

Best Answers over Incomplete Data : Complexity and First-Order Rewritings

Amélie Gheerbrant and Cristina Sirangelo

Université de Paris, IRIF, CNRS, F-75013 Paris, France

{amelie, cristina}@irif.fr

Abstract

Answering queries over incomplete data is ubiquitous in data management and in many AI applications that use query rewriting to take advantage of relational database technology. In these scenarios one lacks full information on the data but queries still need to be answered with certainty. The certainty aspect often makes query answering unfeasible except for restricted classes, such as unions of conjunctive queries. In addition often there are no, or very few, certain answers, thus expensive computation is in vain. Therefore we study a relaxation of certain answers called best answers. They are defined as those answers for which there is no better one (that is, no answer true in more possible worlds). When certain answers exist the two notions coincide. We compare different ways of casting query answering as a decision problem and characterise its complexity for first-order queries, showing significant differences in the behaviour of best and certain answers. We then restrict attention to best answers for unions of conjunctive queries and produce a practical algorithm for finding them based on query rewriting techniques.

1 Introduction

Answering queries over incomplete databases is crucial in many different scenarios such as data integration [Lenzerini, 2002], data exchange [Arenas *et al.*, 2014], inconsistency management [Bertossi, 2011], data cleaning [Geerts *et al.*, 2013], ontology-based data access (OBDA) [Bienvenu and Ortiz, 2015], and many others. The common thread running through all these applications lies in computing *certain answers* [Amendola and Libkin, 2018; Libkin, 2018a], which is the standard way of answering queries over incomplete databases. Intuitively this produces answers that can be obtained from all the possible complete databases a given incomplete database represents. However, computing such query answers then relies on sophisticated algorithms that are often difficult to implement. It is well known that restricting to unions of conjunctive queries allows to overcome the difficulty by using *naïve evaluation* [Imieliński and Lipski, 1984]. This amounts to evaluating queries over incomplete databases

as if nulls were usual data values, thus merely using the standard database query engine to compute certain answers.

In general though it is a common occurrence that few if any certain answers can be found. If there are no certain answers, it is still useful to provide a user with some answers, with suitable guarantees. To address this need, a framework to measure how close an answer is to certainty has recently been proposed [Libkin, 2018b], setting the foundations to both a quantitative and a qualitative approach. We focus on the qualitative notion of *best answers*. Those are a refinement of certain answers based on comparing query answers; one that is supported by a larger set of complete interpretations is better. Best answers are those answers for which there is no better one. They always exist and when certain answers exist the two notions simply coincide.

Best answers is a natural notion, but we still know little about it. Identifying the set of best answers among some given family of sets of answers is known to be complete in data complexity for the class $P^{NP[\log n]}$, which is considered as “mildly” harder than both NP and CONP. However this very formulation of the decision problem is non standard. Traditionally, in databases we rather focus on problems stating that some given result belongs to the set of answers, or that some given set is the set of answers. Certain answers as a decision problem is typically formulated in the first way. So do these variations matter? We fully answer the question for both best and certain answers of first-order queries, showing significant differences in their computational behaviour.

Despite the high complexity of finding best answers in general, one gains tractability when restricting to unions of conjunctive queries. This is a common class of queries, usually well behaved computationally, even for certain answers. Finding best answers for them was shown to be tractable in [Libkin, 2018b] via an adaptation of techniques used in the context of certain answers [Gheerbrant and Libkin, 2015]. Those are essentially resolution based algorithms for first-order formulas; this makes them hard to implement in the database context. To overcome this we develop new query rewriting techniques. In particular we show that best answers to any union of conjunctive queries can be computed by issuing a new first-order query directly on the incomplete database. Query rewritings are standard in the context of, e.g., consistent query answering, OBDA, query answering using views etc., i.e., all contexts where only partial information is

available about the data to be queried [Calvanese *et al.*, 2000; Calvanese *et al.*, 2007; Cali *et al.*, 2013; Cali *et al.*, 2003b]. First-order rewritings are particularly useful, as they allow to use the power of standard database query engines. In fact when they exist, the rewritten queries can be implemented in any relational query engine by expressing them in SQL, with no need to implement ad-hoc algorithms.

2 Preliminaries

We represent missing information in relational databases in the standard way using nulls [Abiteboul *et al.*, 1995; Imieliński and Lipski, 1984; van der Meyden, 1998]. Databases are populated by *constants* and *nulls*, coming respectively from two countably infinite sets Const and Null . We denote nulls by \perp , sometimes with sub- or superscript. We also allow them to repeat, thus adopting the model of *marked* nulls, as customary in the context of applications such as OBDA or data integration and exchange. A relational schema, or vocabulary σ , is a set of relation names with associated arities. A database D over σ associates to each relation name of arity k in σ , a k -ary relation which is a finite subset of $(\text{Const} \cup \text{Null})^k$. Sets of constants and nulls occurring in D are denoted by $\text{Const}(D)$ and $\text{Null}(D)$. The *active domain* of D is $\text{adom}(D) = \text{Const}(D) \cup \text{Null}(D)$. A *complete* database has no nulls.

A *valuation* $v : \text{Null}(D) \rightarrow \text{Const}$ on a database D is a map that assigns constant values to nulls occurring in D . By $v(D)$ [resp. $v(\bar{a})$] we denote the result of replacing each null \perp by $v(\perp)$ in D [resp. in the tuple \bar{a}]. The semantics $\llbracket D \rrbracket$ of an incomplete database D is the set $\{v(D) \mid v \text{ is a valuation on } D\}$ of all complete databases it can represent. $\mathbf{V}(D)$ denotes the set of all valuations defined on D .

An m -ary query of active domain $C \subseteq \text{Const}$ is a map that associates with a database D a subset of $(\text{adom}(D) \cup C)^m$. The active domain of a query will be denoted as $\text{adom}(Q)$. To answer a query Q over an incomplete database D we follow [Lipski, 1984; Libkin, 2018b] and adopt a slight generalisation of the usual intersection based certain answers notion. We define the set of *certain answers* to Q over D as $\square(Q, D) = \{\bar{a} \text{ over } \text{adom}(D) \cup \text{adom}(Q) \mid \forall v : v(\bar{a}) \in Q(v(D))\}$. The only difference with the usual notion is that we allow answers to contain nulls.

Following [Libkin, 2018b], given a query Q , a database D , and a tuple \bar{a} over $\text{adom}(D) \cup \text{adom}(Q)$, we let the *support* of \bar{a} be the set of all valuations that witness it:

$$\text{Supp}(Q, D, \bar{a}) = \{v \in \mathbf{V}(D) \mid v(\bar{a}) \in Q(v(D))\}.$$

Supports thus measure how close a tuple is to certainty. We consider one answer to be *better* than another if it has more support. That is, given a database D , a k -ary query Q , and k -tuples \bar{a}, \bar{b} over $\text{adom}(D) \cup \text{adom}(Q)$, we let

$$\bar{a} \triangleleft_{Q,D} \bar{b} \Leftrightarrow \text{Supp}(Q, D, \bar{a}) \subset \text{Supp}(Q, D, \bar{b}).$$

The set of *best answers* to Q over D is defined as the set of answers for which there is no better one : $\mathbf{Best}(Q, D) = \{\bar{a} \mid \neg \exists \bar{b} : \bar{a} \triangleleft_{Q,D} \bar{b}\}$.

We focus on *first-order* (FO) queries of vocabulary σ written here in the logical notation using Boolean connectives

\wedge, \vee, \neg and quantifiers \exists, \forall . We write $\varphi(\bar{x})$ for an FO-formula φ with free variables \bar{x} . With slight abuse of notation, \bar{x} will denote both a tuple of variables and the set of variables occurring in it. The set of constants used by φ is as usual denoted by $\text{adom}(\varphi)$, and gives the active domain of the associated query. We interpret FO-formulas under active domain semantics, i.e. we consider D as a relational structure with universe $\text{adom}(D) \cup \text{adom}(\varphi)$. Thus an FO formula $\varphi(\bar{x})$ represents a query (of active domain $\text{adom}(\varphi)$) mapping each database D into the set of tuples $\{\bar{t} \text{ over } \text{adom}(D) \cup \text{adom}(\varphi) \mid D \models \varphi(\bar{t})\}$.

To evaluate FO-formulas with free variables we use assignments ν from variables to constants in the active domain. Note that with a little abuse of notation we write $D \models \varphi(\bar{t})$ for $D \models_{\nu} \varphi(\bar{x})$ under the assignment ν sending \bar{x} to \bar{t} .

Here it is important to note that the query associated to φ is a mapping defined on all databases D , possibly with nulls. If D contains nulls, $D \models \varphi(\bar{t})$ is to be intended “naïvely”, i.e. nulls of D are treated as new constants in the domain of D , distinct from each other, and distinct from all the other constants in D and φ . For example the query $\varphi(x, y) = \exists z (R(x, z) \wedge R(z, y))$, on the database $D = \{R(1, \perp_1), R(\perp_1, \perp_2), R(\perp_3, 2)\}$ selects only the tuple $(1, \perp_2)$.

We consider the \exists, \wedge, \vee -fragment of FO known as *unions of conjunctive queries* and its \exists, \wedge -fragment known as *conjunctive queries*.

Example 2.1. Let $Q(x) = \exists y (R(y) \wedge S(y, x))$ and $D = \{R(\perp_1), R(1), S(\perp_2, \perp_2)\}$. We have $\text{Supp}(Q, D, \perp_2) = \{v \in \mathbf{V}(D) \mid v(\perp_2) = 1 \text{ or } v(\perp_1) = v(\perp_2)\}$, $\text{Supp}(Q, D, 1) = \{v \in \mathbf{V}(D) \mid v(\perp_2) = 1\}$ and $\text{Supp}(Q, D, \perp_1) = \{v \in \mathbf{V}(D) \mid v(\perp_1) = v(\perp_2)\}$. It follows that $\square(Q, D) = \emptyset$ and $\mathbf{Best}(Q, D) = \{(\perp_2)\}$.

In order to study the complexity of best answer computation we shall need two classes in the second level of the polynomial hierarchy. Both of these contain NP and coNP, and are contained in $\Sigma_2^P \cap \Pi_2^P$. The class DP consists of languages $L_1 \cap L_2$ where $L_1 \in \text{NP}$ and $L_2 \in \text{coNP}$. This class has appeared in database applications [Fagin *et al.*, 2005; Barceló *et al.*, 2014]. The class $\text{P}^{\text{NP}[\log n]}$ consists of problems that can be solved in polynomial time with a logarithmic number of calls to an NP oracle [Eiter and Gottlob, 1997]. Equivalently, it can be described as the class of problems solved in P with an NP oracle where calls to the oracle are done in parallel, i.e., independent of each other. This class has appeared in the context of AI, modal logic, OBDA [Gottlob, 1995; Eiter and Gottlob, 1997; Calvanese *et al.*, 2006; Bienvenu and Bourgaux, 2016], data exchange [Arenas *et al.*, 2013].

3 Complexity of Best and Certain Answers

A natural way to cast computing best answers as a decision problem would be to proceed along the lines of certain answers [Imieliński and Lipski, 1984] and ask whether a given tuple is a best answer :

PROBLEM:	BESTANSWER
INPUT:	A query Q , a database D , a tuple \bar{a}
QUESTION:	Is $\bar{a} \in \text{Best}(Q, D)$?

Instead, [Libkin, 2018b] considers the following variant of the problem, asking whether the set of best answers belongs to a specified family of sets :

PROBLEM:	BESTANSWER [∈]
INPUT:	A query Q , a database D , a family \mathcal{X} of sets of tuples
QUESTION:	Is $\text{Best}(Q, D) \in \mathcal{X}$?

This suggests yet another alternative formulation of the problem, asking if a given set is the best answer :

PROBLEM:	BESTANSWER ⁼
INPUT:	A query Q , a database D , a set X of tuples,
QUESTION:	Is $\text{Best}(Q, D) = X$?

BESTANSWER[∈] was shown in [Libkin, 2018b] to be $\text{P}^{\text{NP}[\log n]}$ -complete in data complexity. We show that the other alternatives are computationally equivalent. Interestingly, the situation is very different with certain answers, as we show next.

Theorem 3.1. *For FO queries the problems BESTANSWER and BESTANSWER⁼ are $\text{P}^{\text{NP}[\log n]}$ -complete in data complexity.*

Proof (sketch). The upper bound for BESTANSWER⁼ immediately follows from the upper bound for BESTANSWER[∈] (take the family \mathcal{X} to be a singleton $\{X\}$). As for BESTANSWER we only need a slight modification of the upper bound proof in [Libkin, 2018b]. To check whether $\bar{a} \in \text{Best}(Q, D)$ we proceed as follows. Since the query is fixed, and has therefore fixed arity k , in polynomial time we can enumerate all the k -tuples of $\text{adom}(D)$. Then, using parallel calls to the NP oracle, we can check for each such tuple \bar{b} whether $\text{Supp}(Q, D, \bar{a}) \subseteq \text{Supp}(Q, D, \bar{b})$ and whether $\text{Supp}(Q, D, \bar{b}) \subseteq \text{Supp}(Q, D, \bar{a})$. With this information, in polynomial time we know whether $\bar{a} \triangleleft_{Q, D} \bar{b}$ for some \bar{b} .

We prove the two remaining lower bounds, reducing from the same $\text{P}^{\text{NP}[\log n]}$ -complete problem [Wagner, 1990]: given an undirected graph G , is its chromatic number $\chi(G)$ odd? With each undirected graph $G = \langle N, E \rangle$ with nodes N and edges E , we associate a database D_G over binary relations L, E and unary relations C, O as follows. We use a null \perp_n in D_G for each node n of G . For each edge $\{n, n'\}$ of G , we have pairs $(\perp_n, \perp_{n'})$ and $(\perp_{n'}, \perp_n)$ in the relation E of D_G . In relation C we have m constants $\{c_1, \dots, c_m\}$ (intuitively representing possible colors), where m is the number of nodes of G . Relation O of D_G is defined as $\{c_i \mid i \text{ is odd}\}$ and L is a linear ordering on them, i.e., $(c_i, c_j) \in L$ iff $i \leq j$, for $i, j \leq m$.

Remark that any valuation v of D_G that maps each null into a constant of C represents an assignment of colours in

$\{c_1, \dots, c_m\}$ to nodes of G . Then we define a query

$$\varphi(x) = \begin{aligned} & C(x) \\ & \wedge \forall y, z (E(y, z) \rightarrow L(y, x)) \\ & \wedge \forall y (L(y, x) \rightarrow \exists z E(y, z)) \\ & \wedge \neg \exists y E(y, y). \end{aligned}$$

For any valuation v , $\varphi(c)$ holds in $v(D_G)$ iff 1) $c = c_j$ for some $j = 1..m$ (ensured by the first conjunct). 2) For such a c_j , the valuation v maps each null into $\{c_1, \dots, c_j\}$ (second conjunct), i.e. v represents an assignment of colours to nodes of G , using at most the first j colours. 3) Each color $\{c_1, \dots, c_j\}$ is used by v , i.e. v represents an assignment of colours to nodes of G , using precisely the first j colours (third conjunct). 4) There are no loops in E (fourth conjunct).

Thus, for a valuation v , the formula $\varphi(c_j)$ is true in $v(D_G)$ iff v represents a colouring of G using precisely the first j colours $\{c_1, \dots, c_j\}$ (which in the sequel we refer to as an *exact j -colouring* of G).

Next we define :

$$Q(x) = C(x) \wedge (\varphi(x) \vee \exists y (O(y) \wedge L(x, y) \wedge \varphi(y)))$$

For a valuation v , we have that $Q(c_i)$ holds in $v(D_G)$ iff either v represents an exact i -coloring of G ; or v represents an exact j -coloring of G with j odd, and $i \leq j$. In other words valuations representing exact j -colorings, with j even, support only the maximal color c_j ; while valuations representing exact j -colorings, with j odd, support all colors $\{c_1 \dots c_j\}$.

With this in place we can conclude the reduction for the BESTANSWER problem :

Claim. $c_1 \in \text{Best}(Q, D_G)$ iff the chromatic number of G is odd.

An adaptation of this encoding can be used to reduce the odd chromatic number problem to BESTANSWER⁼ as well. \square

Now that we showed that all three formulations of best answers actually collapse computationally, another natural question arises. Does a similar result hold for certain answers ? We obtain the decision problems CERTAINANSWER, CERTAINANSWER⁼ and CERTAINANSWER[∈] by replacing everywhere in the statements of the above decision problems BESTANSWER by CERTAINANSWER and $\text{Best}(Q, D)$ by $\square(Q, D)$. It is well known that data complexity of CERTAINANSWER is CONP -complete for FO-queries [Abiteboul et al., 1991]. We complete the picture as follows.

Theorem 3.2. *For FO queries, CERTAINANSWER⁼ is DP-complete and CERTAINANSWER[∈] is $\text{P}^{\text{NP}[\log n]}$ -complete in data complexity.*

Proof. To prove membership of CERTAINANSWER⁼ in DP, notice that for a query Q , this problem is the intersection of two languages $L_1 \cap L_2$ where $L_1 = \{(D, X) \mid X \subseteq \square(Q, D)\}$ and $L_2 = \{(D, X) \mid \bar{X} \subseteq \square(Q, \bar{D})\}$. L_1 is known to be in CONP : we guess a tuple $\bar{a} \in X$ and a valuation $v \in V(D)$ with $v(\bar{a}) \notin Q(v(D))$. Similarly, L_2 is in NP : we guess a tuple $\bar{b} \in \bar{X}$ and a valuation $v' \in V(D)$ with $v'(\bar{b}) \in Q(v'(D))$.

Type of problem	$\bar{a} \in \text{Answer}$	$X = \text{Answer}$	$\text{Answer} \in \mathcal{X}$
Certain Answer	coNP-complete [Abiteboul <i>et al.</i> , 1991]	DP complete	$\text{P}^{\text{NP}[\log n]}$ -complete
Best Answer	$\text{P}^{\text{NP}[\log n]}$ -complete	$\text{P}^{\text{NP}[\log n]}$ -complete	$\text{P}^{\text{NP}[\log n]}$ -complete [Libkin, 2018b]

Figure 1: Summary of data complexity results for FO queries

To prove membership of $\text{CERTAINANSWER}^\infty$ in $\text{P}^{\text{NP}[\log n]}$, suppose the query Q is k -ary, and we are given a family of sets of k -ary tuples $\mathcal{X} = \{X_1, \dots, X_n\}$ and a database D . For each $X_i \in \mathcal{X}$, we use the NP oracle to decide in parallel whether $X_i = \sqcap(Q, D)$ (for each X_i the two calls to the oracle do not depend on each other and they can also be done in parallel).

For DP-hardness, we reduce from the problem of checking whether $\chi(G)$, the chromatic number of an undirected graph G , equals 4 [Rothe, 2003] and for $\text{P}^{\text{NP}[\log n]}$ -hardness, we reduce from the related problem of checking whether $\chi(G)$ is odd. With such a graph G , we associate the same database D_G as in the proof of Theorem 3.1. Using the exact coloring formula φ in the proof of Theorem 3.1, we define a query

$$Q(x) := C(x) \wedge \forall y (\varphi(y) \rightarrow L(x, y))$$

We claim that $\sqcap(Q, D) = \{c_1, \dots, c_n\}$ iff $\chi(G) = n$, which entails $\sqcap(Q, D) = \{c_1, \dots, c_4\}$ iff $\chi(G) = 4$ and $\sqcap(Q, D) \in \{\{c_1, \dots, c_j\} \mid j \text{ is odd and } 1 \leq j \leq |G|\}$ iff $\chi(G)$ is odd. Recall that $v(D_G) \models \varphi(c_i)$ iff c_i is a color in $\{c_1, \dots, c_{|G|}\}$ and v represents an exact i -coloring of the graph. Now $v(D_G) \models Q(c_j)$ iff c_j is a color and there is no $i < j$ such that v represents an exact i -coloring of the graph, which holds exactly whenever $c_j \in \{c_1, \dots, c_{\chi(G)}\}$. \square

4 First-Order Rewritings for Best Answers

Considering arbitrary FO-queries brought us an intrinsic intractability result for all variants of best answers. This motivates restricting to unions of conjunctive queries, for which a polynomial time evaluation algorithm (in data complexity) already exists [Libkin, 2018b]. The resolution based procedure is however in sharp contrast with naive evaluation, which allows to compute certain answers to unions of conjunctive queries via usual model checking. We thus initiate a descriptive complexity analysis of the best answers problem, showing that for unions of conjunctive queries, it can essentially be reduced - modulo a preprocessing of the query - to (naive) evaluation of an FO-formula.

Given a union of conjunctive queries Q , our starting point towards an FO-rewriting for best answers is finding an FO-formula $Q_\subseteq(\bar{x}, \bar{y})$ encoding the inclusion of supports, i.e. selecting tuples \bar{s}, \bar{t} over $\text{adom}(D) \cup \text{adom}(Q)$ iff $\text{Supp}(Q, D, \bar{s}) \subseteq \text{Supp}(Q, D, \bar{t})$. From Q_\subseteq one can easily define an FO-formula selecting precisely all best answers to Q on D :

$$\text{best}_Q(\bar{x}) := \forall \bar{y} (Q_\subseteq(\bar{x}, \bar{y}) \rightarrow Q_\subseteq(\bar{y}, \bar{x})) \quad (1)$$

We start by putting each conjunctive query in a normal form which eliminates repetition of variables, by introducing new equality atoms.

Definition 4.1 (NRV normal form). *A conjunctive query Q is in non repeating variable normal form (NRV normal form) if it is of the form $Q(\bar{x}) = \exists \bar{y}, \bar{z} (q(\bar{y}, \bar{z}) \wedge e(\bar{y}, \bar{z}) \wedge \bar{x} = \bar{z})$ where variables in $\bar{x}\bar{y}\bar{z}$ are pairwise distinct, \bar{x} and \bar{z} have the same length, and:*

- $q(\bar{y}, \bar{z})$ is a conjunction of relational atoms, where each free variable in \bar{y}, \bar{z} has at most one occurrence in q ,
- $e(\bar{y}, \bar{z})$ is a conjunction of equality atoms, possibly using constants.

We say that $q(\bar{y}, \bar{z})$ is the relational subquery of Q , and $e(\bar{y}, \bar{z}) \wedge \bar{x} = \bar{z}$ is the equality subquery of Q .

Example 4.2. *The query $Q(x)$ from Example 2.1 is equivalent to $\exists y_1 y_2 z (R(y_1) \wedge S(y_2, z) \wedge y_1 = y_2 \wedge z = x)$, which is in NRV normal form.*

Clearly every conjunctive query Q is equivalent to a query in NRV normal form; moreover Q can be easily rewritten in NRV normal form (in linear time in the size of the query). Thus in what follows we assume w.l.o.g. that conjunctive queries are given in NRV normal form. Intuitively the NRV normal form allows us to separate the two ingredients of a conjunctive query: the existence of facts in some relations of the database on the one side, and a set of equality conditions on data values occurring in these facts, on the other side. The existence of facts does not depend on the valuation of nulls, and thus can be directly tested on the incomplete database. Instead equality atoms in an NRV normal form imply conditions that valuations need to satisfy in order for the query to hold. We can thus first concentrate on the support of equality subqueries. This will be encoded in FO and then integrated in the rewriting of the whole conjunctive query.

We introduce a notion of equivalence of database elements w.r.t. to a set of equalities. Intuitively equivalent elements of a tuple \bar{t} are the ones which should be collapsed into a single value in order for a valuation of \bar{t} to satisfy all the given equalities.

Definition 4.3. *Given a database D , a conjunction of equality atoms $\gamma(\bar{y})$ and an assignment $\nu : \bar{y} \cup \text{adom}(\gamma) \rightarrow \text{adom}(D) \cup \text{adom}(\gamma)$ preserving constants, we say that $u, u' \in \text{adom}(D) \cup \text{adom}(\gamma)$ are equivalent w.r.t. γ and ν and write $u \equiv_\gamma^\nu u'$, if either $u = u'$ or (u, u') belongs to the reflexive symmetric transitive closure of $\{(\nu(x), \nu(w)) \mid x = w \in \gamma\}$.*

The relation \equiv_γ^ν is clearly an equivalence relation over $\text{adom}(D) \cup \text{adom}(\gamma)$, where each element outside the range of ν forms a singleton equivalence class.

Example 4.4. *Let γ be $x_1 = x_2 \wedge x_2 = x_3 \wedge x_4 = x_5 \wedge x_6 = 1$. Let ν assign \perp_i to x_i for $i \leq 5$, and \perp_5 to x_6 . The equivalence classes of \equiv_γ^ν are $\{\perp_i \mid i \leq 3\}$ and $\{1, \perp_4, \perp_5\}$, plus one singleton for each other domain element.*

In what follows we denote by \sim_γ the reflexive symmetric transitive closure of $\{(x, w) \mid x = w \in \gamma\}$. Note that this is an equivalence relation among variables and constants of γ .

We will be using the following formula to provide a syntactic characterisation of \equiv'_γ , where m is the number of equivalence classes of \sim_γ :¹

$$\text{equiv}_\gamma(\bar{y}, z, z') := z = z' \vee \bigvee_{\substack{u_1, v_1 \dots u_m, v_m \in \bar{y} \cup \text{adom}(\gamma) \\ u_i \sim_\gamma v_i \text{ for all } 1 \leq i \leq m}} (z = u_1 \wedge z' = v_m \wedge \bigwedge_{1 \leq i < m} v_i = u_{i+1})$$

Proposition 4.5. *Given an incomplete database D , a conjunction of equality atoms $\gamma(\bar{y})$ and an assignment $\nu(\bar{y}) = \bar{t}$ over $\text{adom}(D) \cup \text{adom}(\gamma)$, given s, s' in $\bar{t} \cup \text{adom}(\gamma)$, we have that $D \models \text{equiv}_\gamma(\bar{t}, s, s')$ if and only if $s \equiv'_\gamma s'$.*

Intuitively this holds because each disjunct of $\text{equiv}_\gamma(\bar{t}, s, s')$ corresponds to a possible derivation of (s, s') in the reflexive symmetric transitive closure of $\{(\nu(x), \nu(w)) \mid x = w \in \gamma\}$, and one can prove that there is a bound only depending on γ on the number of steps of this derivation.

Example 4.6. *Let $\gamma := y_1 = y_2 \wedge z = x$ be the equality subquery of the query $Q(x)$ in Example 4.2. Up to logical equivalence, $\text{equiv}_\gamma(y_1, y_2, z, x, w, w')$ contains precisely the disjuncts $w = w'$, $w = y_1 \wedge w' = y_2$, $w = z \wedge w' = x$, $w = y_1 \wedge w' = x \wedge y_2 = z$, plus all disjuncts obtained from them by applying one or more of the following transformations: switch w and w' , switch y_1 and y_2 , switch x and z . Let D be the database from Example 2.1, then we have for instance $D \models \text{equiv}_\gamma(1, \perp_2, \perp_2, 1, a, a')$ and $D \models \text{equiv}_\gamma(1, \perp_2, \perp_2, \perp_2, a, a')$ for all $a, a' \in \{1, \perp_2\}$. Similarly $D \models \text{equiv}_\gamma(\perp_1, \perp_2, \perp_2, 1, a, a')$ for all $a, a' \in \{1, \perp_1, \perp_2\}$.*

As a consequence of Proposition 4.5, for fixed γ and \bar{t} , the relation $\{(s, s') \mid D \models \text{equiv}_\gamma(\bar{t}, s, s')\}$ is an equivalence relation over $\text{adom}(D) \cup \text{adom}(\gamma)$ where each element of $\text{adom}(D)$ neither in \bar{t} nor in $\text{adom}(\gamma)$ forms a singleton equivalence class.

The formula equiv_γ is a key ingredient towards a rewriting of a conjunctive query; in fact, as formalized in the following lemma, it selects precisely the pairs of elements of a tuple that a valuation needs to collapse to satisfy a set of equalities.

Lemma 4.7. *Let $\gamma(\bar{y})$ be a conjunction of equality atoms, D a database and $\nu(\bar{y}) = \bar{t}$ an assignment over $\text{adom}(D) \cup \text{adom}(\gamma)$. Assume v is a valuation of nulls. Then $v(D) \models \gamma(v(\bar{t}))$ if and only if $v(s) = v(s')$ for all s, s' such that $D \models \text{equiv}_\gamma(\bar{t}, s, s')$.*

Example 4.8. *Let γ and ν be as in Example 4.4, then Lemma 4.7 implies that a valuation $v(D) \models \gamma(v(\bar{t}))$ iff $v(\perp_i) = v(\perp_j)$ for all $i, j = 1..3$, and $v(\perp_i) = 1$ for all $i = 4, 5$.*

Formulas we write in the remainder of this section are over signature $\sigma \cup \text{Null}$, where σ is the database schema. In any

¹Queries we write in the sequel can be domain dependent. So it is important to recall that we always use active domain semantics.

incomplete database D over $\sigma \cup \text{Null}$, Null is always interpreted by the set of nulls occurring in D (in accordance with the semantics of the SQL construct `IS NULL`). I.e. we allow rewritings to test whether a database element is null or not.

For $\gamma(\bar{y})$ a conjunction of equality atoms, using equiv_γ we define a new formula $\text{comp}_\gamma(\bar{y})$ stating the existence of a valuation that collapses all equivalent elements of a tuple:

$$\text{comp}_\gamma(\bar{y}) := \forall z z' (\text{equiv}_\gamma(\bar{y}, z, z') \wedge \neg \text{Null}(z) \wedge \neg \text{Null}(z') \rightarrow z = z')$$

Notice that if $D \models \text{comp}_\gamma(\bar{t})$ then for each $s \in \text{adom}(D) \cup \text{adom}(\gamma)$ there exists at most one constant c such that $D \models \text{equiv}_\gamma(\bar{t}, s, c)$. In fact if for constants c_1 and c_2 , $D \models \text{equiv}_\gamma(\bar{t}, s, c_1)$ and $D \models \text{equiv}_\gamma(\bar{t}, s, c_2)$, by transitivity $D \models \text{equiv}_\gamma(\bar{t}, c_1, c_2)$, implying $c_1 = c_2$.

Example 4.9. *Let D and γ be as in Example 4.6. Consider $\text{comp}_\gamma(y_1, y_2, z, x)$. Given the tuples selected by equiv_γ in Example 4.6, we can conclude that $D \models \text{comp}_\gamma(1, \perp_2, \perp_2, 1)$.*

Proposition 4.10. *Let $\gamma(\bar{y})$ be a conjunction of equality atoms, D a database and $\nu(\bar{y}) = \bar{t}$ an assignment over $\text{adom}(D) \cup \text{adom}(\gamma)$, then $D \models \text{comp}_\gamma(\bar{t})$ if and only if there exists a valuation v of nulls such that $v(D) \models \gamma(v(\bar{t}))$.*

We are now ready to define a formula capturing the inclusion of supports between two conjunctions of equality atoms, which will be a crucial ingredient in our rewriting. Let $\gamma(\bar{x})$ and $\gamma'(\bar{y})$ be conjunctions of equality atoms with $\text{adom}(\gamma) = \text{adom}(\gamma')$. We define:

$$\text{imply}_{\gamma, \gamma'}(\bar{x}, \bar{y}) := \forall z z' (\text{equiv}_{\gamma'}(\bar{y}, z, z') \rightarrow \text{equiv}_\gamma(\bar{x}, z, z'))$$

Example 4.11. *Let γ and D be as in Example 4.6. Let $\gamma' := y'_1 = y'_2 \wedge z' = x'$, then it follows from Example 4.6 that $D \models \text{imply}_{\gamma, \gamma'}(\perp_1 \perp_2 \perp_2 1, 1 \perp_2 \perp_2 \perp_2)$ and $D \models \text{imply}_{\gamma, \gamma'}(1 \perp_2 \perp_2 1, 1 \perp_2 \perp_2 \perp_2)$.*

Using Proposition 4.10 and Lemma 4.7 we obtain:

Proposition 4.12. *Let $\gamma(\bar{x}), \gamma'(\bar{y})$ be conjunctions of equality atoms with $\text{adom}(\gamma) = \text{adom}(\gamma')$, D a database and $\nu(\bar{y}) = \bar{t}$, $\nu'(\bar{y}) = \bar{t}'$ assignments over $\text{adom}(D) \cup \text{adom}(\gamma)$. Then $D \models \text{imply}_{\gamma, \gamma'}(\bar{t}, \bar{t}') \vee \neg \text{comp}_\gamma(\bar{t})$ iff for all valuations v , $v(D) \models \gamma(v(\bar{t}))$ implies $v(D) \models \gamma'(v(\bar{t}'))$.*

By combining Propositions 4.12 and 4.10 we also get:

Corollary 4.13. *Let $\gamma(\bar{y}), \gamma'(\bar{y})$ be conjunctions of equality atoms with $\text{adom}(\gamma) = \text{adom}(\gamma')$, D a database and $\nu(\bar{y}) = \bar{t}$, $\nu'(\bar{y}) = \bar{t}'$ assignments over $\text{adom}(D) \cup \text{adom}(\gamma)$. If $D \models \text{comp}_\gamma(\bar{t}) \wedge \text{imply}_{\gamma, \gamma'}(\bar{t}, \bar{t}')$, then $D \models \text{comp}_{\gamma'}(\bar{t}')$.*

We now go back to an arbitrary union of conjunctive queries of vocabulary σ in NRV-normal form:

$$Q(\bar{x}) := \bigvee_{1 \leq i \leq n} Q_i(\bar{x})$$

where each Q_i is in NRV-normal form with relational subquery $q_i(\bar{y}_i, \bar{z}_i)$ and equality subquery $eq_i(\bar{x}, \bar{y}_i, \bar{z}_i)$.

Lemma 4.14. *Let D be a database, v a valuation of D and $Q(\bar{x})$ a union of conjunctive queries in NRV-normal form, then $v \in \text{Supp}(Q, D, \bar{r})$ if and only there exists i, \bar{s} and \bar{t} such that $D \models q_i(\bar{s}, \bar{t}) \wedge \text{comp}_{e_{q_i}}(\bar{r}\bar{s}\bar{t})$ and $v(D) \models e_{q_i}(v(\bar{r}\bar{s}\bar{t}))$.*

We are now ready to define the FO-formula encoding the inclusion of supports.

$$Q_{\subseteq}(\bar{x}, \bar{x}') := \bigwedge_{1 \leq i \leq n} (\forall \bar{y}\bar{z}((q_i(\bar{y}, \bar{z}) \wedge \text{comp}_{e_{q_i}}(\bar{x}, \bar{y}, \bar{z})) \rightarrow \bigvee_{1 \leq j \leq n} \exists \bar{y}'\bar{z}'(q_j(\bar{y}', \bar{z}') \wedge \text{imply}_{e_{q_i}, e_{q_j}}(\bar{x}\bar{y}\bar{z}, \bar{x}'\bar{y}'\bar{z}'))))$$

Combining Lemmas 4.7, 4.14, Propositions 4.10, 4.12 and Corollary 4.13 we get :

Proposition 4.15. $D \models Q_{\subseteq}(\bar{s}, \bar{t})$ iff $\text{Supp}(Q, D, \bar{s}) \subseteq \text{Supp}(Q, D, \bar{t})$.

Recall that from Q_{\subseteq} one can easily define a first order rewriting $\text{best}_Q(\bar{x})$ for best answers as in (1).

Theorem 4.16. *Given Q a union of conjunctive queries over schema σ and an incomplete database D , $\bar{t} \in \text{Best}(Q, D)$ iff $D \models \text{best}_Q(\bar{t})$.*

Example 4.17. *For Q, D, γ, γ' as in Example 2.1 and 4.11 :*

$$Q_{\subseteq}(x, x') := \forall y_1 y_2 z ((R(y_1) \wedge S(y_2, z) \wedge \text{comp}_{\gamma}(y_1, y_2, z, x)) \rightarrow \exists y'_1 y'_2 z' (R(y'_1) \wedge S(y'_2, z') \wedge \text{imply}_{\gamma, \gamma'}(y_1 y_2 z x, y'_1 y'_2 z' x'))).$$

This allows to derive for instance $\text{Supp}(Q, D, 1) \subseteq \text{Supp}(Q, D, \perp_2)$ (as observed in Example 2.1). In fact the subquery $R(y_1) \wedge S(y_2, z) \wedge \text{comp}_{\gamma}(y_1, y_2, z, x)$ with free variables y_1, y_2, z, x selects on D tuples $(1, \perp_2, \perp_2, 1), (\perp_1, \perp_2, \perp_2, 1)$, and no other tuple with last element 1. Moreover as shown in Example 4.11

$$\begin{aligned} D &\models \text{imply}_{\gamma, \gamma'}(\perp_1 \perp_2 \perp_2 1, 1 \perp_2 \perp_2 \perp_2) \\ D &\models \text{imply}_{\gamma, \gamma'}(1 \perp_2 \perp_2 1, 1 \perp_2 \perp_2 \perp_2) \end{aligned}$$

Thus $D \models Q_{\subseteq}(1, \perp_2)$. Similarly one can show $D \models Q_{\subseteq}(\perp_1, \perp_2)$ and therefore $D \models \text{best}_Q(\perp_2)$.

As a corollary of Theorem 4.16, for a union of conjunctive queries Q one can compute $\text{Best}(Q, D)$ by first computing the formula $\text{best}_Q(\bar{x})$ from Q , then evaluating best_Q on D . Since data complexity of FO query evaluation is DLOGSPACE (and in particular AC⁰), this gives the following corollary :

Corollary 4.18. *For each fixed union of conjunctive queries Q , the data complexity of BESTANSWER is DLOGSPACE.*

Note that it was known from [Libkin, 2018b] that the data complexity of computing best answers for unions of conjunctive queries is polynomial time. In terms of combined complexity (i.e. when either Q, D and \bar{a} are in the input), the rewriting approach (i.e. the procedure of computing best_Q from Q and then evaluating best_Q on D), can be easily shown to be in PSPACE. In fact it is well known that a first order query φ can be evaluated on a database D in space at most $qr(\varphi) \log |D| + \log |\varphi|$, where $qr(\varphi)$ is the quantifier rank of φ . Note that although best_Q has size exponential

in Q , the quantifier rank of best_Q is linear in the size of Q . Thus whether $\bar{a} \in \text{best}(Q, D)$ can be checked using space $O(|Q|, |D|)$. Moreover one can easily check that best_Q can be computed from Q in space polynomial in the size of $|Q|$. Since space bounded computations can be composed without storing the intermediate output, computing best_Q from Q and then evaluating best_Q on D can be done overall in PSPACE in the size of $|Q|$ and $|D|$. The rewriting approach thus implies a PSPACE upper bound for the combined complexity of BESTANSWER for unions of conjunctive queries. However we can show that the problem actually stands in the third level of the polynomial hierarchy.

Theorem 4.19. *For unions of conjunctive queries, combined complexity of BESTANSWER is Π_3^P -complete. Hardness already holds for conjunctive queries.*

Therefore under standard complexity theoretic assumptions, our rewriting approach is not optimal in terms of combined complexity, as it is often the case with generic approaches. However it has the advantage of exploiting standard FO query evaluation, which despite the PSPACE combined complexity, is highly optimised in database systems and works well in practice.

5 Future Work

Constraints (e.g., keys and functional dependencies) are known to raise the complexity of finding certain answers [Calì *et al.*, 2003a]. They appear in another model of incompleteness - conditional tables [Imieliński and Lipski, 1984] - that in general leads to higher complexity of query evaluation [Abiteboul *et al.*, 1995] but are nonetheless useful in several applications [Arenas *et al.*, 2013]. We would like to explore how our rewriting techniques interact with integrity constraints.

In another direction, while we focused on FO-rewritings, we could consider more expressive rewriting languages such as Datalog or fixed point logics, as it is common in the context of OBDA, query answering using views, or consistent query answering [Bienvenu and Ortiz, 2015; Francis *et al.*, 2015; Bertossi, 2011]. These more expressive logics are likely to be able to express rewritings of larger classes of queries than unions of conjunctive queries.

We would also like to investigate how our techniques can be extended to different semantics of incompleteness. We used here the closed-world semantics [Abiteboul *et al.*, 1995; Imieliński and Lipski, 1984; van der Meyden, 1998], in which data values are the only missing information, but there are other possible semantics, e.g. needed in order to cope with data inconsistencies [Calì *et al.*, 2003a], where query rewritings could still be found.

Finally, we would like to investigate how techniques developed in this paper could be extended to study rewritings of certain answers.

Acknowledgements

We thank the anonymous referees and Leonid Libkin for their useful feedback. This work was supported by contract ANR-18-CE40-0031 QUID.

References

- [Abiteboul *et al.*, 1991] Serge Abiteboul, Paris C. Kanelakis, and Gösta Grahne. On the representation and querying of sets of possible worlds. *TCS*, 78(1):158–187, 1991.
- [Abiteboul *et al.*, 1995] Serge Abiteboul, Richard Hull, and Victor Vianu. *Foundations of Databases*. Addison-Wesley, Boston, MA, USA, 1995.
- [Amendola and Libkin, 2018] Giovanni Amendola and Leonid Libkin. Explainable certain answers. In *IJCAI*, pages 1683–1690, Stockholm, Sweden, August 2018.
- [Arenas *et al.*, 2013] Marcelo Arenas, Jorge Pérez, and Juan L. Reutter. Data exchange beyond complete data. *J. ACM*, 60(4):28:1–28:59, 2013.
- [Arenas *et al.*, 2014] Marcelo Arenas, Pablo Barcelo, Leonid Libkin, and Filip Murlak. *Foundations of Data Exchange*. Cambridge University Press, 2014.
- [Barceló *et al.*, 2014] Pablo Barceló, Leonid Libkin, and Miguel Romero. Efficient approximations of conjunctive queries. *SIAM J. Comput.*, 43(3):1085–1130, 2014.
- [Bertossi, 2011] Leopoldo E. Bertossi. *Database Repairing and Consistent Query Answering*. Synthesis Lectures on Data Management. Morgan & Claypool Publishers, 2011.
- [Bienvenu and Bourgaux, 2016] Meghyn Bienvenu and Camille Bourgaux. Inconsistency-tolerant querying of description logic knowledge bases. In *Tutorial Lectures of RW Summer School*, pages 156–202. LNCS, 2016.
- [Bienvenu and Ortiz, 2015] Meghyn Bienvenu and Magdalena Ortiz. Ontology-mediated query answering with data-tractable description logics. In *Reasoning Web. Web Logic Rules - 11th International Summer School*, pages 218–307, Berlin, Germany, August 2015. Springer.
- [Calì *et al.*, 2003a] Andrea Calì, Domenico Lembo, and Riccardo Rosati. On the decidability and complexity of query answering over inconsistent and incomplete databases. In *PODS*, pages 260–271, June 2003.
- [Calì *et al.*, 2003b] Andrea Calì, Domenico Lembo, and Riccardo Rosati. Query rewriting and answering under constraints in data integration systems. In *IJCAI*, pages 16–21, Acapulco, Mexico, August 2003. Morgan Kaufmann.
- [Calì *et al.*, 2013] Andrea Calì, Diego Calvanese, and Maurizio Lenzerini. Rewrite and conquer: Dealing with integrity constraints in data integrations. In *Seminal Contributions to Information Systems Engineering*, pages 353–359, 2013.
- [Calvanese *et al.*, 2000] Diego Calvanese, Giuseppe De Giacomo, Maurizio Lenzerini, and Moshe Y. Vardi. What is view-based query rewriting? In *KRDB*, pages 17–27, Berlin, Germany, August 2000.
- [Calvanese *et al.*, 2006] Diego Calvanese, Giuseppe De Giacomo, Domenico Lembo, Maurizio Lenzerini, and Riccardo Rosati. Epistemic first-order queries over description logic knowledge bases. In *DL*, Windermere, Lake District, UK, May-June 2006.
- [Calvanese *et al.*, 2007] Diego Calvanese, Giuseppe De Giacomo, Maurizio Lenzerini, and Moshe Y. Vardi. View-based query processing: On the relationship between rewriting, answering and losslessness. *TCS*, 371(3):169–182, March 2007.
- [Eiter and Gottlob, 1997] Thomas Eiter and Georg Gottlob. The complexity class theta_2^P : Recent results and applications in AI and modal logic. In *FCT*, pages 1–18, Kraków, Poland, September 1997. Springer.
- [Fagin *et al.*, 2005] Ronald Fagin, Phokion G. Kolaitis, and Lucian Popa. Data exchange: getting to the core. *ACM TODS*, 30(1):174–210, 2005.
- [Francis *et al.*, 2015] Nadime Francis, Luc Segoufin, and Cristina Sirangelo. Datalog rewritings of regular path queries using views. *Logical Methods in Computer Science*, 11(4), 2015.
- [Geerts *et al.*, 2013] Floris Geerts, Giansalvatore Mecca, Paolo Papotti, and Donatello Santoro. The LLUNATIC data-cleaning framework. *PVLDB*, 6(9):625–636, 2013.
- [Gheerbrant and Libkin, 2015] Amélie Gheerbrant and Leonid Libkin. Certain answers over incomplete xml documents: Extending tractability boundary. *ACM ToCS*, 57(4):892–926, July 2015.
- [Gottlob *et al.*, 2015] Georg Gottlob, Marco Manna, and Andreas Pieris. Polynomial rewritings for linear existential rules. In *IJCAI*, pages 2992–2998, Buenos Aires, Argentina, August 2015. AAAI Press.
- [Gottlob, 1995] Georg Gottlob. NP trees and Carnap’s modal logic. *Journal of the ACM*, 42(2):421–457, March 1995.
- [Imieliński and Lipski, 1984] Tomasz Imieliński and Witold Lipski. Incomplete information in relational databases. *Journal of the ACM*, 31(4):761–791, Oct 1984.
- [Lenzerini, 2002] Maurizio Lenzerini. Data integration: A theoretical perspective. In *ACM PODS*, pages 233–246, Madison, USA, June 2002.
- [Libkin, 2018a] Leonid Libkin. Certain answers as objects of knowledge. *Artificial Intelligence*, 232:195–207, 2018.
- [Libkin, 2018b] Leonid Libkin. Certain answers meet zero-one laws. In *ACM PODS*, pages 1–19, Houston, USA, June 2018.
- [Lipski, 1984] Witold Lipski. On relational algebra with marked nulls. In *PODS*, pages 201–203, Waterloo, Ontario, Canada, April 1984. ACM.
- [Rothe, 2003] Jörg Rothe. Exact complexity of exact-four-colorability. *Inf. Process. Lett.*, 87(1):7–12, 2003.
- [van der Meyden, 1998] Ron van der Meyden. Logical approaches to incomplete information: a survey. In *Logics for databases and information systems*, pages 307–356, Norwell, MA, USA, 1998. Kluwer Academic Publishers.
- [Wagner, 1990] Klaus W. Wagner. Bounded query classes. *SIAM J. Comput.*, 19(5):833–846, 1990.

A Proofs

Proof of Theorem 3.1

The upper bound for $\text{BESTANSWER}^=$ immediately follows from the upper bound for BESTANSWER^{\leq} (take the family \mathcal{X} to be a singleton $\{X\}$). As for BESTANSWER we only need a slight modification of the upper bound proof in [Libkin, 2018b]. To check whether $\bar{a} \in \text{Best}(Q, D)$ we proceed as follows. Since the query is fixed, and has therefore fixed arity k , in polynomial time we can enumerate all the k -tuples of $\text{adom}(D)$. Then, using parallel calls to the NP oracle, we can check for each such tuple \bar{b} whether $\text{Supp}(Q, D, \bar{a}) \subseteq \text{Supp}(Q, D, \bar{b})$ and whether $\text{Supp}(Q, D, \bar{b}) \subseteq \text{Supp}(Q, D, \bar{a})$. With this information, in polynomial time we know whether $\bar{a} \prec_{Q,D} \bar{b}$ for some \bar{b} .

We prove the two remaining lower bounds, reducing from the same $\text{P}^{\text{NP}[\log n]}$ -complete problem [Wagner, 1990]: given an undirected graph G , is its chromatic number $\chi(G)$ odd? With each undirected graph $G = \langle N, E \rangle$ with nodes N and edges E , we associate a database D_G over binary relations L, E and unary relations C, O as follows. We use a null \perp_n in D_G for each node n of G . For each edge $\{n, n'\}$ of G , we have pairs $(\perp_n, \perp_{n'})$ and $(\perp_{n'}, \perp_n)$ in the relation E of D_G . In relation C we have m constants $\{c_1, \dots, c_m\}$ (intuitively representing possible colors), where m is the number of nodes of G . Relation O of D_G is defined as $\{c_i \mid i \text{ is odd}\}$ and L is a linear ordering on them, i.e., $(c_i, c_j) \in L$ iff $i \leq j$, for $i, j \leq m$.

Remark that any valuation v of D_G that maps each null into a constant of C represents an assignment of colours in $\{c_1, \dots, c_m\}$ to nodes of G . Then we define a query

$$\begin{aligned} \varphi(x) = & C(x) \\ & \wedge \forall y, z (E(y, z) \rightarrow L(y, x)) \\ & \wedge \forall y (L(y, x) \rightarrow \exists z E(y, z)) \\ & \wedge \neg \exists y E(y, y). \end{aligned}$$

For any valuation v , $\varphi(c)$ holds in $v(D_G)$ iff 1) $c = c_j$ for some $j = 1..m$ (ensured by the first conjunct); 2) for such a c_j , the valuation v maps each null into $\{c_1, \dots, c_j\}$ (second conjunct), i.e. v represents an assignment of colours to nodes of G , using at most the first j colours. 3) Each color $\{c_1, \dots, c_j\}$ is used by v , i.e. v represents an assignment of colours to nodes of G , using precisely the first j colours (third conjunct). 4) There are no loops in E (fourth conjunct).

Thus, for a valuation v , the formula $\varphi(c_j)$ is true in $v(D_G)$ iff v represents a colouring of G using precisely the first j colours $\{c_1, \dots, c_j\}$ (which in the sequel we refer to as an *exact j -colouring* of G).

Next we define :

$$Q(x) = C(x) \wedge (\varphi(x) \vee \exists y (O(y) \wedge L(x, y) \wedge \varphi(y)))$$

For a valuation v , we have that $Q(c_i)$ holds in $v(D_G)$ iff either v represents an exact i -coloring of G ; or v represents an exact j -coloring of G with j odd, and $i \leq j$. In other words valuations representing exact j -colorings, with j even, support only the maximal color c_j ; while valuations representing exact j -colorings, with j odd, support all colors $\{c_1..c_j\}$.

With this in place we can conclude the reduction for the BESTANSWER problem :

Claim. $c_1 \in \text{Best}(Q, D_G)$ iff the chromatic number of G is odd.

Proof. . Let χ_G be the chromatic number of G . Then there exist no exact colorings of G which are prefixes of $\{c_1, \dots, c_{\chi_G}\}$, while $\{c_1, \dots, c_{\chi_G}\}$ is an exact coloring of G .

Assume first that χ_G is even. Then there exist no valuations representing the exact coloring $\{c_1\}$. Thus the support of c_1 is the set of valuation representing an exact coloring $\{c_1..c_j\}$ of G with j odd and $j > \chi_G$. This support is not maximal, In fact the support of c_{χ_G} is :

- the valuations representing the exact coloring $\{c_1..c_{\chi_G}\}$ (there exists at least one);
- the valuations representing an exact coloring $\{c_1..c_j\}$ of G with j odd and $j > \chi_G$.

This support strictly contains the support of c_1 ; in fact valuations in the first item cannot be also in the second.

Assume now that χ_G is odd. Then the support of c_1 is the set of valuations representing an exact coloring $\{c_1..c_j\}$ of G with j odd and $j \geq \chi_G$. We show that this set is maximal, i.e. no color c_k can have a support strictly containing it.

- if $k \leq \chi_G$ then the support of c_k is the set of valuations representing an exact coloring $\{c_1..c_j\}$ of G with j odd, and $j \geq \chi_G$. So same support as c_1 .
- if $k > \chi_G$, the support of c_k cannot contain the valuations representing $\{c_1, \dots, c_{\chi_G}\}$. There exists at least one such valuation and it belongs to the support of c_1 . Thus the support of c_k does not contain the support of c_1 .

□

We now move to $\text{BESTANSWER}^=$. With any undirected graph G we associate a relational structure D'_G obtained from D_G by adding a new colour c_0 in C with $L(c_0, c_i)$ for every $0 \leq i \leq m$. We define a restriction ψ of the original formula φ by disallowing c_0 in colourings : to obtain ψ it suffices to replace $L(y, x)$ in φ by $L(y, x) \wedge y \neq c_0$, and $C(x)$ by $C(x) \wedge x \neq c_0$. Thus it is still true that $\psi(c_j)$ is true in $v(D'_G)$ iff v represents a colouring of G using precisely $\{c_1, \dots, c_j\}$.

We define a new query :

$$Q'(x) := O(x) \wedge (\psi(x) \vee \exists y (O(y) \wedge L(x, y) \wedge \psi(y))) \\ \vee \\ \neg O(x) \wedge (\psi(x) \vee \exists y (O(y) \wedge x + 2 < y \wedge \psi(y))) \\ \vee \\ \neg O(x) \wedge \exists y (x \neq y \wedge L(x, y)) \wedge \forall y \forall z (E(y, z) \\ \rightarrow (y = c_0 \wedge z = c_0))$$

Note that $x + 2 < y$ is used as a shorthand, as it is definable in our language.

$Q'(c_i)$ holds in $v(D'_G)$ iff

- either i is odd and v represents an exact j -colouring of G , with j odd and $i \leq j$;
- or i is even and :
 - either v represents an exact colouring $\{c_1 \dots c_j\}$ of G with j odd, and $i + 2 < j$;
 - or v represents an exact colouring $\{c_1 \dots c_i\}$ of G ;
 - or $i < m$ and $v(\perp_j) = c_0$ for all $1 \leq j \leq m$;

The following claim completes the reduction for $\text{BESTANSWER}^=$:

Claim. $\{c_i \mid i \text{ is even}\} = \text{Best}(Q', D'_G)$ iff $\chi(G)$ is even.

Proof. . In the following, we call v_0 the unique valuation such that $v_0(\perp_j) = c_0$ for all $1 \leq j \leq m$.

First assume that χ_G is even. For all $0 < i \leq m$ odd, $\text{Supp}(c_i)$ is not maximal as $\text{Supp}(c_0) \supset \text{Supp}(c_i) \cup \{v_0\}$. Hence $\text{Best}(Q', D'_G) \subseteq \{c_i \mid i \text{ is even}\}$, so we show $\{c_i \mid i \text{ is even}\} \subseteq \text{Best}(Q', D'_G)$. The inclusion holds whenever $c_i \geq \chi(G)$, as $\text{Supp}(c_i)$ contains all valuations representing exact colorings $\{c_1 \dots c_i\}$ of G , while no other $\text{Supp}(c_j)$ with $i \neq j$ contains them. Now take $c_i < \chi(G)$ with i even, then $\text{Supp}(c_i)$ contains v_0 together with all exact odd colourings (if there are any). First assume that there exists odd exact colourings of G , so there are $\chi(G) + 1$ ones and valuations representing them are not contained in $\text{Supp}(\chi(G))$. Also $v_0 \notin \text{Supp}(c_k)$ with k odd and $k < \chi(G)$. It follows that $\text{Supp}(c_i)$, which is the union of v_0 and of all valuations representing odd exact colorings is maximal. Now assume that there is no exact odd colouring. This corresponds to the special case $\chi(G) = m$ where $\text{Supp}(c_m)$ contains only the exact colourings $\{c_1 \dots c_m\}$ of G , but not v_0 ; while $\text{Supp}(c_j) = \emptyset$ whenever j odd. In such a case $\text{Supp}(c_i) = \{v_0\}$ is also maximal.

We assume now $\chi(G)$ is odd and show $\{c_i \mid i \text{ is even}\} \neq \text{Best}(Q', D'_G)$. First notice that $\text{Supp}(c_1)$ is maximal whenever $\chi(G) = 1$, as neither $\text{Supp}(c_0)$, nor any $\text{Supp}(c_i)$ with $i \geq 2$ contain valuations representing the exact $\{c_1\}$ colourings. So we assume $\chi(G) \geq 3$, from which it follows that there exists a constant $c_{\chi(G)-3}$ in the active domain which support contains v_0 together with all valuations representing exact odd colourings. As $\text{Supp}(c_{\chi(G)-1})$ contains exactly the same set of valuations, to the exclusion of those representing $\{c_1 \dots c_{\chi(G)}\}$ colourings, it follows that $\text{Supp}(c_{\chi(G)-3}) \supset \text{Supp}(c_{\chi(G)-1})$ and so $c_{\chi(G)-1} \notin \text{Best}(Q', D'_G)$. □

Proof of Proposition 4.5

Proof. To start with we naturally extend ν to be the identity on $\text{adom}(\gamma)$. Assume first $D \models \text{equiv}_\gamma(\bar{t}, s, s')$. If $s = s'$ then $s \equiv_\gamma^\nu s'$. Now assume $s \neq s'$. Then there exist variables and/or constants $u_1, v_1 \dots u_m, v_m \in \bar{y} \cup \text{adom}(\gamma)$ with $u_i \sim_\gamma v_i$ for all i , such that $s = \nu(u_1)$, $s' = \nu(v_m)$ and $\nu(v_i) = \nu(u_{i+1})$ for all $i < m$. Clearly $u_i \sim_\gamma v_i$ implies $\nu(u_i) \equiv_\gamma^\nu \nu(v_i)$. Then $\nu(u_i) \equiv_\gamma^\nu \nu(u_{i+1})$ for all $i < m$. We conclude by transitivity that $s = \nu(u_i) \equiv_\gamma^\nu \nu(u_m) \equiv_\gamma^\nu \nu(v_m) = s'$, and therefore $s \equiv_\gamma^\nu s'$.

Assume now that $s \equiv_\gamma^\nu s'$. If $s = s'$ then clearly $D \models \text{equiv}_\gamma(\bar{t}, s, s')$. Thus assume $s \neq s'$. We proceed by induction on the number of transitive closure steps needed to derive (s, s') starting for the base relation $\{(\nu(x), \nu(w)) \mid x = w \in \gamma\}$. In the base case $(s, s') = (\nu(x), \nu(w))$ for some equality $x = w \in \gamma$. Then D satisfies the following disjunct of $\text{equiv}_\gamma(\bar{t}, s, s')$: take $u_1 = x$, $v_1 = w$, $u_i = v_i = w$ for all $i = 2..m$ (this is a disjunct since $u_i \sim_\gamma v_i$ for all $i = 1..m$). The disjunct is satisfied since $s = \nu(x) = \nu(u_1)$, $s' = \nu(w) = \nu(v_m)$, and for all $i = 1..m - 1$, $\nu(v_i) = \nu(u_{i+1}) = \nu(w)$.

In the general inductive case, there exists r such that $(r, s') = (\nu(x), \nu(w))$ for some equality $x = w$ (or $w = x) \in \gamma$, with $s \equiv_\gamma^\nu r$ derived at the previous step. By the induction hypothesis $D \models \text{equiv}_\gamma(\bar{t}, s, r)$. We can assume $s \neq r$ since otherwise (s, s') would be in the base relation. Therefore D satisfies one of the disjuncts of $\text{equiv}_\gamma(\bar{t}, s, r)$. Then there exists a sequence of $m + 1$ pairs in $\bar{y} \cup \text{adom}(\gamma)$

$$(u_1, v_1)(u_2, v_2) \dots (u_m, v_m)(u_{m+1}, v_{m+1})$$

such that

- $u_{m+1} = x$ and $v_{m+1} = w$,
- $u_i \sim_\gamma v_i, i = 1..m + 1$,

- $s = \nu(u_1), r = \nu(v_m), s' = \nu(v_{m+1}),$
- $\nu(v_i) = \nu(u_{i+1}),$ for all $i \leq m,$

We now show that from this sequence of pairs one can construct another one of exactly m pairs, $(u'_i, v'_i), i = 1..m$ still connecting s and s' , i.e. such that :

- (a) $u'_i \sim_\gamma v'_i, i = 1..m$
- (b) $s = \nu(u'_1), s' = \nu(v'_m)$
- (c) $\nu(v'_i) = \nu(u'_{i+1}),$ for all $i < m.$

The idea is to first cut the sequence $(u_i, v_i), i = 1..m + 1,$ removing at least one pair, then pad it to size m if necessary.

In order to cut the original sequence, remark that it contains $m + 1$ pairs where m is the number of \sim_γ equivalence classes. Thus there exist $i < j$ such that $u_i \sim_\gamma u_j$. We remove from the sequence all elements between u_i and v_j (excluded), the new sequence is

$$(u_1, v_1) \dots (u_{i-1}, v_{i-1})(u_i, v_j)(u_{j+1}, v_{j+1}) \dots (u_{m+1}, v_{m+1})$$

Note that this sequence satisfies (a) (b) and (c) above since $u_i \sim_\gamma u_j \sim_\gamma v_j$. Let the new sequence contain k pairs. We know $k \leq m$ because we have removed at least one pair from the original sequence (recall $i < j$). If $k < m$ we pad the sequence on the right with $m - k$ pairs (v_{m+1}, v_{m+1}) . The new sequence still satisfies (a), (b) and (c), therefore the corresponding disjunct of $\text{equiv}_\gamma(\bar{t}, s, s')$ is satisfied by D . □

Proof of Lemma 4.7

Proof. \Rightarrow Assume $v(D) \models \gamma(v(\bar{t}))$ and let s, s' such that $D \models \text{equiv}_\gamma(\bar{t}, s, s')$. By Proposition 4.5, $s \equiv_\gamma^\nu s'$. Hence either $s = s'$ and $v(s) = v(s')$ follows immediately, or there exists a sequence $r_1 \dots r_n$ of values of \bar{y} under ν such that $r_1 = s, r_n = s'$ and $\forall i < n, r_i$ and r_{i+1} are values of variables or constants of the same \equiv_γ -equivalence class. We show that $\forall i, j, v(r_i) = v(r_j)$. We proceed by induction on i and assume that $\forall j > i, v(r_j) = v(r_n)$. As r_i and r_{i+1} are values of variables or constants of the same \equiv_γ -equivalence class, there exist $u_1^i \sim_\gamma u_2^i$ in $\bar{y}\bar{c}$ such that $r_i = \nu(u_1^i)$ and $r_{i+1} = \nu(u_2^i)$. Now as $v(D) \models \gamma(v(\bar{t}))$, by definition of \sim_γ , also $v(r_i) = v(r_{i+1})$.

\Leftarrow Assume $\forall s, s', D \models \text{equiv}_\gamma(\bar{t}, s, s')$ implies $v(s) = v(s')$. We show that $\forall y, y' \in \bar{y}\bar{c}$ with $y \sim_\gamma y'$ we have $v(\nu(y)) = v(\nu(y'))$, from which it follows that $v(D) \models \gamma(v(\bar{t}))$. Let $y, y' \in \bar{y}\bar{c}$ with $y \sim_\gamma y'$ thus $\nu(y) \equiv_\gamma^\nu \nu(y')$ (here ν is naturally extended \bar{c} as the identity). By Proposition 4.5 it follows that $D \models \text{equiv}_\gamma(\bar{t}, \nu(y), \nu(y'))$ and so by assumption $v(\nu(y)) = v(\nu(y'))$ □

Proof of Proposition 4.10

Proof. \Rightarrow Assume $D \models \text{comp}_\gamma(\bar{t})$. Then $\forall c, c' \in \bar{t}$ constants, $c \equiv_\gamma^\nu c'$ implies $c = c'$. As \equiv_γ^ν is an equivalence relation, \equiv_γ^ν -equivalence classes form a partition of $\bar{t}\bar{c}$. In each \equiv_γ^ν -class there is at most one constant, so we define a valuation $v(\bar{t}\bar{c})$ mapping all nulls of a class to the unique constant of that class (or to any constant if the class does not contain any). We claim that under this valuation $v(D) \models \gamma(v(\bar{t}))$. Indeed let $w = w'$ an equality of γ , then $w \sim_\gamma w'$, which implies $D \models \text{equiv}_\gamma(\bar{t}, \nu(w), \nu(w'))$. Then by construction of $v, v(\nu(w)) = v(\nu(w'))$, which shows $v(D) \models \gamma(v(\bar{t}, \bar{c}))$.

\Leftarrow Assume $v(D) \models \gamma(v(\bar{t}, \bar{c}))$. Then partition elements of $\bar{t}\bar{c}$ according to their v values (all elements of $\bar{t}\bar{c}$ contained in one component of the partition having the v same value). Clearly any element of \bar{c} , as well as any constant of \bar{t} can only belong to the partition component associated to its value. Therefore in each partition component there is at most one constant value. Now let $s, s' \in \bar{t}$ such that $D \models \text{equiv}_\gamma(\bar{t}, s, s')$. As $v(D) \models \gamma(v(\bar{t}, \bar{c}))$ by Lemma 4.7 we have $v(D) \models v(s) = v(s')$. Moreover $D \models \text{equiv}_\gamma(\bar{t}, s, s')$ also implies that s and s' are in the same partition component with respect to v and therefore if s, s' are both constants, then $s = s'$. Hence $D \models \text{equiv}_\gamma(\bar{t}, s, s') \wedge \neg \text{Null}(s') \wedge \neg \text{Null}(s) \rightarrow s = s',$ i.e., $D \models \text{comp}_\gamma(\bar{t})$. □

Proof of Lemma A.2

Before we show Lemma 4.12 we first show that, in order to test inclusion of supports of two equality formulas, one can restrict to single valuations collapsing just what is needed.

Definition A.1 (Tight valuation). *Let $\gamma(\bar{y})$ be a conjunction of equality atoms, D a database and $\nu(\bar{y}) = \bar{t}$ an assignment over $\text{adom}(D) \cup \text{adom}(\gamma)$. A valuation v of D is called tight for ν and γ if, for all $s, s' \in \text{adom}(D) \cup \text{adom}(\gamma)$, we have $v(s) = v(s')$ iff $D \models \text{equiv}_\gamma(\bar{t}, s, s')$.*

By Lemma 4.7, any tight valuation v^* for ν and γ satisfies $v^*(D) \models \gamma(v^*(\bar{t}))$. It is also easy to see that a tight valuation for ν and γ exists whenever there is a valuation v with $v(D) \models \gamma(v(\bar{t}))$. In fact if such a v exists, by Proposition 4.10, $D \models \text{comp}_\gamma(\bar{t})$. Then for each $s \in \text{adom}(D) \cup \text{adom}(\gamma)$ there is at most one constant c such that $D \models \text{equiv}_\gamma(\bar{t}, s, c)$. In addition we associate to each equivalence class \mathcal{C} of the relation $\{(s, s') \mid D \models \text{equiv}_\gamma(\bar{t}, s, s')\}$, a new fresh constant $c_{\mathcal{C}}$ outside $\text{adom}(D) \cup \text{adom}(\gamma)$. Then a tight valuation v^* for \bar{t} and γ can be defined as follows. For $s \in \text{adom}(D)$, if $D \models \text{equiv}_\gamma(\bar{t}, s, c)$, for some constant c , then $v^*(s) = c$; otherwise $v^*(s) = c_{\mathcal{C}}$ where \mathcal{C} is the equivalence class of s . We can also characterise in terms of tight valuations the fact that for all valuations $v, v(D) \models \gamma(v(\bar{t})) \Rightarrow v(D) \models \gamma'(v(\bar{t}'))$.

Lemma A.2. Let D be a database, $\gamma(\bar{y}), \gamma'(\bar{y})$ conjunctions of equality atoms with $\text{adom}(\gamma) = \text{adom}(\gamma'), \nu(\bar{y}) = \bar{t}, \nu'(\bar{y}) = \bar{t}'$ assignments over $\text{adom}(D) \cup \text{adom}(\gamma)$ and v^* a tight valuation of D w.r.t. ν and γ . Then $v^*(D) \models \gamma'(v^*(\bar{t}'))$ iff for all valuations $v, v(D) \models \gamma(v(\bar{t}))$ implies $v(D) \models \gamma'(v(\bar{t}'))$.

Proof. \Rightarrow Assume $v^*(D) \models \gamma'(v^*(\bar{t}'))$ and let v be a valuation such that $v(D) \models \gamma(v(\bar{t}))$. We want to show $v(D) \models \gamma'(v(\bar{t}'))$. By Lemma 4.7 it is enough to show that $\forall s, s' \in \bar{t}', D \models \text{equiv}_{\gamma'}(\bar{t}', s, s')$ implies $v(D) \models v(s) = v(s')$. So let $s, s' \in \bar{t}'$ such that $D \models \text{equiv}_{\gamma'}(\bar{t}', s, s')$. As $v^*(D) \models \gamma'(v^*(\bar{t}'))$, by Lemma 4.7, $v^*(D) \models v^*(s) = v^*(s')$. Now v^* is tight w.r.t. ν and γ , so $D \models \text{equiv}_{\gamma}(\bar{t}, s, s')$. As $v(D) \models \gamma(v(\bar{t}))$, by Lemma 4.7 it follows that $v(D) \models v(s) = v(s')$.

\Leftarrow Assume for all valuations $v, v(D) \models \gamma(v(\bar{t}))$ implies $v(D) \models \gamma'(v(\bar{t}'))$. By Lemma 4.7, v^* being tight for ν and γ , we have $v^*(D) \models \gamma(v^*(\bar{t}))$ and so by our assumption $v^*(D) \models \gamma'(v^*(\bar{t}'))$. \square

Proof of Lemma 4.12

Proof. \Rightarrow Assume $D \models \text{imply}_{\gamma, \gamma'}(\bar{t}, \bar{t}') \vee \neg \text{comp}_{\gamma}(\bar{t})$. If $D \models \neg \text{comp}_{\gamma}(\bar{t})$ then by Proposition 4.10, there is no valuation v such that $v(D) \models \gamma(v(\bar{t}))$ and so the implication trivially holds. Now assume $D \models \text{imply}_{\gamma, \gamma'}(\bar{t}, \bar{t}')$, i.e. :

$$D \models \forall z z' (\text{equiv}_{\gamma'}(\bar{t}', z, z') \rightarrow \text{equiv}_{\gamma}(\bar{t}, z, z'))$$

We want to show that for all valuations v of nulls such that $v(D) \models \gamma(v(\bar{t}))$, one also has $v(D) \models \gamma'(v(\bar{t}'))$. So assume there is such a valuation, then in particular there is one which is tight w.r.t. γ and ν . By Lemma A.2 it is enough to show that $v^*(D) \models \gamma'(v^*(\bar{t}'))$, where v^* is a tight valuation of D w.r.t. γ and ν ; which by Lemma 4.7, is equivalent to showing $\forall s, s' \in \bar{t}, D \models \text{equiv}_{\gamma'}(\bar{t}', s, s')$ implies $v^*(D) \models v^*(s) = v^*(s')$. So take $s, s' \in \bar{t}$ with $D \models \text{equiv}_{\gamma'}(\bar{t}', s, s')$. By our assumption it follows that $D \models \text{equiv}_{\gamma}(\bar{t}, s, s')$. Hence by definition of tightness we have $v^*(D) \models v^*(s) = v^*(s')$.

\Leftarrow Assume for all valuations v of nulls such that $v(D) \models \gamma(v(\bar{t}))$, one also has $v(D) \models \gamma'(v(\bar{t}'))$. By Proposition 4.10, if there is no such valuation, then $D \models \neg \text{comp}_{\gamma}(\bar{t})$. So assume now there is one such valuation. This entails that in particular, there exists v^* which is tight w.r.t. γ and ν . By Lemma A.2 we thus have $v^*(D) \models \gamma'(v^*(\bar{t}'))$. Hence, by Lemma 4.7, $\forall s, s' \in \bar{t}, D \models \text{equiv}_{\gamma'}(\bar{t}', s, s')$ implies $v^*(D) \models v^*(s) = v^*(s')$. By definition of tightness $\forall s, s' \in \bar{t}, D \models \text{equiv}_{\gamma'}(\bar{t}', s, s')$ implies $D \models \text{equiv}_{\gamma}(\bar{t}, s, s')$ and so $D \models \text{imply}_{\gamma, \gamma'}(\bar{t}, \bar{t}')$. \square

Proof of Corollary 4.13

Proof. Assume $D \models \text{comp}_{\gamma}(\bar{t}) \wedge \text{imply}_{\gamma, \gamma'}(\bar{t}, \bar{t}')$, i.e., $D \models \forall z z' (\text{equiv}_{\gamma}(\bar{t}, z, z') \wedge \neg \text{Null}(z) \wedge \neg \text{Null}(z') \rightarrow z = z')$ and $D \models \forall z z' (\text{equiv}_{\gamma'}(\bar{t}', z, z') \rightarrow \text{equiv}_{\gamma}(\bar{t}, z, z'))$. Now let $s, s' \in \text{adom}(D) \cup \text{adom}(\gamma)$ with $D \models \text{equiv}_{\gamma'}(\bar{t}', s, s') \wedge \neg \text{Null}(s) \wedge \neg \text{Null}(s')$. As $D \models \text{imply}_{\gamma, \gamma'}(\bar{t}, \bar{t}')$ it follows that $D \models \text{equiv}_{\gamma}(\bar{t}, s, s')$ and so $D \models \neg \text{Null}(s) \wedge \neg \text{Null}(s') \rightarrow s = s'$ now follows from $D \models \text{comp}_{\gamma}(\bar{t})$. Hence $D \models \text{comp}_{\gamma'}(\bar{t}')$. \square

Proof of Lemma 4.14

Proof. \Leftarrow Assume $\exists i \bar{s} \bar{t} D \models q_i(\bar{s}, \bar{t}) \wedge \text{comp}_{eq_i}(\bar{r} \bar{s} \bar{t})$ and $v(D) \models eq_i(v(\bar{r} \bar{s} \bar{t}))$. By preservation of $q_i(\bar{s}, \bar{t})$ under homomorphism we have $v(D) \models q_i(v(\bar{s}, \bar{t}))$. Thus $v(D) \models q_i(v(\bar{s}, \bar{t})) \wedge eq_i(v(\bar{r} \bar{s} \bar{t}))$, i.e., $v \in \text{Supp}(Q, D, \bar{r})$.

\Rightarrow Assume $v \in \text{Supp}(Q, D, \bar{r})$, i.e., $v(D) \models Q(v(\bar{r}))$ and so there exist some $Q_i(\bar{x}) := q_i(\bar{y}, \bar{z}) \wedge eq_i(\bar{x}, \bar{y}, \bar{z})$ and tuples $\bar{a}, \bar{b} \in \text{Const}$ such that $v(D) \models q_i(\bar{a}, \bar{b}) \wedge eq_i(v(\bar{r}), \bar{a}, \bar{b})$. As $eq_i(\bar{x}, \bar{y}, \bar{z})$ contains $\bar{z} = \bar{x}$, we have $v(\bar{r}) = \bar{b}$. For each atom α in $q_i(\bar{a}, \bar{b})$, fix an arbitrary tuple β in D with $v(\beta) = \alpha$. As all variables occurring in $q_i(\bar{y}, \bar{z})$ are pairwise distinct, the set of all such β yields an assignment ν sending \bar{y}, \bar{z} to $\text{adom}(D)$ with $D \models q_i(\nu(\bar{y}), \nu(\bar{z}))$. So there exist $\bar{s} = \nu(\bar{y}), \bar{t} = \nu(\bar{z}) \in \text{adom}(D)$ with $v(\bar{s}) = \bar{a}, v(\bar{t}) = \bar{b}$ and $D \models q_i(\bar{s}, \bar{t})$. By assumption then $v(D) \models eq_i(v(\bar{r}, \bar{s}, \bar{t}))$ and by Proposition 4.10 it follows that $D \models \text{comp}_{eq_i}(\bar{r} \bar{s} \bar{t})$ and so $D \models q_i(\bar{s}, \bar{t}) \wedge \text{comp}_{eq_i}(\bar{r} \bar{s} \bar{t})$. \square

Proof of Lemma 4.15

Proof. \Rightarrow Assume $D \models Q_{\subseteq}(\bar{s}, \bar{t})$ and let $v \in \text{Supp}(Q, D, \bar{s})$ be a valuation of D . By Lemma 4.14 $\exists i \bar{a} \bar{b} D \models q_i(\bar{a} \bar{b}) \wedge \text{comp}_{eq_i}(\bar{s} \bar{a} \bar{b})$ and $v(D) \models eq_i(v(\bar{s} \bar{a} \bar{b}))$. So by our assumption there exists $j, \bar{a}' \bar{b}'$ with $D \models q_j(\bar{a}' \bar{b}') \wedge \text{imply}_{eq_i, eq_j}(\bar{s} \bar{a} \bar{b}, \bar{t} \bar{a}' \bar{b}')$ and by Corollary 4.13 $D \models \text{comp}_{eq_j}(\bar{t} \bar{a}' \bar{b}')$. Now let t_1, t_2 such that $D \models \text{equiv}_{eq_j}(\bar{t} \bar{a}' \bar{b}', t_1, t_2)$. By $D \models \text{imply}_{eq_i, eq_j}(\bar{s} \bar{a} \bar{b}, \bar{t} \bar{a}' \bar{b}')$ we have $D \models \text{equiv}_{eq_i}(\bar{s} \bar{a} \bar{b}, t_1, t_2)$ and by Lemma 4.7, $v(t_1) = v(t_2)$. But then, again by Lemma 4.7, $v(D) \models eq_i(v(\bar{t} \bar{a}' \bar{b}'))$ and by Lemma 4.14 it follows that $v \in \text{Supp}(Q, D, \bar{t})$.

\Leftarrow Assume $\text{Supp}(Q, D, \bar{s}) \subseteq \text{Supp}(Q, D, \bar{t})$ and let i, \bar{a}, \bar{b} with $D \models q_i(\bar{a}, \bar{b}) \wedge \text{comp}_{eq_i}(\bar{s}, \bar{a}, \bar{b})$. By Proposition 4.10 there exists a valuation v (that we assume w.l.o.g. to be tight) such that $v(D) \models eq_i(v(\bar{s} \bar{a} \bar{b}))$ and so by Lemma 4.14 $v \in \text{Supp}(Q, D, \bar{s})$. Hence by our assumption we also have $v \in \text{Supp}(Q, D, \bar{t})$ and so by Lemma 4.14 there exists $j, \bar{a}' \bar{b}'$ with $D \models q_j(\bar{a}' \bar{b}') \wedge \text{comp}_{eq_j}(\bar{t} \bar{a}' \bar{b}')$ and $v(D) \models eq_j(v(\bar{t} \bar{a}' \bar{b}'))$. As v is tight, by Lemma A.2 it follows from $v(D) \models eq_j(v(\bar{t} \bar{a}' \bar{b}'))$ that $\forall v$ with $v(D) \models eq_i(v(\bar{s} \bar{a} \bar{b}))$, also $v(D) \models eq_i(v(\bar{t} \bar{a}' \bar{b}'))$. Now by Proposition 4.12 $D \models \text{imply}_{eq_i, eq_j}(\bar{s} \bar{a} \bar{b}, \bar{t} \bar{a}' \bar{b}') \vee \neg \text{comp}_{eq_i}(\bar{s}, \bar{a}, \bar{b})$. But $D \models \text{comp}_{eq_i}(\bar{s}, \bar{a}, \bar{b})$, so $D \models \exists \bar{y} \bar{z} (q_j(\bar{y}, \bar{z}) \wedge \text{imply}_{eq_i, eq_j}(\bar{s} \bar{a} \bar{b}, \bar{t} \bar{y} \bar{z}))$. \square

Proof of Theorem 4.16

Proof. By Proposition 4.15 $D \models \text{best}_Q(\bar{t})$ if and only if $\forall s \text{Supp}(Q, D, \bar{t}) \subseteq \text{Supp}(Q, D, \bar{s})$ implies $\text{Supp}(Q, D, \bar{s}) \subseteq \text{Supp}(Q, D, \bar{t})$. Notice that this holds exactly whenever $\neg \exists \bar{s}$ with $\text{Supp}(Q, D, \bar{t}) \subset \text{Supp}(Q, D, \bar{s})$, i.e., whenever $\bar{t} \in \text{Best}(Q, D)$. \square

Proof of Theorem 4.19 (Sketch)

Proof. For membership, first note that one can check in Π_2^p whether $\text{Supp}(Q, D, \bar{a}) \subseteq \text{Supp}(Q, D, \bar{b})$ on input given by a database D , a UCQ Q , and tuples \bar{a} and \bar{b} . In fact in order to check $\text{Supp}(Q, D, \bar{a}) \not\subseteq \text{Supp}(Q, D, \bar{b})$ one guesses a valuation v of D , then calls an NP oracle to check $v(\bar{a}) \in Q(v(D))$ and $v(\bar{b}) \notin Q(v(D))$.

On input given by a database D , a UCQ Q , and a tuple \bar{a} one can check $\bar{a} \notin \text{Best}(Q, D)$ as follows. First guess a tuple \bar{b} over $\text{adom}(D)$ of the same arity as \bar{a} ; then, using two calls to a Σ_2^p oracle, check that $\text{Supp}(Q, D, \bar{a}) \subseteq \text{Supp}(Q, D, \bar{b})$ and $\text{Supp}(Q, D, \bar{b}) \not\subseteq \text{Supp}(Q, D, \bar{a})$.

For hardness, we reduce from $\forall \exists \forall 3DNF$, which is known to be Π_3^p -complete. We take as input a $\forall \exists \forall 3DNF$ -formula of the form

$$F := \forall z_1, \dots, z_l \exists x_1 \dots x_k \forall y_1 \dots y_p \bigvee_{i=1}^n \text{conj}_i$$

where each conj_i is a conjunction of 3 (not necessarily distinct) literals over variables $z_1, \dots, z_l, x_1, \dots, x_k, y_1, \dots, y_p$.

We construct a database D_F with $\text{adom}(D_F) = \{0, 1, \text{good}, \text{bad}\} \cup \{i, \bar{i}, \perp_i, \bar{\perp}_i, |i = 1..k\}$, and a conjunctive query $Q_F(z_1, \dots, z_l, z)$ such that $(\bar{0}, \text{good}) \in \text{Best}(Q_F, D_F)$ if and only if F is true.

D_F is of signature $\{S^4, C^2, A^2, B^3\}$ as follows :

- The extension of S and A and B are fixed and do not depend on F :
 - S contains tuple $(1, 1, 1, \text{good})$, and tuples $(b_1, b_2, b_3, \text{good})$ and $(b_1, b_2, b_3, \text{bad})$ for every $b_1, b_2, b_3 \in \{0, 1\}$ with $(b_1, b_2, b_3) \neq (1, 1, 1)$. Intuitively S encodes the possible truth assignment of each disjunct of F . Note that only the satisfying assignment (i.e. $(1, 1, 1)$) appears together with the only constant good , all the others appear both with good and bad .
 - A contains only two tuples : $(0, 1)$ and $(1, 0)$. Intuitively A will be used to encode truth values for pairs of literals $(w, \neg w)$, $w \in y_1, \dots, y_p, z_1, \dots, z_l$.
 - B contains tuples $(0, 0, \text{bad})$, $(1, 1, \text{bad})$ and tuples (b_1, b_2, good) and (b_1, b_2, bad) for every $b_1, b_2 \in \{0, 1\}$, $b_1 \neq b_2$. Intuitively B encodes assignments for pairs of literals $(w, \neg w)$, $w \in \{x_1, \dots, x_k\}$. Note that here inconsistent pairs (i.e. same truth value) are possible, but these are the only ones which do not appear together with constant good .
- The extension of C depends on F and contains tuples $\{(\perp_i, i) | i = 1..k\}$ and $\{(\bar{\perp}_i, \bar{i}) | i = 1..k\}$. Intuitively a valuation (b, i) (resp. (b, \bar{i})) of one of these tuples, with $b \in \{0, 1\}$, will encode truth value b for the literal x_i (resp. $\neg x_i$) of F .

Q_F is defined as follows. For each variable w of F , the conjunctive query Q_F will use variables w and \bar{w} (either quantified or free). For a literal α of F the corresponding variable of Q_F will be denoted as $\text{enc}(\alpha)$. More precisely if $\alpha = w$ is a positive literal then $\text{enc}(\alpha) := w$, otherwise if $\alpha = \neg w$ then $\text{enc}(\alpha) := \bar{w}$.

$$\begin{aligned} Q_F(z_1, \dots, z_l, z) := & \exists x_1, \dots, x_k, \bar{x}_1, \dots, \bar{x}_k, y_1, \dots, y_p, \bar{y}_1, \dots, \bar{y}_p, \bar{z}_1, \dots, \bar{z}_p \\ & \bigwedge_{i=1, \dots, k} B(x_i, \bar{x}_i, z) \wedge \bigwedge_{i=1, \dots, p} A(y_i, \bar{y}_i) \wedge \bigwedge_{i=1, \dots, l} A(z_i, \bar{z}_i) \wedge \\ & \bigwedge_{i=1, \dots, k} (C(x_i, i) \wedge C(\bar{x}_i, \bar{i})) \wedge \\ & \bigwedge_{(\alpha_1 \wedge \alpha_2 \wedge \alpha_3) \in F} S(\text{enc}(\alpha_1), \text{enc}(\alpha_2), \text{enc}(\alpha_3), z) \end{aligned}$$

We can prove that all tuples of the form (\bar{t}, good) (which we refer to as good tuples) have the same support. This is given by the set of all consistent boolean valuations (i.e. valuations of $\perp_i, \bar{\perp}_i$ in $\{0, 1\}$ such that $v(\perp_i) \neq v(\bar{\perp}_i)$ for all i). Moreover we can prove that if there exists a (\bar{t}, bad) whose support contains all consistent boolean valuations then the support of (\bar{t}, bad) strictly contains the support of good tuples. Therefore any good tuple (including $(\bar{0}, \text{good})$) is a best answer iff for all tuples \bar{t} there exists a consistent boolean valuation which is not in the support of (\bar{t}, bad) . We can finally show that the last holds iff F is true. \square