practical importance. As Proposition 12 shows, queries with CWI can be answered directly in the database instance  $D$ , when  $D$  is regarded as a two-valued Herbrand structure of  $\mathfrak{D}$ .<sup>2</sup>

**Proposition 12** A locally closed database  $\mathfrak{D}$  has CWI on  $Q[\bar{x}]$  iff  $\text{Cert}_{\mathfrak{D}}(Q[\bar{x}]) = \text{Poss}_{\mathfrak{D}}(Q[\bar{x}]) = \{\bar{a} \mid D \models Q[\bar{a}]\}$ .

*Proof.* See the appendix.  $\square$

Deciding whether there is CWI on a query  $Q[\bar{x}]$  in a specific database  $\mathfrak{D} = (D, \mathcal{L})$  is not tractable (Cortés-Calabuig et al. 2007). The next proposition shows that the more ambitious problem, of whether there is CWI on  $Q[\bar{x}]$  in *all* locally closed databases containing a fixed set  $\mathcal{L}$  of local closed-world assumptions, is not even decidable.

**Proposition 13** The question whether all the locally closed databases  $(\cdot, \mathcal{L})$  convey CWI on a query  $Q[\bar{x}]$  is undecidable.

*Proof.* See the appendix.  $\square$

So far, the results in this section give little reason for optimism regarding practical applicability of local closed-world assumptions. But, as it turns out, in many applications there is no need to have *all* certain answers to a query. Often, it suffices to have a sufficiently large subset of the answers. E.g., if a company searches an (incomplete) database for a provider of some urgently required service, it will be happy if it finds *some* candidate providers; this list does not need to be complete. Likewise, in many applications, it would not harm if the answers to a *possible* query contain a few extra

<sup>2</sup>This was Levy's motivation to study CWI.

“impossible” elements. For instance, if a company wants to advertise one of its services and queries the above database for a group of potential clients, it would not care to receive some additional companies that could not really be possibly interested. Thus, one reasonable strategy to solve the complexity problem would be to develop tractable approximate methods. This is the approach followed in the next section.

The other, more conventional approach to the complexity problem, is to restrict the expressiveness of the language so that efficient query processing is possible. As it turns out, below we obtain such results as well, though in a slightly indirect way: we will show that for certain classes of queries and local closed-world assumptions, the approximate methods are *optimal* in the sense that they compute exactly the certain and possible answers to queries. Thus, these combinations of queries and local closed-world assumptions provide tractable sub-languages.

## Approximative Reasoning

### Approximations by Three-Valued Structures

The basic idea of approximative reasoning is to compute a 3-valued structure that provides a ‘good approximation’ of all models of  $\mathcal{D}$  and then to evaluate queries with respect to this structure. The underlying semantics is, therefore, a 3-valued one, where the truth values  $\mathcal{T}\mathcal{H}\mathcal{R}\mathcal{E}\mathcal{E} = \{\mathbf{t}, \mathbf{f}, \mathbf{u}\}$  stand for true, false, and unknown (respectively). These values are usually arranged in two orders: the *truth order*,  $\leq$ , which is a linear order given by  $\mathbf{f} \leq \mathbf{u} \leq \mathbf{t}$ , and the *precision order*,  $\leq_p$ , which is a partial order on  $\mathcal{T}\mathcal{H}\mathcal{R}\mathcal{E}\mathcal{E}$  in which  $\mathbf{u}$  is the least element, and  $\mathbf{t}$  and  $\mathbf{f}$  are incomparable maximal elements. The connectives are defined according to the truth order: Conjunction  $\wedge$ , disjunction  $\vee$  and the negation operator  $\neg$  are defined, respectively, by the  $\leq$ -glb,  $\leq$ -lub, and the  $\leq$ -involution (that is,  $\neg\mathbf{t} = \mathbf{f}$ ,  $\neg\mathbf{f} = \mathbf{t}$ , and  $\neg\mathbf{u} = \mathbf{u}$ ) on  $\mathcal{T}\mathcal{H}\mathcal{R}\mathcal{E}\mathcal{E}$ .

The notions of 3-valued (Herbrand) structures and (Herbrand) models are defined with respect to  $\mathcal{T}\mathcal{H}\mathcal{R}\mathcal{E}\mathcal{E}$  in the standard way. The three-valued Herbrand interpretations of  $\Sigma$  are denoted  $\mathcal{L}^c$  and the subset of two-valued structures is denoted  $\mathcal{L}$ . A truth order  $\leq$  and a precision order  $\leq_p$  are also definable on  $\mathcal{L}^c$  by pointwise extensions of the corresponding orders in  $\mathcal{T}\mathcal{H}\mathcal{R}\mathcal{E}\mathcal{E}$ . Clearly,  $\leq$  is a lattice order and  $\leq_p$  is a chain-complete order on  $\mathcal{L}^c$ .

Truth assignment in the three-valued case is defined through the same recursive rules as two-valued truth assignment (but with respect to different orders). For instance,

$$\begin{aligned} (\psi \wedge \phi)^{\mathcal{K}} &= \min_{\leq}(\psi^{\mathcal{K}}, \phi^{\mathcal{K}}); \\ (\neg\psi)^{\mathcal{K}} &= (\psi^{\mathcal{K}})^{-1}; \\ (\forall x \psi[x])^{\mathcal{K}} &= \min_{\leq} \{(\psi[a])^{\mathcal{K}} \mid a \in \text{Dom}(\mathcal{K})\}; \end{aligned}$$

In what follows we simulate three-valued truth assignments by two-valued truth assignments as follows: given a vocabulary  $\Sigma$ , we introduce for each predicate  $P \in \mathcal{R}(\Sigma)$  two new predicate symbols  $P^c$  and  $P^{c\neg}$  (intuitively standing for ‘certainly  $P$ ’ and ‘certainly not  $P$ ’, respectively). Denote by  $\Sigma'$  the set of all constant and predicate symbols of  $\Sigma$  together with all the new predicate symbols.

**Definition 14** A 2-valued  $\Sigma'$ -structure  $I$  simulates a three-valued  $\Sigma$ -structure  $\mathcal{K}$ , iff  $\mathcal{K}$  and  $I$  have the same domain, they assign the same interpretations to constant and function symbols, and for each predicate  $P \in \mathcal{R}(\Sigma)$  it holds that  $(P^c)^I = \{\bar{d} \mid P(\bar{d})^{\mathcal{K}} = \mathbf{t}\}$  and  $(P^{c\neg})^I = \{\bar{d} \mid P(\bar{d})^{\mathcal{K}} = \mathbf{f}\}$ .

In the following definition,  $P_i^c, P_i^p, P_i^{c\neg}$  and  $P_i^{p\neg}$  are symbols representing respectively the certain and the possible tuples of  $P_i$ , and the tuples that certainly and possibly do not belong to  $P_i$ . Accordingly,  $\Phi^c$  and  $\Phi^p$  represent the certain instances and the possible instances of  $\Phi$  when interpreted as a query. As noted in Proposition 16, these formulas can be used to compute three-valued answers for  $\Phi$ .

**Definition 15** Given a database vocabulary  $\Sigma$ , we introduce, for each element in  $\mathcal{R}(\Sigma) = \{P_1, \dots, P_n\}$ , four new predicate symbols  $P_i^c, P_i^p, P_i^{c\neg}$  and  $P_i^{p\neg}$  of the same arity as  $P_i$ . Now, each formula  $\Phi$  with predicate symbols amongst  $P_1, \dots, P_n$  is associated with the following two formulas:

- $\Phi^c$  is the formula obtained by substituting  $P_i^c(\bar{t})$  for each positive occurrence<sup>3</sup> of  $P_i(\bar{t})$  in  $\Phi$ , and substituting  $\neg P_i^{c\neg}(\bar{t})$  for each negative occurrence of  $P_i(\bar{t})$  in  $\Phi$ .
- $\Phi^p$  is the formula obtained by substituting  $P_i^p(\bar{t})$  for each positive occurrence of  $P_i(\bar{t})$  in  $\Phi$ , and substituting  $\neg P_i^{p\neg}(\bar{t})$  for each negative occurrence of  $P_i(\bar{t})$  in  $\Phi$ .

Note that  $(\neg P(\bar{t}))^c = \neg\neg P^{c\neg}(\bar{t}) \equiv P^{c\neg}(\bar{t})$ . Also,  $P^p(\bar{t})$  and  $\neg P^{c\neg}(\bar{t})$  are equivalent and so are  $P^{p\neg}(\bar{t})$  and  $\neg P^c(\bar{t})$ . Moreover,  $\Phi^c$  contains only positive occurrences of  $P_i^c(\bar{t})$  and  $P_i^{c\neg}(\bar{t})$ . Similarly,  $\Phi^p$  contains only positive occurrences of  $P_i^p(\bar{t})$  and  $P_i^{p\neg}(\bar{t})$ .

The following proposition is well known.

**Proposition 16** If  $I$  simulates  $\mathcal{K}$ , then for each formula  $\varphi[\bar{x}]$  and a suitable tuple of domain elements  $\bar{d}$ ,  $\varphi[\bar{d}]^{\mathcal{K}} = \mathbf{t}$  iff  $(\varphi[\bar{d}]^c)^I = \mathbf{t}$  and  $\varphi[\bar{d}]^{\mathcal{K}} = \mathbf{f}$  iff  $((\neg\varphi[\bar{d}])^c)^I = \mathbf{f}$ .

This implies tractability of three-valued truth evaluation and query answering.

**Corollary 17** Given a finite three-valued  $\Sigma$ -structure  $\mathcal{K}$ , for each formula  $\varphi[\bar{x}]$ ,  $\{\bar{d} \mid (\varphi[\bar{d}])^{\mathcal{K}} = \mathbf{t}\}$ ,  $\{\bar{d} \mid (\varphi[\bar{d}])^{\mathcal{K}} = \mathbf{f}\}$ , and  $\{\bar{d} \mid (\varphi[\bar{d}])^{\mathcal{K}} = \mathbf{u}\}$  can be computed in polynomial time in the size of  $\mathcal{K}$ .

We now consider approximation theory. The definitions in the rest of this section and in the next subsection are taken from (Cortés-Calabuig et al. 2006).

**Definition 18** Let  $\Gamma$  be a satisfiable theory based on  $\Sigma$  and containing  $\text{UNA}(\Sigma)$  and  $\text{DCA}(\Sigma)$ . We say that a 3-valued Herbrand  $\Sigma$ -interpretation  $\mathcal{K}$  approximates  $\Gamma$  (from below), iff for every 2-valued Herbrand model  $M$  of  $\Gamma$ ,  $\mathcal{K} \leq_p M$ . The optimal approximation for  $\Gamma$  is the 3-valued Herbrand structure  $\mathcal{O}_\Gamma = \text{glb}_{\leq_p} \{M \mid M \models \Gamma\}$ , where  $M$  ranges over all the 2-valued Herbrand models of  $\Gamma$ .

Note that  $\mathcal{O}_\Gamma$  is the most precise among all 3-valued Herbrand  $\Sigma$ -structures approximating  $\Gamma$  and it is well-defined,

<sup>3</sup>A predicate occurs *positively* (alternatively *negatively*) in a formula, iff it appears in the scope of an even (alternatively odd) number of negation symbols.

since the set of  $\Gamma$ 's Herbrand models is non-empty and every nonempty set  $S \subseteq \mathcal{L}^c$  has a greatest  $\leq_p$ -lower bound. Now, by the fact that all models of a theory containing  $\text{UNA}(\Sigma) \wedge \text{DCA}(\Sigma)$  are isomorphic to Herbrand structures, we have:

**Proposition 19** *Let  $\mathcal{K}$  be an approximation of  $\Gamma$ . For any sentence  $\varphi$ , if  $\varphi^{\mathcal{K}} = \mathbf{t}$ , then  $\Gamma \models \varphi$  and if  $\varphi^{\mathcal{K}} = \mathbf{f}$ , then  $\Gamma \models \neg\varphi$ .*

**Definition 20** *For a 3-valued  $\Sigma$ -interpretation  $\mathcal{K}$  and a query  $\mathcal{Q}[\bar{x}]$  in  $\Sigma$ , define the certain answers and the possible answers of  $\mathcal{Q}[\bar{x}]$  w.r.t.  $\mathcal{K}$  by the following sets (respectively):*

- $\text{Cert}_{\mathcal{K}}(\mathcal{Q}[\bar{x}]) = \{\bar{a} \mid \mathcal{Q}[\bar{a}]^{\mathcal{K}} = \mathbf{t}\}$ ,
- $\text{Poss}_{\mathcal{K}}(\mathcal{Q}[\bar{x}]) = \{\bar{a} \mid \mathcal{Q}[\bar{a}]^{\mathcal{K}} \geq \mathbf{u}\}$ .

As computing truth values of sentences is polynomial, we have

**Proposition 21** *For each finite three-valued  $\Sigma$ -structure  $\mathcal{K}$  and  $\Sigma$ -query  $\mathcal{Q}[\bar{x}]$ , the sets  $\text{Cert}_{\mathcal{K}}(\mathcal{Q}[\bar{x}])$  and  $\text{Poss}_{\mathcal{K}}(\mathcal{Q}[\bar{x}])$  are polynomially computable in the size of  $\mathcal{K}$ .*

### Query Answering by Fixpoint Computations

From Proposition 21 it is clear that a tractable method to compute 3-valued approximations induces a tractable and sound approximative query answering method. Such a method is defined in (Cortés-Calabuig et al. 2006) as follows:

**Definition 22** *Given a locally closed database  $\mathfrak{D} = (D, \mathcal{L})$ , the operator  $\text{App}_{\mathfrak{D}} : \mathcal{L}^c \rightarrow \mathcal{L}^c$  maps a three-valued structure  $\mathcal{K}$  to a three-valued structure  $\mathcal{K}' = \text{App}_{\mathfrak{D}}(\mathcal{K})$  such that, for every predicate  $P$  of  $\mathcal{R}(\Sigma)$  and every tuple  $\bar{a}$ ,*

$$P(\bar{a})^{\mathcal{K}'} = \begin{cases} \mathbf{t} & \text{if } P(\bar{a}) \in D, \\ \mathbf{f} & \text{if there exists } \text{LCWA}(P(\bar{x}), \Psi_P[\bar{x}]) \in \mathcal{L} \\ & \text{so that } \Psi_P[\bar{a}]^{\mathcal{K}} = \mathbf{t} \text{ and } P(\bar{a}) \notin D, \\ \mathbf{u} & \text{otherwise.} \end{cases}$$

The idea here is to start from the structure with total ignorance (i.e., a valuation that assigns  $\mathbf{u}$  to every ground atom), and to iterate  $\text{App}_{\mathfrak{D}}$ , thereby gradually extending the definite knowledge using the database and its LCWAs. Clearly,  $\text{App}_{\mathfrak{D}}$  is  $\leq_p$ -monotone. Thus, by (an extension of) the well-known Knaster-Tarski theorem, we have

**Proposition 23**  *$\text{App}_{\mathfrak{D}}$  is a  $\leq_p$ -monotone operator on the chain complete poset  $\mathcal{L}^c$ , thus it has a  $\leq_p$ -least fixpoint.*

**Definition 24** *Denote by  $\mathcal{C}_{\mathfrak{D}}$  the  $\leq_p$ -least fixpoint of  $\text{App}_{\mathfrak{D}}$ .*

**Note 25** *As the number of iterations for reaching  $\mathcal{C}_{\mathfrak{D}}$  is at most polynomial in the size of the database, and each iteration takes polynomial time in the size of the database, it follows that  $\mathcal{C}_{\mathfrak{D}}$  can be computed in polynomial time in  $|D|$ .*

The following proposition shows that  $\mathcal{C}_{\mathfrak{D}}$  is a sound approximation of  $\mathfrak{D}$ .

**Proposition 26**  *$\mathcal{C}_{\mathfrak{D}}$  approximates  $\mathfrak{D}$  and for every  $\mathcal{Q}[\bar{x}]$  we have  $\text{Cert}_{\mathcal{C}_{\mathfrak{D}}}(\mathcal{Q}[\bar{x}]) \subseteq \text{Cert}_{\mathfrak{D}}(\mathcal{Q}[\bar{x}]) \subseteq \text{Poss}_{\mathfrak{D}}(\mathcal{Q}[\bar{x}]) \subseteq \text{Poss}_{\mathcal{C}_{\mathfrak{D}}}(\mathcal{Q}[\bar{x}])$ .*

### Fixpoint Queries for LCWA

We now come to the main contribution of this paper, namely: computing fixpoint queries for locally closed databases, and analyzing the accuracy of approximative query answering in the context defined above.

A substantial drawback of query answering with  $\mathcal{C}_{\mathfrak{D}}$  is the need to recompute it each time the database changes. In what follows we partially avoid this by using fixpoint formulas that symbolically describe the construction of  $\mathcal{C}_{\mathfrak{D}}$ . Using these expressions, certain or possible answers to queries can be computed by transforming the query into a fixpoint query or a query with respect to some datalog program. This means, in practice, that it suffices to compute the relations that are relevant for the query rather than computing all the relations in  $\mathcal{C}_{\mathfrak{D}}$ . Moreover, goal directed methods such as magic sets (Bancilhon et al. 1986) or tabling (Swift 1999), will often need only fractions of those relations.

**Definition 27** *Let  $\mathfrak{D} = (D, \mathcal{L})$  be a locally closed database. For a query  $\mathcal{Q}[\bar{x}]$  we introduce two new variables  $Q^c$  and  $Q^{c\neg}$ , the arity of which is the number of free variables of  $\mathcal{Q}[\bar{x}]$ , and define the set  $\Delta_{\mathcal{Q}, \mathcal{L}}$  by*

$$\left\{ \begin{array}{l} Q^c(\bar{x}) \leftarrow \mathcal{Q}[\bar{x}]^c \\ Q^{c\neg}(\bar{x}) \leftarrow (\neg\mathcal{Q}[\bar{x}])^c \end{array} \right\} \cup \bigcup \left\{ \begin{array}{l} P_i^c(\bar{x}_i) \leftarrow P_i(\bar{x}_i) \\ P_i^{c\neg}(\bar{x}_i) \leftarrow \neg P_i(\bar{x}_i) \wedge (\Psi_{P_i}[\bar{x}_i])^c \end{array} \right\}$$

where the right union is over the database predicates  $P_i$ , and  $\Psi_{P_i}$  is the window of expertise of  $P_i$ .

Intuitively,  $Q^c$  is meant to represent the collection of certain instances of  $\mathcal{Q}[\bar{x}]$ , and  $Q^{c\neg}$  represents the certain instances of  $\neg\mathcal{Q}[\bar{x}]$ . This is captured by the fixpoint computations on  $\Delta_{\mathcal{Q}, \mathcal{L}}$  in the logic  $\text{LFP}^{\text{simult}}$ , described below. The fixpoint expressions that are simultaneously computed in this logic are of the form

$$[\text{lfp}_{R_i, \Delta}](\bar{t})$$

for some predicates  $\{R_1, \dots, R_n\}$ , where  $\bar{t}$  is a tuple of terms whose arity is the same as that of  $R_i$ , and  $\Delta$  is a collection of rules of the form

$$\{R_j(\bar{x}_j) \leftarrow \varphi_j[\bar{x}_j] \mid 1 \leq j \leq n\}^4$$

**Definition 28** *Let  $\Gamma_{\Delta}$  be the standard immediate consequence fixpoint operator on  $\Delta$ . A fixpoint expression  $[\text{lfp}_{R_i, \Delta}](\bar{t})$  is true in a structure  $\mathfrak{A}$  and variable assignment  $\nu$ , if  $\bar{t}^{\mathfrak{A}, \nu} \in \mathcal{R}_i$ , where  $\mathcal{R}_i$  is the  $i$ 'th argument in the least fixpoint  $(\mathcal{R}_1, \dots, \mathcal{R}_n)$  of  $\Gamma_{\Delta}$ , associated to  $\Delta$  and  $\mathfrak{A}$ . Now, given a locally closed database  $\mathfrak{D} = (D, \mathcal{L})$ , define the certain query answer for  $\mathcal{Q}[\bar{x}]$  as  $[\text{lfp}_{Q^c, \Delta_{\mathcal{Q}, \mathcal{L}}}]^c(\bar{x})$  and the possible query answer for  $\mathcal{Q}[\bar{x}]$  as  $\neg[\text{lfp}_{Q^{c\neg}, \Delta_{\mathcal{Q}, \mathcal{L}}}]^{c\neg}(\bar{x})$ , where both of these expressions are evaluated in  $D$ .*

It is worth noting that  $\Delta$  is an extended datalog program as defined in (Van Gelder 1993) or a positive definition as defined in FO[ID] (Denecker and Ternovska 2007) and that its semantics, i.e., its least fixpoint, coincides with the well-founded model of  $\Delta$ . It follows that  $Q^c$  is the collection of

<sup>4</sup>Here, each  $\varphi_j$  – the definition of  $R_j$  – is a formula over  $\Sigma$  and  $\{R_1, \dots, R_n\}$ , where  $R_1, \dots, R_n$  may occur positively in  $\varphi_j$ , and the length of  $\bar{x}_j$  is the same as the arity of  $R_j$ .

certain instances of  $Q[\bar{x}]$  and  $Q^{c\bar{\neg}}$  is the collection of certain instances of  $\neg Q[\bar{x}]$ . Those instances are represented by  $[\text{lfp}_{Q^c, \Delta_{Q, \mathcal{L}}}](\bar{x})$  and  $[\text{lfp}_{Q^{c\bar{\neg}}, \Delta_{Q, \mathcal{L}}}](\bar{x})$ , respectively.

**Example 29** Consider again the local closed-world assumption in Example 5:  $\text{LCWA}(\text{Tel}(x, y), \text{Dept}(x, \text{CS}))$ . Assume that there is no closure for the relation Dept, that is,  $\text{LCWA}(\text{Dept}(x, y), \mathbf{f})$ . Let  $Q = \text{Tel}(\text{BD}, 3962836)$  (as in Example 1). Then:

$$\Delta_{Q, \mathcal{L}} = \left\{ \begin{array}{l} Q^c \leftarrow \text{Tel}^c(\text{BD}, 3962836). \\ Q^{c\bar{\neg}} \leftarrow \text{Tel}^{c\bar{\neg}}(\text{BD}, 3962836). \\ \text{Tel}^c(x, y) \leftarrow \text{Tel}(x, y). \\ \text{Tel}^{c\bar{\neg}}(x, y) \leftarrow \neg \text{Tel}(x, y) \wedge \text{Dept}^c(x, \text{CS}). \\ \text{Dept}^c(x, y) \leftarrow \text{Dept}(x, y). \\ \text{Dept}^{c\bar{\neg}}(x, y) \leftarrow \neg \text{Dept}(x, y) \wedge \mathbf{f}. \end{array} \right\}$$

As  $\text{lfp}_{Q^{c\bar{\neg}}, \Delta_{Q, \mathcal{L}}}$  is true in  $D$ ,  $\text{Tel}(\text{BD}, 3962836)$  is certainly false.

**Proposition 30** Given a locally closed database  $(D, \mathcal{L})$  and a query  $Q[\bar{x}]$ . Let  $(\mathcal{R}_Q^c, \mathcal{R}_Q^{c\bar{\neg}}, \mathcal{R}_1^c, \mathcal{R}_1^{c\bar{\neg}}, \dots, \mathcal{R}_n^c, \mathcal{R}_n^{c\bar{\neg}})$  be the relations defined by  $\Delta_{Q, \mathcal{L}}$  in  $D$ . Then, for all  $1 \leq i \leq n$ ,

$$\begin{aligned} \mathcal{R}_i^c &= \{\bar{d} \mid P_i(\bar{d})^{c\bar{\neg}} = \mathbf{t}\}, & \mathcal{R}_i^{c\bar{\neg}} &= \{\bar{d} \mid P_i(\bar{d})^c = \mathbf{f}\}, \\ \mathcal{R}_Q^c &= \{\bar{d} \mid Q(\bar{d})^{c\bar{\neg}} = \mathbf{t}\}, & \mathcal{R}_Q^{c\bar{\neg}} &= \{\bar{d} \mid Q(\bar{d})^c = \mathbf{f}\}. \end{aligned}$$

*Proof.* See the appendix.  $\square$

**Note 31** Dually, one may define the set  $\Delta'_{Q, \mathcal{L}}$  by

$$\left\{ \begin{array}{l} Q^p(\bar{x}) \leftarrow Q[\bar{x}]^p \\ Q^{p\bar{\neg}}(\bar{x}) \leftarrow (\neg Q[\bar{x}])^p \end{array} \right\} \cup \bigcup \left\{ \begin{array}{l} P_i^p(\bar{x}_i) \leftarrow P_i(\bar{x}_i) \vee (\neg \Psi_{P_i}[\bar{x}_i])^p \\ P_i^{p\bar{\neg}}(\bar{x}_i) \leftarrow \neg P_i(\bar{x}_i) \end{array} \right\}$$

and then consider the expressions of the greatest fixpoints  $\neg[\text{gfp}_{Q^{p\bar{\neg}}, \Delta'_{Q, \mathcal{L}}}](\bar{x})$  and  $[\text{gfp}_{Q^p, \Delta'_{Q, \mathcal{L}}}](\bar{x})$  for representing the certain and the possible answers of  $Q[\bar{x}]$ , respectively. This follows from the fact that if  $(\mathcal{R}_1, \dots, \mathcal{R}_n)$  is the least fixpoint of  $\Gamma_\Delta$ , and  $(\mathcal{R}'_1, \dots, \mathcal{R}'_n)$  is the greatest fixpoint of  $\Gamma_{\Delta'}$ , the relations  $\mathcal{R}_i, \mathcal{R}'_i$  are complements, and so  $[\text{lfp}_{\mathcal{R}_i, \Delta}](\bar{t})$  and  $\neg[\text{gfp}_{\mathcal{R}'_i, \Delta'}](\bar{t})$  are logically equivalent.

**Note 32** It is interesting to note that for consistent databases, the approximate propagation mechanism, as implemented by the rules above, may be used for imitating common types of integrity constraints in database systems. For instance, in Examples 5 and 29, a primary key constraint on telephone-numbers, that is:

$$\forall xyz(\text{Tel}(x, z) \wedge \text{Tel}(y, z)) \supset x = y,$$

can be expressed as the rewriting rule

$$\text{Tel}^{c\bar{\neg}}(x, z) \leftarrow \exists y \text{Tel}^c(y, z) \wedge x \neq y.$$

Similarly, a foreign key constraint of the form

$$\forall xz(\text{Tel}(x, z) \supset \exists y \text{Dept}(x, y))$$

can be expressed as the rewriting rule

$$\text{Tel}^{c\bar{\neg}}(x, y) \leftarrow \forall z \text{Dept}^{c\bar{\neg}}(x, z).$$

Observe that these extra rules allow us to derive more negative information, e.g.,  $\text{Tel}^{c\bar{\neg}}(\text{LD}, 5845213)$ .

## The Accuracy of Approximate Query Answering

The results above give us a tractable method for computing possible and certain answers to queries by first computing  $\mathcal{C}_\mathcal{D}$  and then evaluating queries against it, using standard database techniques. Tractability, however, has a price. As the next example shows, in certain cases we lose accuracy.

**Example 33** Below, we abbreviate the optimal approximation of  $\mathcal{D} = (D, \mathcal{L})$  by  $\mathcal{O}_\mathcal{D}$  (instead of  $\mathcal{O}_{\mathcal{M}(\mathcal{D})}$ ).

1. Let  $D = \emptyset$  and  $\mathcal{L} = \{\text{LCWA}(P, P)\}$ . As the meaning of  $\mathcal{D}$  is equivalent to  $\neg P$  in this case, we have that  $P^{\mathcal{O}_\mathcal{D}} = \mathbf{f}$ . Yet,  $P^{\mathcal{C}_\mathcal{D}} = \mathbf{u}$ , as  $P$  is unknown in the initial step of the constructed approximation of  $\mathcal{D}$ , and remains unknown when applying  $\text{App}_\mathcal{D}$  (see Definition 22).
2. Let  $D = \emptyset$  and  $\mathcal{L} = \{\text{LCWA}(Q, P \vee \neg P)\}$ . This database has models in which  $P$  is true and others in which  $P$  is false but, because of its LCWA,  $Q$  is false in all of them. Thus,  $P^{\mathcal{C}_\mathcal{D}} = \mathbf{u}$  and  $Q^{\mathcal{C}_\mathcal{D}} = \mathbf{f}$ . However, since  $P \vee \neg P$  evaluates to  $\mathbf{u}$  in each structure  $\mathcal{K}$  for which  $P^\mathcal{K} = \mathbf{u}$ , we have that  $Q^{\mathcal{C}_\mathcal{D}} = \mathbf{u}$ . The answer for the query  $\neg Q$  in  $\mathcal{C}_\mathcal{D}$  is therefore  $\mathbf{u}$ , while it is  $\mathbf{t}$  with respect to  $\mathcal{D}$  or  $\mathcal{O}_\mathcal{D}$ .
3. Let  $D = \emptyset$  and  $\mathcal{L} = \{\text{LCWA}(P, R), \text{LCWA}(Q, R \supset \neg P)\}$ . Here,  $\mathcal{M}(\mathcal{D}) = (R \supset \neg P) \wedge ((R \supset \neg P) \supset \neg Q)$ , which obviously entails  $\neg Q$ , and so  $Q^{\mathcal{C}_\mathcal{D}} = \mathbf{f}$ . The fact that in this case the window of expertise of the second LCWA is exactly the meaning of the first LCWA is not captured by  $\mathcal{C}_\mathcal{D}$ , and so  $Q^{\mathcal{C}_\mathcal{D}} = \mathbf{u}$ .

In what follows, we consider some important cases in which accuracy of the approximation method is guaranteed.

As the first item of Example 33 shows, the precision of  $\mathcal{C}_\mathcal{D}$  may be lost when cycles exist in the local closed world assumptions. This leads to the next definition.

**Definition 34** The LCWA dependency graph of  $\mathcal{L}$  is the directed graph on  $\mathcal{R}(\Sigma)$ , containing a directed edge from predicate  $P$  to  $Q$  iff there exists  $\text{LCWA}(Q(\bar{x}), \Psi[\bar{x}]) \in \mathcal{L}$  such that  $P$  occurs in  $\Psi$ . A hierarchically closed database is a locally closed database  $\mathcal{D} = (D, \mathcal{L})$  in which the LCWA dependency graph of  $\mathcal{L}$  is cycle-free.

We may extend the LCWA-dependency graph with a query  $Q[\bar{x}]$  by adding it as a node and adding edges to it from each predicate  $P$  that occurs in  $Q[\bar{x}]$ . The transitive closure of the LCWA-dependency graph  $\mathcal{G}$  of  $\mathcal{D}$  is denoted  $\prec$ , and consists of all pairs  $(P, Q)$  such that there is a path from  $P$  to  $Q$  in  $\mathcal{G}$ . Thus, a hierarchically closed database is one such that  $\prec$  is cycle free. Let us consider also the smaller graph containing a directed edge from predicate  $P$  to  $Q$  iff there exists  $\text{LCWA}(Q(\bar{x}), \Psi[\bar{x}]) \in \mathcal{L}$  such that  $P$  occurs *negatively* in  $\Psi$ , and let us denote its transitive closure by  $\prec^-$ . As before, we may extend  $\prec^-$  for a given query by adding an edge to the query from all predicates that occur negatively in it and by closing the graph under the transitivity rule. We observe that  $\prec^-$  is cycle free iff for any query  $Q[\bar{x}]$ , the set  $\Delta_{Q, \mathcal{L}}$  (and also  $\Delta'_{Q, \mathcal{L}}$ ) is non-recursive. The reflexive closures of  $\prec$  and  $\prec^-$  are denoted  $\preceq$  and  $\preceq^-$ , respectively.

To study the accuracy of approximate query answering in locally closed databases, we need the following two definitions:

**Definition 35** Let  $\Gamma$  be a consistent theory over  $\Sigma$  containing  $\text{DCA}(\Sigma) \wedge \text{UNA}(\Sigma)$ , and let  $\text{HU}$  be the Herbrand universe of  $\Sigma$ . A query  $Q[\bar{x}]$  is squared in  $\Gamma$  if for every  $\bar{d} \in \text{HU}^n$ , if  $\Gamma \models Q[\bar{d}]$  then  $Q[\bar{d}]^{\text{Cr}} = \mathbf{t}$ .

It is easy to verify that each one of the following conditions is equivalent to the fact that  $Q[\bar{x}]$  is squared in  $\Gamma$ :

- $\forall \bar{d} \in \text{HU}^n$  if  $Q[\bar{d}]^{\text{Cr}} = \mathbf{u}$  then  $\Gamma \cup \{\neg Q[\bar{d}]\}$  is satisfiable,
- $\text{Cert}_{\mathcal{O}_\Gamma}(Q[\bar{x}]) = \text{Cert}_\Gamma(Q[\bar{x}])$ ,
- $\text{Poss}_{\mathcal{O}_\Gamma}(\neg Q[\bar{x}]) = \text{Poss}_\Gamma(\neg Q[\bar{x}])$ .

It follows that the construction of the optimal approximation  $\mathcal{O}_\Gamma$  of  $\Gamma$  allows to compute all the certain answers of queries that are squared in  $\Gamma$ .

In what follows, when  $Q[\bar{x}]$  is squared in a theory  $\Gamma = \mathcal{M}(\mathcal{D})$  (see Definition 9), we shall say that  $Q[\bar{x}]$  is squared in  $\mathcal{D}$ .

**Definition 36** For a query  $Q[\bar{x}]$ , let  $\sigma$  be the downward closed set  $\{R \mid \exists S : S \prec^- Q[\bar{x}] \text{ such that } R \preceq S\}$ .

- A predicate  $P$  is positive free in  $Q[\bar{x}]$  if  $P \notin \sigma$ .
- We denote by  $\mathcal{D}_Q$  be the extension of  $\mathcal{D} = (D, \mathcal{L})$ , obtained by adding  $\text{LCWA}(P(\bar{x}), \mathbf{t})$  to  $\mathcal{L}$ , for each predicate  $P$  that is positive free in  $Q[\bar{x}]$ .

The set  $\sigma$  in the last definition consists of all the predicates on which negatively occurring predicates in  $Q[\bar{x}]$  depend. In practice, this is often a small set in which case most predicates are positive free and  $\mathcal{D}_Q$  is almost a complete database. This is useful because databases with more base predicates have more squared queries.

**Example 37** In Example 5, all the predicates are positive free for  $\text{Tel}(x, y)$  and none of them is positive free for  $\neg \text{Tel}(x, y)$ , as  $\text{Tel}$  depends on  $\text{Dept}$ .

**Theorem 38 (Completeness)** Let  $Q[\bar{x}]$  be a query that is squared in  $\mathcal{D}_Q$  and suppose that for every predicate  $P$  that occurs negatively in  $Q$ , we have that  $P^{\mathcal{C}_Q} = P^{\mathcal{O}_Q}$ . Then  $\text{Cert}_{\mathcal{C}_Q}(Q[\bar{x}]) = \text{Cert}_{\mathcal{D}_Q}(Q[\bar{x}])$ .

Proofs of this theorem and of most of the other propositions in this section, are given in the appendix.

Note that if  $Q[\bar{x}]$  is a positive query, i.e., does not contain negative occurrences of predicates, then all predicates are positive free,  $\mathcal{D}_Q$  is a complete database and has a two-valued  $\mathcal{O}_{\mathcal{D}_Q}$ . Hence,  $Q[\bar{x}]$  is squared in  $\mathcal{D}_Q$ . By Theorem 38,  $\text{Cert}_{\mathcal{C}_Q}(Q[\bar{x}])$  is optimal.

**Example 39** Consider again Examples 5 and 29. The query  $Q = \text{Tel}(\text{BD}, 3962836)$  is positive, thus it is squared in  $\mathcal{D}_Q$ , and so Theorem 38 implies that  $\text{Cert}_{\mathcal{C}_Q}(Q)$  is optimal in this case.

To get other concrete results from Theorem 38, we need sufficient conditions for the two assumptions under which it can be applied, namely:

1. squaredness of queries, and
2. optimality of predicate approximations.

**A) proving squaredness of queries:** The propositions below provide syntactic conditions for squared queries.

**Proposition 40**  $Q[\bar{x}]$  is squared in  $\mathcal{D}_Q$  if it is of the form

$$\forall \bar{y} (C_1 \vee \dots \vee C_n),$$

where each  $C_i$  is a conjunction in which

1. each non-literal conjunct  $C_{ij}$  is an arbitrary subformula containing only base predicates of  $\mathcal{D}$  and positive free predicates in  $Q[\bar{x}]$ ,
2. each pair  $C_i, C_j$  ( $i \neq j$ ) is mutually exclusive w.r.t.  $\mathcal{C}_Q$ .<sup>5</sup>

This proposition shows that the class of squared queries includes, among others, formulas of the following forms:

1. literals and conjunctions of literals,
2. positive formulas,
3. decision tree-like formulas, in which the test formulas contain only base predicates of  $\mathcal{D}$  and the leaves consist of conjunctions of database literals and formulas containing base predicates and (positive occurrences of) positive free predicates of  $Q[\bar{x}]$ .

**Example 41** According to Proposition 40, both of the formulas in Example 37 and the query of Example 39 are squared in the database of Example 5.

To define an alternative semantical condition for squaredness, we need the following definition.

**Definition 42** A theory  $\Gamma$  is atomical iff for each 2-valued Herbrand structure  $I$  such that  $\mathcal{O}_\Gamma \leq_p I$ , it holds that  $I \models \Gamma$ .

Note that every consistent atomical theory  $\Gamma$  that includes  $\text{DCA}(\Sigma) \wedge \text{UNA}(\Sigma)$  is equivalent with the first-order theory consisting of  $\text{DCA}(\Sigma) \wedge \text{UNA}(\Sigma)$  and the set of all ground literals entailed by  $\Gamma$ .

**Proposition 43**  $Q[\bar{x}]$  is squared in  $\mathcal{D}_Q$  if  $\mathcal{D}_Q$  is atomical and  $Q[\bar{x}]$  is a formula such that all the predicates with positive and negative occurrences in  $Q[\bar{x}]$  are two-valued in  $\mathcal{C}_Q$ .

The class of queries covered by this proposition allows arbitrary quantification. Since  $\mathcal{D}_Q$  is atomical if  $Q[\bar{x}]$  is a positive query, positive queries are covered by this proposition as well.

It can be shown that  $\mathcal{D}_Q$  is atomical if the database has CWI on the windows of expertise. It appears quite natural that locally closed databases will often meet this condition.

**B) proving partial optimality of  $\mathcal{C}_Q$ :** The second assumption of Theorem 38 requires the optimality of predicate approximations. The next result reduces the problem of proving optimality of  $\mathcal{C}_Q$  in a predicate  $P$  to the problem of proving optimality of  $\mathcal{C}_Q$  in predicates  $Q$  that occur negatively in  $\Psi_P$ .

**Proposition 44** Let  $\mathcal{D}$  be a locally closed database. Then  $P^{\mathcal{C}_Q} = P^{\mathcal{O}_Q}$  if the following conditions are satisfied:

<sup>5</sup>Two conjunctions  $C_1$  and  $C_2$  are called *mutually exclusive* with respect to a three-valued structure  $\mathcal{K}$ , if  $C_1$  has a conjunct whose negation is a conjunct of  $C_2$ , and that contains only predicates that are two-valued in  $\mathcal{K}$ .

- $\Psi_P[\bar{x}]$  is squared in  $\mathfrak{D}$ ,
- for every predicate  $Q$  such that  $Q \preceq^- P$ ,  $P$  does not occur positively in  $\Psi_Q$ ,
- for every predicate  $Q$  that has negative occurrence in  $\Psi_P$ ,  $Q^{\mathcal{C}_{\mathfrak{D}}} = Q^{\mathcal{O}_{\mathfrak{D}}}$ .

By induction on the depth of the strict (well-founded) order  $\preceq^-$  on  $\{Q \mid Q \preceq^- P\}$ , using Proposition 44, the following way of showing optimality on  $P$  is obtained:

**Proposition 45** *Let  $\mathfrak{D}$  be a locally closed database. Then  $P^{\mathcal{C}_{\mathfrak{D}}} = P^{\mathcal{O}_{\mathfrak{D}}}$  if for every predicate  $Q$  such that  $Q \preceq^- P$ , it holds that  $Q \not\prec^- Q$ ,  $\Psi_Q[\bar{x}]$  is squared in  $\mathfrak{D}$ , and  $Q$  does not occur positively in  $\Psi_S$  for any  $S \preceq^- Q$ .*

Note that the conditions of the proposition are satisfied if  $P$ 's window of expertise is a positive formula. It is also satisfied if for every  $Q \preceq^- P$ ,  $\Psi_Q[\bar{x}]$  is squared in  $\mathfrak{D}$  and  $Q \not\prec^- Q$ , i.e., the LCWA dependency graph does not cycle through  $Q$ .

**C) summary:** By Theorem 38 and Proposition 45 we get the following theorem for the accuracy of approximative query answering.

**Theorem 46 (Completeness)** *Let  $\mathcal{Q}[\bar{x}]$  be a query that is squared in  $\mathfrak{D}$  and suppose that for every predicate  $P$  that occurs negatively in  $\mathcal{Q}$  it holds that  $P \not\prec^- P$ ,  $\Psi_P[\bar{x}]$  is squared in  $\mathfrak{D}$ , and  $P$  does not occur positively in  $\Psi_Q$  for any  $Q \preceq^- P$ . Then  $\text{Cert}_{\mathcal{C}_{\mathfrak{D}}}(\mathcal{Q}[\bar{x}]) = \text{Cert}_{\mathfrak{D}}(\mathcal{Q}[\bar{x}])$ .*

**Example 47** Consider again the database instance of Examples 1 and 5, together with the following LCWAs:

- $\text{LCWA}(\text{Tel}(x, y), \text{Dept}(x, \text{CS}))$ ,
- $\text{LCWA}(\text{Dep}(x, y), y = \text{CS} \wedge \exists z \text{Tel}(x, z))$ .

The first assumption is considered in Examples 2 and 5; the second assumption expresses complete knowledge of the database about all the people of the computer science department that have a telephone.

This locally closed database is circular, as the window of expertise in each assumption mentions an object predicate of the other assumption. Yet, Theorem 46 is applicable here:

- The windows of expertise are positive formulas, hence the contained predicates are positive free. It follows that both are squared in  $\mathfrak{D}$ . Also,  $\prec^-$  is the empty graph, hence it is acyclic. It follows that  $\mathcal{C}_{\mathfrak{D}}$  is optimal.
- The queries  $\text{Tel}(x, y)$  and  $\neg \text{Tel}(x, y)$  are literals, hence they are squared. It follows that the approximate methods compute the optimal answers for certain and possible answers of this query.

So let us compute the answers for  $\text{Tel}(x, y)$  by our approximate methods, using the following set  $\Delta_{\mathcal{Q}, \mathcal{L}}$  of certain rules (implicit universal quantification is assumed here):

$$\left. \begin{array}{l} Q^c(x, y) \leftarrow \text{Tel}^c(x, y). \\ Q^{c\neg}(x, y) \leftarrow \text{Tel}^{c\neg}(x, y). \\ \text{Tel}^c(x, y) \leftarrow \text{Tel}(x, y). \\ \text{Tel}^{c\neg}(x, y) \leftarrow \neg \text{Tel}(x, y) \wedge \text{Dept}^c(x, \text{CS}). \\ \text{Dept}^c(x, y) \leftarrow \text{Dept}(x, y). \\ \text{Dept}^{c\neg}(x, y) \leftarrow \neg \text{Dept}(x, y) \wedge y = \text{CS} \wedge \exists z \text{Tel}^c(x, z). \end{array} \right\}$$

Interestingly, this set of rules is non-recursive. It follows that its least and greatest fixpoint coincide. As expected, the

expression  $[\text{lfp}_{Q^c, \Delta_{\mathcal{Q}, \mathcal{L}}}] (x, y)$  for  $\mathcal{Q} = \text{Tel}(x, y)$ , gives the certain answers

Tel(LD, 6531421), Tel(LD, 09-23314),  
Tel(BD, 5985625), Tel(TD, 5845213),

which are the database tuples of  $\text{Tel}$ .

For the certainly false tuples of the query, we unfold the expression  $[\text{lfp}_{Q^{c\neg}, \Delta_{\mathcal{Q}, \mathcal{L}}}] (x, y)$ , yielding the database query  $\neg \text{Tel}(x, y) \wedge \text{Dept}(x, \text{CS})$ . Its answer is the following set:

$\{\text{Tel}(\text{BD}, y) \mid y \neq 5985625\} \cup \{\text{Tel}(\text{TD}, y) \mid y \neq 5845213\}$ .

The set of possible but uncertain answers is the complement of the union of the certain answers and the certainly false answers. This set is specified by the database query  $\neg \text{Tel}(x, y) \wedge (\text{Tel}(x, y) \vee \neg \text{Dept}(x, \text{CS}))$ , which can be simplified to  $\neg \text{Tel}(x, y) \wedge \neg \text{Dept}(x, \text{CS})$ . Its answer is the set:

$\{\text{Tel}(\text{LD}, y) \mid y \neq 6531421 \wedge y \neq 09-23314\} \cup \{\text{Tel}(\text{DF}, y)\}$ .

As noted above, Theorem 46 guarantees that these query answers are precise.

To summarize, the results in this section allow to prove the optimality of the approximate certain answers in the context of queries  $\mathcal{Q}[\bar{x}]$  and (a subset of the) windows of expertise in the form of Proposition 40: e.g., positive formulas, or conjunctions of literals, or decision-tree like formulas with base predicates in the tests, etc. Even larger classes of queries can be optimally answered with respect to atomically locally closed databases (Proposition 43). Clearly, this is a rich class of queries and a rich class of databases, for example allowing many forms of cycles in the LCWA-dependency graph. In contrast, the optimality theorem in (Cortés-Calabuig et al. 2007) is for queries in the form of conjunctions of literals and on hierarchically closed databases with conjunctions of literals in the windows of expertise.

We currently lack experience to evaluate the precision of the approximate query answering methods beyond the conditions in the above optimality propositions. On one hand, we have constructed examples showing a drastic loss of precision of the approximate answers. On the other hand, an analysis of the optimality proofs suggests that in many other applications where the conditions do not hold, the approximate methods should still be quite precise.

## Related Works

Incompleteness in relational databases has been investigated almost since their inception back in the seventies. This issue is continuously arising whenever a new database paradigm is introduced. The general problem of dealing with incompleteness in relational databases has already been discussed in (Imielinski and Lipski 1981) and in (Grahne 1984). Reiter (1986) provides an early semantic characterization of a database containing null values, and defines a sound algorithm for querying such databases. At the representation level, Motro (1989) is perhaps the first to introduce a language – similar in spirit to ours – to represent partial completeness using logical views. This is followed by



the works of Levy (1996), Etzioni et al. (1997), and Doherty et al. (2000). Incompleteness of databases is also addressed in (Grahne 2002) from a data integration perspective. In (Grahne and Mendelson 1999) a framework for dealing with incompleteness in mediator-based systems is introduced, based on tableaux techniques for query answering. Incompleteness in data exchange has recently been considered in (Libkin 2006).

Evidently, the concept of a locally closed database has strong ties also to non-monotonic reasoning, so there is no wonder that it can be expressed by a variety of non-monotonic formalisms. A detailed analysis on the relationship between the LCWA and circumscription is provided in (Cortés-Calabuig et al. 2005) and a discussion on how to express LCWAs using (general) logic programs as presented in (Gelfond et al. 1991) is given in (Cortés-Calabuig et al. 2006). The idea of using logic programs for querying incomplete databases is also investigated in (Baral et al. 1998) and (Bonatti et al. 1996). In the former, incompleteness is expressed by extending a relational database to a set of literals, and it is shown how queries are expanded from complete databases to be applicable to incomplete databases, using general logic programs. In (Bonatti et al. 1996) incomplete information is represented by disjunctive databases, and query answering is defined by providing semantics to such databases using an extension of logic programs that is based on autoepistemic logic.

## Conclusions

The ability to correctly and efficiently reason with partially complete databases is a major goal whose importance is obvious. However, as we have shown, the corresponding decision problems are not tractable and sometimes are even not decidable. To handle this, we introduced a rewriting technique to compute certain or possible answers in polynomial time from such databases. Our approach is based on approximating all models of the database's theory by means of three-valued structures, which are implicitly represented by fixpoint queries. We have shown that this approach is applicable on hierarchically closed databases, an extended notion of the one introduced in (Cortés-Calabuig et al. 2007). Moreover, in this new setting, certain types of important integrity constraints can be incorporated in the framework.

We have also provided a 'toolbox' to prove optimality of the approximate query answering methods for quite a broad class of queries and databases. Our results suggest that the approximate reasoning methods may often be quite precise, and frequently optimal. We believe that, in practice, the approximate methods may offer very often a high degree of accuracy, even when applied beyond the conditions that guarantee optimality.

Finally, we view the current work as a further step towards the more ambitious goal of providing a unifying framework for efficient approximative query answering in incomplete databases in the presence of arbitrary integrity constraints. Towards this goal, efficient reasoning methods developed for first order logic theories as the one recently investigated in (Wittocx et al. 2008), seem a promising path to explore.

In the other direction, we believe the results concerning approximative reasoning presented in this work and the way they were obtained can be adapted to other contexts in which efficient query answering is imperative, such as OWL ontologies or databases of a less structured nature (XML).

## Acknowledgments

This work was in part funded by the projects GOA/2003/08 "Inductive Knowledge Bases" and G.0357.06 "Design, implementation and application of model generation techniques for ID-logic", Research Foundation - Flanders (FWO-Vlaanderen).

## References

- Bancilhon, F.; Maier, D.; Sagiv, Y.; and Ullman, J. 1986. Magic sets and other strange ways to implement logic programs. In *Proc. 5th PODS*, 1–15. ACM Press.
- Baral, C.; Gelfond, M.; and Koshelva, O. 1998. Expanding queries to incomplete databases by interpolating general logic programs. *J. Log. Program.* 35(3):195–230.
- Bonatti, P.; and Eiter, T. 1996. Querying disjunctive databases through nonmonotonic logics. *Theor. Comput. Sci.* 160(1–2):321–363.
- Cortés-Calabuig, A.; Denecker, M.; Arieli, O.; Van Nuffelen, B.; and Bruynooghe, M. 2005. On the local closed-world assumption of data-sources. In *Proc. 8th LPNMR*, LNCS 3662, 145–157. Springer.
- Cortés-Calabuig, A.; Denecker, M.; Arieli, O.; and Bruynooghe, M. 2006. Representation of partial knowledge and query answering in locally complete databases. In *Proc. 13th LPAR*, LNCS 4246, 407–421. Springer.
- Cortés-Calabuig, A.; Denecker, M.; Arieli, O.; and Bruynooghe, M. 2007. Approximate query answering in locally closed databases. In *Proc. 22nd AAI*, 397–402. AAI Press.
- Denecker, M.; and Ternovska, E. 2007. Inductive situation calculus. *Artificial Intelligence* 171(5-6):332–360.
- Doherty, P.; Łukaszewicz, W.; and Szalas, A. 2000. Efficient reasoning using the local closed-world assumption. In *Proc. 9th AIMS*, LNCS 2407, 49–58. Springer.
- Etzioni, O.; Golden, K.; and Weld, D. 1997. Sound and efficient closed-world reasoning for planning. *Artificial Intelligence* 89(1–2):113–148.
- Gelfond, M.; Lifschitz, V. 1991. Classical negation in logic programs and disjunctive databases. *New Generation Computing* 9:367–387.
- Grahne, G. 1984. Dependency satisfaction in databases with incomplete information. In *Proc. 10th VLDB*, 37–45. Morgan Kaufmann.
- Grahne, G. 2002. Information integration and incomplete information. *IEEE Data Engineering Bulletin* 25(3):46–52.
- Grahne, G., Mendelson A. 1999. Tableau techniques for querying information sources through global schemas. In *Proc. 7th ICDT*, LNCS 1540, 332–347. Springer.

Imielinski, T.; and Lipski Jr, W. 1981. On representing incomplete information in a relational data base. In *Proc. 7th VLDB*, 388–397. IEEE Press.

Levy, A. 1996. Obtaining complete answers from incomplete databases. In *Proc. 22nd VLDB*, 402–412. Morgan Kaufmann.

Libkin, L. 2006. Data exchange and incomplete information. In *Proc. 25th PODS*, 60–69. ACM Press.

Motro, A. 1989. Integrity = Validity + Completeness. *ACM Trans. Database Syst.* 14(4):480–502.

Reiter, R. 1982. Towards a logical reconstruction of relational database theory. In *On Conceptual Modelling (Inter-vale Workshop)*, 191–233.

Reiter, R. 1986. A sound and sometimes complete query evaluation algorithm for relational databases with null values. *J. ACM* 33(2):349–370.

Swift, T. 1999. Tabling for non-monotonic programming. *Annals of Mathematics and Artificial Intelligence* 25(3-4):201–240.

Trakhtenbrot, B. 1963. Impossibility of an algorithm for the decision problem in finite classes. *American Mathematical Society Transaction* 3(2):1–5.

Van Gelder, A. 1993. The alternating fixpoint of logic programs with negation. *Journal of Computer and System Sciences* 47(1):185–221.

Wittocx, J.; Mariën, M.; and Denecker, M. 2008. Grounding with bounds. In *Proc. 23rd AAAI*, 572–577. AAAI Press.

## Appendix: Proofs

**Proof of Proposition 12.** Obviously, when  $\mathfrak{D}$  has complete information about  $Q[\bar{x}]$ , certain and possible answers coincide, i.e.,  $Cert_{\mathfrak{D}}(Q[\bar{x}]) = Poss_{\mathfrak{D}}(Q[\bar{x}])$ . Thus, since  $D$  is a model of  $\mathfrak{D}$ , we have:  $\{\bar{a} \mid D \models Q[\bar{a}]\} \subseteq Poss_{\mathfrak{D}}(Q[\bar{x}]) = Cert_{\mathfrak{D}}(Q[\bar{x}]) \subseteq \{\bar{a} \mid D \models Q[\bar{a}]\}$ . Conversely, if  $Cert_{\mathfrak{D}}(Q[\bar{x}]) = Poss_{\mathfrak{D}}(Q[\bar{x}])$ , then obviously either  $\mathfrak{D} \models Q[\bar{x}]$  or  $\mathfrak{D} \models \neg Q[\bar{x}]$ , i.e. there is CWI on  $Q[\bar{x}]$ .  $\square$

**Proof of Proposition 13.** Consider the query  $Q[\bar{x}] = P(c)$  and  $\mathcal{L} = \{\mathcal{LCWA}(P(c), \varphi)\}$ , where  $\varphi$  is a sentence not containing  $P$ . We observe that a database  $(D, \mathcal{L})$  has no CWI on  $P(c)$  iff  $\neg\varphi$  has a finite model. It follows that there is CWI on  $P(c)$  in all databases  $(D, \mathcal{L})$  iff  $\varphi$  is satisfied in all finite structures. This is a validity checking problem of a first-order formula with respect to the class of finite structures. By Trakhtenbrot’s theorem (Trakhtenbrot 1963), this problem is undecidable.  $\square$

**Proof outline of Proposition 30.** By induction on the number of iterations in the computations of the operators  $App_{\mathfrak{D}}$  (Definition 22) and  $\Gamma_{\Delta_{\mathfrak{Q}, \mathcal{L}}}$  (Definition 28), one shows that the structure that is obtained by the latter in a certain iteration simulates (in the sense of Definition 14) the structure

that is obtained by the former in the same iteration. Suppose now that  $App_{\mathfrak{D}}$  reaches a fixpoint after  $\alpha$  iterations. This fixpoint is  $\mathcal{C}_{\mathfrak{D}}$ . The structure  $I_{\alpha}$  obtained by  $\Gamma_{\Delta_{\mathfrak{Q}, \mathcal{L}}}$  at this iteration simulates  $\mathcal{C}_{\mathfrak{D}}$ , and it is a fixpoint on all the predicates  $P_i^c$  and  $P_i^{c\bar{c}}$ . After one more iteration  $\Gamma_{\Delta_{\mathfrak{Q}, \mathcal{L}}}$  reaches a fixpoint also on the predicates  $Q^c$  and  $Q^{c\bar{c}}$ . By Proposition 16, it holds that  $\bar{d} \in \mathcal{R}_{\mathfrak{Q}}^c$  iff  $(Q[\bar{d}]^c)^{I_{\alpha}} = \mathbf{t}$  iff  $(Q[\bar{d}])^{\mathcal{C}_{\mathfrak{D}}} = \mathbf{t}$ . Likewise,  $\bar{d} \in \mathcal{R}_{\mathfrak{Q}}^{c\bar{c}}$  iff  $(Q[\bar{d}])^{\mathcal{C}_{\mathfrak{D}}} = \mathbf{f}$ .  $\square$

**Proof of Theorem 38.** First, we show the following lemmas:

**Lemma 48** *Let  $\mathcal{K}$  be an approximation of  $\mathfrak{D}$  such  $P(\bar{a})^{\mathcal{K}} = \mathbf{t}$  iff  $P(\bar{a}) \in D$ . Let  $Q[\bar{x}]$  be a query squared in  $\mathfrak{D}$  such that for every database predicate  $P$  that occurs negatively in  $Q$ ,  $P^{\mathcal{K}} = P^{\mathcal{O}_{\mathfrak{D}}}$ . It holds that  $Cert_{\mathfrak{D}}(Q[\bar{x}]) = Cert_{\mathcal{K}}(Q[\bar{x}])$ .*

*Proof.* We have that  $Cert_{\mathcal{K}}(Q[\bar{x}]) \subseteq Cert_{\mathfrak{D}}(Q[\bar{x}])$ . To show the other direction, assume that  $\mathfrak{D} \models Q[\bar{d}]$ . Since  $Q[\bar{x}]$  is squared in  $\mathfrak{D}$ , it follows that  $Q[\bar{d}]^{\mathcal{O}_{\mathfrak{D}}} = \mathbf{t}$ . For every atom  $P(\bar{a})$  of a predicate that occurs in  $Q[\bar{x}]$ , if  $P(\bar{a})^{\mathcal{K}} \neq P(\bar{a})^{\mathcal{O}_{\mathfrak{D}}}$ , then  $P$  occurs only positively in  $Q[\bar{x}]$  and, since  $\mathcal{K} \leq_p \mathcal{O}_{\mathfrak{D}}$ , it holds that  $P(\bar{a})^{\mathcal{K}} = \mathbf{u}$  and  $P(\bar{a})^{\mathcal{O}_{\mathfrak{D}}} \neq \mathbf{u}$ . By the assumption we made about  $\mathcal{K}$ , it holds that  $P(\bar{a})^{\mathcal{O}_{\mathfrak{D}}} = \mathbf{t}$  iff  $P(\bar{a}) \in D$  iff  $P(\bar{a})^{\mathcal{K}} = \mathbf{t}$ . Hence, if  $P(\bar{a})^{\mathcal{K}} \neq P(\bar{a})^{\mathcal{O}_{\mathfrak{D}}}$  then  $P(\bar{a})^{\mathcal{K}} = \mathbf{u}$  and  $P(\bar{a})^{\mathcal{O}_{\mathfrak{D}}} = \mathbf{f}$ . It follows from a well-known monotonicity property of the standard three-valued truth assignment, that  $Q[\bar{d}]^{\mathcal{K}} \geq Q[\bar{d}]^{\mathcal{O}_{\mathfrak{D}}} = \mathbf{t}$ .

**Lemma 49** *If  $Q[\bar{x}]$  is squared in  $\mathfrak{D}_{\mathfrak{Q}}$ , then  $Q[\bar{x}]$  is squared in  $\mathfrak{D}$ .*

*Proof.* Note that the models of  $\mathfrak{D}_{\mathfrak{Q}}$  correspond exactly to the models  $M'$  of  $\mathfrak{D}$  such that  $P^{M'} = P^D$ , for every positive free predicate  $P$  of  $Q[\bar{x}]$ . Also, for every model  $M$  of  $\mathfrak{D}$ , the structure  $M'$  obtained by setting  $P^{M'} = P^D$  if  $P$  is positive free and  $P^{M'} = P^M$  otherwise, is a model of  $\mathfrak{D}_{\mathfrak{Q}}$ . It follows that  $Q[\bar{d}]^{M'} \leq Q[\bar{d}]^M$ , as for all predicates  $P$  that occur negatively in  $Q[\bar{x}]$ , it holds that  $P^{M'} = P^M$  while predicates  $P$  with only positive occurrences have  $P^{M'} \subseteq P^M$ . Thus,  $\mathfrak{D} \models Q[\bar{d}]$  iff  $\mathfrak{D}_{\mathfrak{Q}} \models Q[\bar{d}]$ , and so  $Cert_{\mathfrak{D}_{\mathfrak{Q}}}(Q[\bar{x}]) = Cert_{\mathfrak{D}}(Q[\bar{x}])$ .

Observe, also, that  $\mathcal{O}_{\mathfrak{D}}$  is an approximation of  $\mathfrak{D}_{\mathfrak{Q}}$  that satisfies the condition on  $\mathcal{K}$  in Lemma 48. It follows that  $Cert_{\mathcal{O}_{\mathfrak{D}}}(Q[\bar{x}]) = Cert_{\mathfrak{D}_{\mathfrak{Q}}}(Q[\bar{x}])$  and so, since  $Q[\bar{x}]$  is squared in  $\mathfrak{D}_{\mathfrak{Q}}$ ,  $Cert_{\mathcal{O}_{\mathfrak{D}}}(Q[\bar{x}]) = Cert_{\mathfrak{D}_{\mathfrak{Q}}}(Q[\bar{x}])$ . By what we have obtained above, it follows that  $Cert_{\mathcal{O}_{\mathfrak{D}}}(Q[\bar{x}]) = Cert_{\mathfrak{D}}(Q[\bar{x}])$ , which implies that  $Q[\bar{x}]$  is squared in  $\mathfrak{D}$ .

To complete the proof of Theorem 38, let  $Q[\bar{x}]$  be a query that is squared in  $\mathfrak{D}_{\mathfrak{Q}}$ . By Lemma 49,  $Q[\bar{x}]$  is squared in  $\mathfrak{D}$ , and so by Lemma 48 for  $\mathcal{K} = \mathcal{C}_{\mathfrak{D}}$  the theorem is obtained.  $\square$

We now turn to Propositions 40 and 43 that identify conditions for squaredness. For this, we need the following definitions:

**Definition 50** *The supervaluation  $sv_{\mathcal{K}}(\varphi)$  of a sentence  $\varphi$  with respect to a three-valued structure  $\mathcal{K}$  is defined as*

$$sv_{\mathcal{K}}(\varphi) = lub_{\leq_p} \{\varphi^M \mid \mathcal{K} \leq_p M\}.$$

where  $M$  ranges over two-valued structures.

**Definition 51** A query  $\mathcal{Q}[\bar{x}]$  is literal-based in a satisfiable theory  $\Gamma$  containing  $DCA \wedge UNA$  iff

$$Cert_{\Gamma}(\mathcal{Q}[\bar{x}]) = \{\bar{d} \in HU^n \mid sv_{\mathcal{O}_{\Gamma}}(\mathcal{Q}[\bar{d}]) = \mathbf{t}\}.$$

$\mathcal{Q}[\bar{x}]$  is Kleene-precise in a three-valued structure  $\mathcal{K}$  with domain  $Dom$  iff

$$Cert_{\mathcal{K}}(\mathcal{Q}[\bar{x}]) = \{\bar{d} \in Dom^n \mid sv_{\mathcal{K}}(\mathcal{Q}[\bar{d}]) = \mathbf{t}\}.$$

Again, in the sequel, if  $\mathcal{Q}[\bar{x}]$  is literal-based in a theory  $\Gamma = \mathcal{M}(\mathfrak{D})$  (see Definition 9), we shall sometimes say that  $\mathcal{Q}[\bar{x}]$  is literal-based in  $\mathfrak{D}$ .

Note that since a supervaluation is more precise than standard Kleene truth assignment, it follows that the certain answers of a Kleene-precise query under standard Kleene truth assignment coincides with the supervaluation.

**Lemma 52** A query  $\mathcal{Q}[\bar{x}]$  is squared in  $\Gamma$  iff  $\mathcal{Q}[\bar{x}]$  is literal-based in  $\Gamma$  and Kleene-precise in  $\mathcal{O}_{\Gamma}$ .

*Proof.* For every  $\Gamma$  and  $\mathcal{Q}[\bar{x}]$ , the following inequalities hold:

$$\begin{aligned} Cert_{\mathcal{O}_{\Gamma}}(\mathcal{Q}[\bar{x}]) &\subseteq \{\bar{d} \in HU^n \mid sv_{\mathcal{O}_{\Gamma}}(\mathcal{Q}[\bar{d}]) = \mathbf{t}\} \\ &\subseteq Cert_{\Gamma}(\mathcal{Q}[\bar{x}]). \end{aligned}$$

If  $\mathcal{Q}[\bar{x}]$  is literal-based in  $\Gamma$  and Kleene-precise in  $\mathcal{O}_{\Gamma}$ , the inequalities turn into equalities. Vice versa, if  $\mathcal{Q}[\bar{x}]$  is squared, the three terms are equal, and it follows that the query is literal-based and Kleene-precise.  $\square$

By Lemma 52, then, for showing Propositions 40 and 43, it is enough to show that  $\mathcal{Q}[\bar{x}]$  is literal-based in  $\mathfrak{D}_{\mathcal{Q}}$  and is Kleene-precise in  $\mathcal{O}_{\mathfrak{D}_{\mathcal{Q}}}$ .

**Proof of Proposition 40.** Let  $\Gamma = \mathcal{M}(\mathfrak{D}_{\mathcal{Q}})$ . Since  $\Gamma$  contains  $DCA(\Sigma) \wedge UNA(\Sigma)$ , the formula  $\forall x \varphi[\bar{x}]$  will be literal-based in  $\Gamma$  if  $\varphi[\bar{x}]$  is literal-based in  $\Gamma$ . Hence, it suffices to prove the proposition for the case that  $\mathcal{Q}[\bar{x}]$  is quantifier free.

Now, given that all predicates of non-literal conjuncts  $C_{ij}$  are two-valued in  $\mathcal{O}_{\Gamma}$ , it is obvious that  $\Gamma \models C_{ij}[\bar{d}]$  iff  $sv_{\mathcal{O}_{\Gamma}}(C_{ij}[\bar{d}]) = \mathbf{t}$ , iff  $C_{ij}[\bar{d}]^{\mathcal{O}_{\Gamma}} = \mathbf{t}$  for each  $\bar{d} \in HU^n$ .

Since each pair  $C_i[\bar{x}], C_j[\bar{x}]$  of  $\mathcal{Q}[\bar{x}]$  is mutually exclusive in  $\mathcal{O}_{\Gamma}$ , it follows that for all  $\bar{d}$ , there is at most one conjunct  $C_i$  such that  $\Gamma \models C_i[\bar{d}]$ . It is then easy to verify that

$$\begin{aligned} \Gamma \models C_i[\bar{d}] &\text{ iff } \Gamma \models C_{ij}[\bar{d}], \text{ for each conjunct } C_{ij} \text{ of } C_i, \\ &\text{ iff } sv_{\mathcal{O}_{\Gamma}}(C_{ij}[\bar{d}]) = \mathbf{t} \text{ for each conjunct } C_{ij} \\ &\text{ iff } (C_{ij}[\bar{d}])^{\mathcal{O}_{\Gamma}} = \mathbf{t} \text{ for each conjunct } C_{ij} \\ &\text{ iff } sv_{\mathcal{O}_{\Gamma}}(C_i[\bar{d}]) = \mathbf{t} \\ &\text{ iff } (C_i[\bar{d}])^{\mathcal{O}_{\Gamma}} = \mathbf{t}. \end{aligned}$$

Thus  $\mathcal{Q}[\bar{x}]$  is literal-based in  $\Gamma$  and Kleene-precise in  $\mathcal{O}_{\Gamma}$ .  $\square$

**Proof outline of Proposition 43.** As  $\mathfrak{D}_{\mathcal{Q}}$  is atomical, every  $\mathcal{Q}[\bar{x}]$  is literal-based in it. By Lemma 52 it remains to show that  $\mathcal{Q}[\bar{x}]$  is Kleene-precise in  $\mathcal{O}_{\mathfrak{D}_{\mathcal{Q}}}$ . All we need to do is to show that if  $\mathcal{Q}[\bar{d}]^{\mathcal{O}_{\mathfrak{D}_{\mathcal{Q}}}} = \mathbf{u}$  then  $sv_{\mathcal{O}_{\mathfrak{D}_{\mathcal{Q}}}}(\mathcal{Q}[\bar{d}]) = \mathbf{u}$ . We construct two-valued structures  $K$  and  $K'$  from  $\mathcal{O}_{\mathfrak{D}_{\mathcal{Q}}}$  where  $K$  maps each unknown atom  $P(\bar{a})$  of  $\mathcal{O}_{\mathfrak{D}_{\mathcal{Q}}}$  to  $\mathbf{t}$  if  $P$  occurs positively in  $\varphi$  and to  $\mathbf{f}$  otherwise; likewise,  $K'$  follows the inverse strategy and maps unknown atoms  $P(\bar{a})$  to  $\mathbf{f}$  if  $P$

occurs positively in  $\varphi$  and to  $\mathbf{t}$  otherwise. By a routine induction, one shows that if  $\mathcal{Q}[\bar{d}]^{\mathcal{O}_{\mathfrak{D}_{\mathcal{Q}}}} = \mathbf{u}$ , then  $\mathcal{Q}[\bar{d}]^K = \mathbf{t}$  and  $\mathcal{Q}[\bar{d}]^{K'} = \mathbf{f}$ , and hence,  $sv_{\mathcal{O}_{\mathfrak{D}_{\mathcal{Q}}}}(\mathcal{Q}[\bar{d}]) = \mathbf{u}$ .  $\square$

**Proof of Proposition 44.** By the soundness of  $\mathcal{C}_{\mathfrak{D}}$ , it holds for every  $\bar{a} \in Dom^n$ , that if  $P(\bar{a})^{\mathcal{C}_{\mathfrak{D}}} \neq \mathbf{u}$ , then  $P(\bar{a})^{\mathcal{C}_{\mathfrak{D}}} = P(\bar{a})^{\mathcal{O}_{\mathfrak{D}}}$ . So, let us assume that  $P(\bar{a})^{\mathcal{C}_{\mathfrak{D}}} = \mathbf{u}$ . Observe that, in this case,  $P(\bar{a}) \notin D$ .

To show that  $P(\bar{a})^{\mathcal{O}_{\mathfrak{D}}} = \mathbf{u}$ , we need to construct two models  $M, M'$  of  $\mathfrak{D}$  such that  $M \models P(\bar{a})$  and  $M' \models \neg P(\bar{a})$ . Since  $P(\bar{a}) \notin D$ , we can take  $M' = D$  which is indeed a model of  $\mathfrak{D}$ . Let us now construct  $M$ . By construction of  $\mathcal{C}_{\mathfrak{D}}$ , it holds that  $P(\bar{a}) \notin D$  and  $\Psi_P[\bar{a}]^{\mathcal{C}_{\mathfrak{D}}} \neq \mathbf{t}$ . The sentence  $\Psi_P[\bar{a}]$  is squared in  $\mathfrak{D}$  and by assumption, it holds that  $Q^{\mathcal{C}_{\mathfrak{D}}} = Q^{\mathcal{O}_{\mathfrak{D}}}$  for each negatively occurring predicate  $Q$  in this formula. Therefore, the conditions of Proposition 48 are satisfied, so there is a model  $N$  of  $\mathfrak{D}$  such that  $N \models \neg \Psi_P[\bar{a}]$ . If  $N \models \neg P(\bar{a})$ , then we can take  $M = N$ . Otherwise, we transform  $N$  into a model  $M$  in which  $P(\bar{a})$  is false.

Consider the set  $S_P = \{Q \in \mathcal{R}(\Sigma) \mid \text{for some } Q' \preceq^- Q, P \in^+ \Psi_{Q'}\}$ . It consists of all predicates in whose window of expertise  $P$  has a positive occurrence, and all predicates that negatively depend on these. By the condition of the proposition, it holds that  $P \notin S_P$ .

Define  $M$  as the structure obtained by modifying  $N$  as follows:

- $P^M = P^N \cup \{\bar{a}\}$ , i.e.,  $P(\bar{a})$  is made true;
- $Q^M = Q^D$  for  $Q \in S_P$ .

This modification increases  $P$  and decreases all predicates of  $S_P$ ; i.e.,  $P^N \leq P^M$ ,  $Q^M \leq Q^N$  for  $Q \in S_P$ , and  $Q^N = Q^M$  otherwise. Thus, formulas with only positive occurrences of  $P$  and only negative occurrences of predicates  $Q \in S_P$  have a larger truth value in  $M$  than in  $N$ .

To verify that  $M$  is a model of  $\mathfrak{D}$ , it suffices to check that  $M$  satisfies all local closed world assumptions. Consider any instance of a local closed world assumption:

$$\varphi \equiv \neg \Psi_Q[\bar{d}] \vee \neg Q(\bar{d}) \vee (Q(\bar{d}) \in D)$$

Each of these formulas is satisfied in  $N$ . Let us verify that it is satisfied in  $M$  as well. There are four cases:

- $Q = P$  and  $\bar{d} = \bar{a}$ : in this case,  $N \models \neg \Psi_P[\bar{a}]$ . The formula  $\neg \Psi_P[\bar{a}]$  contains only positive occurrences of  $P$  and only negative occurrences of predicates  $Q \in S_P$ , hence  $\mathbf{t} = (\neg \Psi_P[\bar{a}])^N \leq (\neg \Psi_P[\bar{a}])^M$ .
- $Q = P$  and  $\bar{d} \neq \bar{a}$ : we have  $(\neg P(\bar{d}) \vee (P(\bar{d}) \in D))^N = (\neg P(\bar{d}) \vee (P(\bar{d}) \in D))^M$  and  $\neg \Psi_P[\bar{a}]$  contains only positive occurrences of  $P$  and only negative occurrences of predicates  $Q \in S_P$ . It follows that  $\mathbf{t} = \varphi^N \leq \varphi^M$ .
- $Q \in S_P$ :  $M$  satisfies  $\neg Q(\bar{d}) \vee (Q(\bar{d}) \in D)$ .
- $Q \notin S_P$  and  $Q \neq P$ :  $\varphi$  contains only positive occurrences of  $P$  and only negative occurrences of predicates of  $S_Q$ ; hence  $\mathbf{t} = \varphi^N \leq \varphi^M$ .  $\square$