



**TECHNISCHE  
UNIVERSITÄT  
DRESDEN**

Technische Universität Dresden  
Institute for Theoretical Computer Science  
Chair for Automata Theory

## **LTCS-Report**

### **Most Probable Explanations for Probabilistic Database Queries (Extended Version)**

İsmail İlkan Ceylan and Stefan Borgwardt and Thomas  
Lukasiewicz

LTCS-Report 17-11

This is an extended version of the article in: Proc.  
of the 26th Int. Joint Conf. on Artificial Intelligence  
(IJCAI'17), pages 950–956, 2017.

Postal Address:  
Lehrstuhl für Automatentheorie  
Institut für Theoretische Informatik  
TU Dresden  
01062 Dresden

<http://lat.inf.tu-dresden.de>

Visiting Address:  
Nöthnitzer Str. 46  
Dresden

# Most Probable Explanations for Probabilistic Database Queries

İsmail İlkan Ceylan and Stefan Borgwardt

Faculty of Computer Science  
Technische Universität Dresden, Germany  
*firstname.lastname@tu-dresden.de*

Thomas Lukasiewicz

Department of Computer Science  
University of Oxford, UK  
*thomas.lukasiewicz@cs.ox.ac.uk*

## Abstract

Forming the foundations of large-scale knowledge bases, probabilistic databases have been widely studied in the literature. In particular, probabilistic query evaluation has been investigated intensively as a central inference mechanism. However, despite its power, query evaluation alone cannot extract all the relevant information encompassed in large-scale knowledge bases. To exploit this potential, we study two inference tasks; namely finding the *most probable database* and the *most probable hypothesis* for a given query. As natural counterparts of *most probable explanations (MPE)* and *maximum a posteriori hypotheses (MAP)* in probabilistic graphical models, they can be used in a variety of applications that involve *prediction* or *diagnosis* tasks. We investigate these problems relative to a variety of query languages, ranging from conjunctive queries to ontology-mediated queries, and provide a detailed complexity analysis.

## 1 Introduction

Research on building *large-scale probabilistic knowledge bases (PKBs)* has resulted in a number of systems including NELL [Mitchell *et al.*, 2015], Yago [Hoffart *et al.*, 2013], DeepDive [Shin *et al.*, 2015], and Google’s Knowledge Vault [Dong *et al.*, 2014]. They have been used in a wide range of areas to automatically build structured knowledge bases. Most of them are based on the long tradition and foundations of *probabilistic databases (PDBs)* [Imieliski and Lipski, 1984; Suciu *et al.*, 2011], which define probability distributions over sets of classical databases.

*Probabilistic query evaluation* is the key inference task underpinning these systems. However, PKBs encompass rich, structured knowledge, for which alternative inference mechanisms beyond probabilistic query evaluation are needed. Inspired by the *maximal posterior probability* computations in Probabilistic Graphical Models (PGMs) [Pearl, 1988; Koller and Friedman, 2009], we investigate the problem of finding *most probable explanations* for probabilistic queries to exploit the potential of such large databases to their full extent.

Computing the *maximal posterior probability* of a distinguished set of variables, given evidence about another set of

variables, is one of the key computational problems in PGMs. In its general form, this problem is studied under the name of *maximum a posteriori hypothesis*<sup>1</sup> (*MAP*), where the hypothesis refers to an instantiation of a set of distinguished variables. The *most probable explanation (MPE)* is a special case of MAP, where the distinguished variables and the evidence variables cover all variables in the model. Maximal posterior probability computations have also been lifted to relational probabilistic models such as Markov Logic Networks (MLNs) [Richardson and Domingos, 2006].

Both MPE and MAP translate to probabilistic databases in a natural way through the rich structure of queries. The *most probable database* problem (analogous to MPE), first proposed in [Gribkoff *et al.*, 2014], is the problem of determining the (classical) database with the largest probability that satisfies a given query. Intuitively, the query defines constraints on the data, and the goal is to find the most probable database that satisfies these constraints. We also introduce a more intricate notion, called *most probable hypothesis*, which only asks for *partial* databases satisfying the query (analogous to MAP). The most probable hypothesis contains only tuples that contribute to the satisfaction condition of the query, which allows to more precisely pinpoint the most likely explanations of the query.

We study the computational complexity of the corresponding decision problems, denoted by MPD and MPH, respectively, for a variety of query languages. Our results provide detailed insights about the nature of these problems. We show that the data complexity of both problems is lower for *existential queries* than for *universal queries*. As expected, MPH usually has a higher complexity than MPD.

We extend our results towards *ontology-mediated queries (OMQs)*, which enrich the well-known class of *unions of conjunctive queries* with the power of ontological rules based on Datalog<sup>±</sup> (also called existential rules) [Baget *et al.*, 2011; Calì *et al.*, 2013]. This follows the tradition of *ontology-based data access* [Poggi *et al.*, 2008], which allows us to query PDBs in a more advanced manner. For OMQs, our analysis shows that the computational complexity of these problems can change significantly w.r.t. the ontology lan-

<sup>1</sup>There exist other variants of these inference tasks [Kwisthout, 2011], and there are different naming conventions across communities; here, we use the terminology from the Bayesian Network (BN) literature [Darwiche, 2009].

guages under consideration and the complexity-theoretic assumptions employed. Our results provide tight complexity bounds for all Datalog<sup>±</sup> languages known in the literature.

## 2 Background and Motivation

We recall the basics of classical query languages and databases and briefly describe the *tuple-independent* PDB model.

### 2.1 Queries and Databases

We consider a relational vocabulary consisting of *finite* sets  $\mathbf{R}$  of *predicates*,  $\mathbf{C}$  of *constants*, and  $\mathbf{V}$  of *variables*. A *term* is a constant or a variable. An *atom* is of the form  $P(s_1, \dots, s_n)$ , where  $P$  is an  $n$ -ary predicate, and  $s_1, \dots, s_n$  are terms. A *tuple* is an atom without variables.

A *conjunctive query* (CQ) is an existentially quantified formula  $\exists \mathbf{x} \phi$ , where  $\phi$  is a conjunction of atoms, written as a comma-separated list. A *union of conjunctive queries* (UCQ) is a disjunction of CQs. If  $\phi$  is an arbitrary Boolean combination of atoms, then we call  $\exists \mathbf{x} \phi$  an *existential query* ( $\exists Q$ ), and  $\forall \mathbf{x} \phi$  a *universal query* ( $\forall Q$ ). A *first-order query* (FOQ) is an arbitrary first-order formula. A query is *Boolean* if it has no free variables.

A *database* is a finite set of tuples. The central problem studied for databases is *query evaluation*: Finding all *answers* to a query  $Q$  over a database  $\mathcal{D}$ , which are assignments of the free variables in  $Q$  to constants such that the resulting first-order formula is satisfied in  $\mathcal{D}$  in the usual sense. In the following, we consider only Boolean queries  $Q$ , and focus on the associated decision problem, i.e., deciding whether  $Q$  is satisfied in  $\mathcal{D}$ , denoted as usual by  $\mathcal{D} \models Q$ .

### 2.2 Probabilistic Databases

The most elementary probabilistic database model is based on the tuple-independence assumption. We adopt this model and refer to [Suciu *et al.*, 2011] for details and alternatives. A probabilistic database induces a set of classical databases (called *worlds*), each associated with a probability value.

Formally, a *probabilistic database* (PDB)  $\mathcal{P}$  is a finite set of *probabilistic tuples* of the form  $\langle t : p \rangle$ , where  $t$  is a tuple and  $p \in [0, 1]$ , and, whenever  $\langle t : p \rangle, \langle t : q \rangle \in \mathcal{P}$ , then  $p = q$ . A PDB  $\mathcal{P}$  assigns to every tuple  $t$  the probability  $p$  if  $\langle t : p \rangle \in \mathcal{P}$ , and otherwise the probability 0. Since we usually consider a single, fixed PDB, we denote this probability assignment simply by  $P$ . We concentrate on the well-known *possible worlds semantics*. Under the *tuple-independence assumption*,  $P$  induces the following *unique probability distribution* over classical databases  $\mathcal{D}$ :

$$P(\mathcal{D}) := \prod_{t \in \mathcal{D}} P(t) \prod_{t \notin \mathcal{D}} (1 - P(t)).$$

By the *worlds induced by*  $\mathcal{P}$  we refer to those databases  $\mathcal{D}$  with  $P(\mathcal{D}) > 0$ . Query evaluation is also enriched by probabilistic information. More formally, the *probability* of a Boolean query  $Q$  w.r.t.  $P$  is  $P(Q) := \sum_{\mathcal{D} \models Q} P(\mathcal{D})$ .

**Example 1.** Consider the PDB  $\mathcal{P}_v$  given in Figure 1 and the conjunctive query

$$Q_1 = \exists x, y \text{Veg}(x) \wedge \text{FriendOf}(x, y) \wedge \text{Veg}(y),$$

| Vegetarian |     | FriendOf |       |     |
|------------|-----|----------|-------|-----|
| alice      | 0.7 | alice    | bob   | 0.7 |
| bob        | 0.9 | alice    | chris | 0.8 |
| chris      | 0.6 | bob      | chris | 0.1 |

  

| Eats  |          |     | Meat     |     |
|-------|----------|-----|----------|-----|
| bob   | spinach  | 0.7 | shrimp   | 0.7 |
| chris | mussels  | 0.8 | mussels  | 0.9 |
| alice | broccoli | 0.2 | seahorse | 0.3 |

Figure 1: The PDB  $\mathcal{P}_v$ , where each table represents a predicate and each row is a probabilistic tuple.

through which we can ask the probability of vegetarians being friends with vegetarians. In the given PDB, *alice*, *bob*, and *chris* are all vegetarians and friends with each other (with some probability). The query

$$Q'_1 = \exists y \text{Veg}(\text{bob}) \wedge \text{FriendOf}(\text{bob}, y) \wedge \text{Veg}(y)$$

is a special case of  $Q_1$  which asks whether *bob* has vegetarian friends. Its probability can be computed as  $P(Q'_1) = 0.9 \cdot 0.1 \cdot 0.6 = 0.054$ . ■

This task is known as *probabilistic query evaluation*, and it has been studied extensively in PDBs. In a celebrated result by [Dalvi and Suciu, 2012], it has been shown that UCQs exhibit a data complexity dichotomy between  $P$  and  $\#P$  for probabilistic query evaluation. However, while dealing with large-scale knowledge bases, alternative inference mechanisms are needed. Consider again our running example, and suppose that we have observed that  $Q'_1$  is true and would like to learn what best explains this observation w.r.t. the PDB  $\mathcal{P}_v$ . We revisit this example after providing a principled approach for dealing with such tasks.

## 3 Most Probable Explanations

We investigate the problem of finding *most probable explanations* for probabilistic database queries under two different semantics. Finding the *most probable database* is to determine the world with the largest probability that satisfies a given query, as formalized by [Gribkoff *et al.*, 2014]:

**Definition 2.** The *most probable database* for a query  $Q$  over a PDB  $\mathcal{P}$  is

$$\arg \max_{\mathcal{D} \models Q} P(\mathcal{D}),$$

where  $\mathcal{D}$  ranges over all worlds induced by  $\mathcal{P}$ .

Intuitively, a PDB defines a probability distribution over exponentially many classical databases, and the most probable database is the element in this collection that has the highest probability while still satisfying the query. This can be seen as the best instantiation of a probabilistic model, and hence analogous to MPE in BNs [Darwiche, 2009].

**Example 3.** Consider again the PDB  $\mathcal{P}_v$  and the query

$$Q_2 = \forall x, y \neg \text{Veg}(x) \vee \neg \text{Eats}(x, y) \vee \neg \text{Meat}(y),$$

which defines the constraints of being a vegetarian, which is violated by the tuples  $\text{Veg}(\text{chris})$ ,  $\text{Eats}(\text{chris}, \text{mussels})$

and `Meat(mussels)`. Hence, the most probable database for  $Q_2$  cannot contain all three of them. It is easy to see that `Veg(chris)` is removed, as it has the lowest probability. Thus, the most probable database (in this case unique) contains all tuples of  $\mathcal{P}_v$  that have a probability above 0.5, except for `Veg(chris)`. Suppose that we have observed  $Q'_1$ , and we are interested in finding an explanation for this observation under the constraint of  $Q_2$ , which is specified by  $Q_3 = Q'_1 \wedge Q_2$ . In this case, the most probable database must contain the tuples `Veg(chris)` and `FriendOf(bob, chris)`, whereas `Eats(chris, mussels)` is omitted. ■

Finding the most probable database is important, as it identifies the most likely state of a PDB relative to a query. However, it has certain limitations, which are analogous to the limitations of MPE in PGMs [Koller and Friedman, 2009]. Most importantly, one is always forced to choose a *complete* database although the query usually affects only a subset of the tuples. In other words, it is usually not the case that the whole database is responsible for the goal query to be satisfied. To be able to more precisely pinpoint the explanations of a query, we introduce the following more intricate notion.

**Definition 4.** The *most probable hypothesis* for a query  $Q$  over a PDB  $\mathcal{P}$  is

$$\arg \max_{\mathcal{H} \models Q} \sum_{\mathcal{D} \models \mathcal{H}} P(\mathcal{D}),$$

where  $\mathcal{H}$  ranges over sets of tuples  $t$  and negated tuples  $\neg t$  such that  $t$  occurs in  $\mathcal{P}$ , and  $\models$  denotes the *open-world* entailment relation, i.e.,  $\mathcal{H} \models Q$  holds iff all worlds induced by  $\mathcal{P}$  that satisfy  $\mathcal{H}$  also satisfy  $Q$ .

The sum inside the maximization evaluates to the product of the probabilities of the (negated) tuples in  $\mathcal{H}$ , and hence we denote it by  $P(\mathcal{H})$ . Note that each  $\mathcal{D} \models \mathcal{H}$  must satisfy  $Q$ , and thus the most probable database is a special case of the most probable hypothesis that has to contain all tuples from  $\mathcal{P}$ . In contrast, the most probable hypothesis contains tuples only if they contribute to the satisfaction of the query.

**Example 5.** The most probable hypothesis  $\mathcal{H}$  for  $Q_3$  consists of `Veg(bob)`, `Veg(chris)`, `FriendOf(bob, chris)`, and  $\neg$ `Eats(chris, mussels)`, which yields a probability of  $P(\mathcal{H}) = 0.9 \cdot 0.6 \cdot 0.1 \cdot (1 - 0.8) = 0.0108$ . Since the most probable hypothesis contains less tuples, it is more informative than the most probable database for  $Q_3$ . ■

Whereas the most probable database represents full knowledge about all facts, which corresponds to the common closed-world assumption for (probabilistic) databases, the most probable hypothesis may leave tuples of  $\mathcal{P}$  unresolved, which can be seen as a kind of open-world assumption (although the tuples that do not occur in  $\mathcal{P}$  are still false).

### Complexity Results

We formulate the decision problems MPD and MPH and investigate their computational complexity (see Table 1). We assume familiarity with computational complexity theory, and the distinction between *data* and *combined complexity*.

**Definition 6.** Let  $Q$  be a query,  $\mathcal{P}$  a PDB, and  $p \in (0, 1]$  a threshold. We denote by MPD (resp., MPH) the problem of deciding whether there exists a database  $\mathcal{D}$  (resp., hypothesis  $\mathcal{H}$ ) that satisfies  $Q$  with  $P(\mathcal{D}) \geq p$  (resp.,  $P(\mathcal{H}) \geq p$ ).

Our first result concerns the well-known class of UCQs, and we observe that both MPD and MPH remain in *polynomial time* in the data complexity and in NP for the combined complexity.

**Theorem 7.** MPD and MPH for UCQs can be decided in polynomial time in the data complexity and are NP-complete in the combined complexity.

Intuitively, polynomial-time data complexity is ensured by the fact that UCQs are monotone queries: as a simple algorithm for MPD, consider all matches (of which there are polynomially many) of a given UCQ, and extend each match with all (negated) tuples from the PDB that have a probability greater than 0.5. This results in (polynomially many) classical databases, one of which is the most probable database. The polynomial data complexity result for MPH follows from similar arguments.

This result may seem immediate, as UCQs are monotone queries, but for MPD we can show an even stronger result that applies to all existential queries, even if negations in front of query atoms are allowed. For MPH, we show an additional hardness in the combined complexity.

**Theorem 8.** MPD for  $\exists Q$  can be decided in polynomial time in the data complexity and is NP-complete in the combined complexity. MPH for  $\exists Q$  can be decided in  $\Sigma_2^P$  in the data complexity and is  $\Sigma_3^P$ -complete in the combined complexity.

Determining the precise data complexity of MPH for  $\exists Q$  is left as an open problem.

For universally quantified queries, MPD becomes harder even in the data complexity. Given the rich structure of the query, it becomes possible to encode non-deterministic choices over all possible databases. It is important to note that a similar hardness result was also obtained in [Gribkoff et al., 2014] (although on a different query).

**Theorem 9.** MPD for  $\forall Q$  is NP-complete in the data complexity.

We observe a similar phenomenon for MPH, where the increase in the complexity is even more dramatic.

**Theorem 10.** MPH for  $\forall Q$  is  $\Sigma_2^P$ -complete in the data complexity.

Note that also MPE and MAP differ in terms of computational complexity: in BNs, the former is NP-complete [Shimony, 1994; Littman, 1999] and the latter  $\text{NP}^{\text{PP}}$ -complete [Park and Darwiche, 2004a]. Differently, MPH remains in  $\Sigma_2^P \subseteq \text{NP}^{\text{PP}}$ , since computing the probability of the hypothesis can be done in polynomial time due to the tuple-independence assumption. As we will examine later, allowing ontological rules to induce correlations on the tuples results in different complexities, which are more comparable to PGMs, which can express conditional dependencies between their variables. Our final result for this section concerns the combined complexity of both problems for  $\forall Q$ .

**Theorem 11.** MPD and MPH for  $\forall Q$  are  $\Sigma_2^P$ -complete in the combined complexity.

If we consider arbitrary FOQs, it is easy to show that the data complexity remains the same as for  $\forall Q$ s, and the combined complexity is still PSPACE, as for classical databases.

Moreover, all our combined complexity results for PDBs hold even in the case where the arity of the predicates is bounded by some constant, commonly known as the *bounded-arity (ba) combined complexity* (see Table 1). The reason is that our hardness results use only predicates of a bounded arity.

## 4 Most Probable Explanations for OMQs

We now study the problems in the presence of ontological rules, which add deductive power to PDBs. The benefits of using of ontologies for large-scale PKB completion are well known [Jung and Lutz, 2012; Borgwardt *et al.*, 2017]. From a broader perspective, extending tuple-independent PDBs with logical rules is an old idea, aiming to induce correlations on a logical level, and thus to relax the tuple independence, resulting in very powerful formalisms [Poole, 1993; 1997].

In the remainder of this paper, we restrict ourselves to UCQs, but in exchange consider additional knowledge encoded through an ontology. In the simplest case, we can formulate *negative constraints (NCs)* like

$$\forall x, y \text{ Veg}(x) \wedge \text{Eat}(x, y) \wedge \text{Meat}(y) \rightarrow \perp,$$

which imposes the same constraints as  $Q_2$ . NCs are a special case of denial constraints over databases [Staworko and Chomicki, 2010]. We also allow to formulate more general ontological knowledge in the form of *tuple-generating dependencies (TGDs)*. For instance, the first-order formulas

$$\begin{aligned} \forall x, y \text{ FriendOf}(x, y) &\rightarrow \text{FriendOf}(y, x), \\ \forall x \text{ Veg}(x) &\rightarrow \exists y \text{ Knows}(x, y) \wedge \text{Veg}(y) \end{aligned}$$

are TGDs stating that the friend relation is symmetric and that every vegetarian knows another vegetarian. In particular, TGDs can express the well-known inclusion dependencies and join dependencies from database theory. Taking all parts together, we are talking about *ontology-mediated queries*, which are UCQs in combination with a so-called Datalog<sup>±</sup> ontology, i.e., a finite set of NCs and TGDs [Cali *et al.*, 2012]. We will pay particular attention to the case where only NCs are allowed, as this is closest to the constraints over PDBs that we described in the previous examples, and it gives us a baseline for the computational complexity.

More formally, an *NC*  $\nu$  is an FO formula  $\forall \mathbf{x} \varphi(\mathbf{x}) \rightarrow \perp$ , where  $\varphi(\mathbf{x})$  is a conjunction of atoms, called the *body* of  $\nu$ , and  $\perp$  is the truth constant *false*. A *TGD*  $\sigma$  is an FO formula  $\forall \mathbf{x} \varphi(\mathbf{x}) \rightarrow \exists \mathbf{y} P(\mathbf{x}, \mathbf{y})$ , where  $\varphi(\mathbf{x})$  is a conjunction of atoms, called the *body* of  $\sigma$ , and  $P(\mathbf{x}, \mathbf{y})$  is an atom, called the *head* of  $\sigma$ . A (Datalog<sup>±</sup>) *program* (or *ontology*)  $\Sigma$  is a finite set of NCs and TGDs.<sup>2</sup> An *ontology-mediated query (OMQ)*  $(Q, \Sigma)$  consists of a program  $\Sigma$  and a Boolean UCQ  $Q$ .

For the semantics, we extend the vocabulary by an infinite set  $\mathbf{N}$  of *nulls*. An instance  $I$  is a possibly infinite set of tuples that may additionally contain nulls; it *satisfies* a TGD or NC  $\sigma$  if  $I \models \sigma$ , where  $\models$  denotes standard FO entailment.  $I$  satisfies a program  $\Sigma$ , written  $I \models \Sigma$ , if  $I$  satisfies each formula in  $\Sigma$ . The set  $\text{mods}(\mathcal{D}, \Sigma)$  of *models* of a database  $\mathcal{D}$  relative to a program  $\Sigma$  is  $\{I \mid I \supseteq \mathcal{D} \text{ and } I \models \Sigma\}$ . A database  $\mathcal{D}$  is *consistent* w.r.t.  $\Sigma$  if  $\text{mods}(\mathcal{D}, \Sigma)$  is non-empty. The OMQ  $(Q, \Sigma)$

is *entailed* by  $\mathcal{D}$ , denoted  $\mathcal{D} \models (Q, \Sigma)$ , if  $I \models Q$  holds for all  $I \in \text{mods}(\mathcal{D}, \Sigma)$ . Observe that consistency of  $\mathcal{D}$  w.r.t.  $\Sigma$  can thus be written as  $\mathcal{D} \not\models (\perp, \Sigma)$ , or equivalently,  $\mathcal{D} \not\models (Q_{\perp}, \Sigma^+)$ , where  $Q_{\perp}$  is the UCQ obtained from the disjunction of the bodies of all NCs in  $\Sigma$ , and  $\Sigma^+$  contains all TGDs of  $\Sigma$ . The *probability* of an OMQ over a PDB is defined as usual, by summing up over the consistent worlds that entail the query.

In general, the entailment problem is undecidable [Beeri and Vardi, 1981], which motivated syntactic fragments of TGDs (but not the NCs) [Cali *et al.*, 2012; Baget *et al.*, 2011; Krötzsch and Rudolph, 2011; Cali *et al.*, 2013; Fagin *et al.*, 2005]. We focus on only a few of these classes here (see Table 1). If there are no TGDs at all, i.e.,  $\Sigma^+ = \emptyset$ , then we denote the query language by  $\text{OMQ}_{\text{NC}}$ . One of the most important classes, *guarded* programs ( $\text{OMQ}_{\text{G}}$ ) allow only such TGDs that have a body atom (the *guard*) that contains all body variables. In the special case of *linear* programs ( $\text{OMQ}_{\text{L}}$ ), the body may consist of only one atom. In *frontier-guarded* programs ( $\text{OMQ}_{\text{FG}}$ ), only the frontier variables, i.e., those shared by the body and the head, need to be guarded. We pay particular attention to the class of *guarded and full* programs ( $\text{OMQ}_{\text{GF}}$ ), which does not allow existentially quantified variables, and is one of the least expressive Datalog<sup>±</sup>-based query languages with a P-complete data complexity. At the upper end of the expressivity spectrum, we observe classes like *weakly guarded* ( $\text{OMQ}_{\text{WG}}$ ) programs, whose data complexity is already EXP-complete. From a theoretical perspective, the class of *acyclic* programs ( $\text{OMQ}_{\text{A}}$ ), which do not allow cyclic dependencies between predicates, is of interest, because it has a non-deterministic combined complexity, which leads to a different behavior.

A key paradigm in OMQ answering is the *FO-rewritability* of queries: an OMQ  $(Q, \Sigma)$  is *FO-rewritable* if there exists a Boolean UCQ  $Q_{\Sigma}$  such that, for all consistent databases  $\mathcal{D}$ , it holds that  $\mathcal{D} \models (Q, \Sigma)$  iff  $\mathcal{D} \models Q_{\Sigma}$ . Intuitively, the rewritten query  $Q_{\Sigma}$  can be answered directly over the database, without referring to the Datalog<sup>±</sup> program. A class  $\text{OMQ}_{\text{X}}$  is *FO-rewritable* if all its OMQs are FO-rewritable; in this case, the OMQ entailment problem for  $\text{OMQ}_{\text{X}}$  has a data complexity of AC<sup>0</sup>. Of the classes mentioned above, only  $\text{OMQ}_{\text{NC}}$ ,  $\text{OMQ}_{\text{L}}$ , and  $\text{OMQ}_{\text{A}}$  have this property.

In addition to data complexity, *ba*-combined complexity, and combined complexity, for OMQs we also consider the *fixed-program (fp) combined complexity*, where the Datalog<sup>±</sup> program is viewed as fixed, but both (P)DB and query are considered to be part of the input. The *ba*-combined complexity is of interest for Description Logics [Baader *et al.*, 2003], some of which can be considered to be Datalog<sup>±</sup> languages with arity at most 2.

### 4.1 Most Probable Databases for OMQs

For ontology-mediated queries, models are restricted to those that are consistent with the ontology. In other words, the constraints are considered separately from the (existential) query; thus, the definitions of MPD and MPH have to be adapted accordingly. The main difference is that we now consider only the *consistent* worlds induced by the PDB, and thus maximize only over consistent worlds.

<sup>2</sup>We omit the universal quantifiers in TGDs and NCs, and use commas (instead of  $\wedge$ ) for conjoining atoms for ease of presentation.

| Query Languages   | Most Probable Database |          |                 |                 | Most Probable Hypothesis |                  |                  |                  |
|-------------------|------------------------|----------|-----------------|-----------------|--------------------------|------------------|------------------|------------------|
|                   | data                   | fp-comb. | ba-comb.        | comb.           | data                     | fp-comb.         | ba-comb.         | comb.            |
| UCQ               | in P                   | —        | NP              | NP              | in P                     | —                | NP               | NP               |
| $\exists$ Q       | in P                   | —        | NP              | NP              | in $\Sigma_2^P$          | —                | $\Sigma_3^P$     | $\Sigma_3^P$     |
| $\forall$ Q       | NP                     | —        | $\Sigma_2^P$    | $\Sigma_2^P$    | $\Sigma_2^P$             | —                | $\Sigma_2^P$     | $\Sigma_2^P$     |
| FOQ               | NP                     | —        | PSPACE          | PSPACE          | $\Sigma_2^P$             | —                | PSPACE           | PSPACE           |
| OMQ <sub>NC</sub> | NP                     | NP       | $\Sigma_2^P$    | $\Sigma_2^P$    | PP                       | NP <sup>PP</sup> | NP <sup>PP</sup> | NP <sup>PP</sup> |
| OMQ <sub>L</sub>  | NP                     | NP       | $\Sigma_2^P$    | PSPACE          | PP                       | NP <sup>PP</sup> | NP <sup>PP</sup> | PSPACE           |
| OMQ <sub>A</sub>  | NP                     | NP       | P <sup>NE</sup> | P <sup>NE</sup> | PP                       | NP <sup>PP</sup> | P <sup>NE</sup>  | P <sup>NE</sup>  |
| OMQ <sub>GF</sub> | NP                     | NP       | $\Sigma_2^P$    | EXP             | NP <sup>PP</sup>         | NP <sup>PP</sup> | NP <sup>PP</sup> | EXP              |
| OMQ <sub>G</sub>  | NP                     | NP       | EXP             | 2EXP            | NP <sup>PP</sup>         | NP <sup>PP</sup> | EXP              | 2EXP             |
| OMQ <sub>FG</sub> | NP                     | NP       | 2EXP            | 2EXP            | NP <sup>PP</sup>         | NP <sup>PP</sup> | 2EXP             | 2EXP             |
| OMQ <sub>WG</sub> | EXP                    | EXP      | EXP             | 2EXP            | EXP                      | EXP              | EXP              | 2EXP             |

Table 1: Complexity results for MPD and MPH for a wide range of queries: all listed results are original contributions of this paper except the data complexity for  $\forall$ Q. This latter result has been strengthened towards the weaker representation OMQ<sub>NC</sub>.

**Definition 12.** The *most probable database* for an OMQ  $(Q, \Sigma)$  over a PDB  $\mathcal{P}$  is

$$\arg \max_{\mathcal{D}=(Q, \Sigma), \text{mods}(\mathcal{D}, \Sigma) \neq \emptyset} P(\mathcal{D}).$$

As before, we consider the decision problem MPD, which checks for the existence of a world  $\mathcal{D}$  as above that exceeds a given probability threshold  $p$ . Our complexity results for this problem are summarized in the lower left part of Table 1.

A naive approach to solve MPD is to guess the database  $\mathcal{D}$ , and then check that it entails the given OMQ, does *not* entail the query  $(Q_{\perp}, \Sigma^+)$ , and exceeds the probability threshold. Since the probability can be computed in polynomial time, the problem can be decided by an NP Turing machine using an oracle to check OMQ entailment.

**Theorem 13.** *If entailment for OMQ<sub>X</sub> is in a complexity class C, then MPD for OMQ<sub>X</sub> is in NP<sup>C</sup> (under the same complexity assumptions).*

By a reduction from 3-colorability, we show that MPD is NP-hard already in the data complexity, even if we only use NCs, i.e., the query and the positive program  $\Sigma^+$  are empty. This strengthens our previous result about  $\forall$ Qs, since NCs can be expressed by universal queries, but are not allowed to use negated atoms.

**Theorem 14.** *MPD for OMQ<sub>NC</sub> is NP-hard in the data complexity, even for  $Q = \top$ .*

We show a matching upper bound even in the fp-combined complexity, by reconsidering the approach from Theorem 13. Since here the query  $(Q_{\perp}, \Sigma^+)$  is fixed, for the non-entailment check, we can refer to the *data complexity* of OMQ entailment. Assuming that this is possible in P, and further that the fp-combined complexity does not exceed NP, then the whole test can be done by a single non-deterministic Turing machine in polynomial time. This insight yields an upper bound of NP for most of the classes we consider.

**Theorem 15.** *If entailment for OMQ<sub>X</sub> is in NP in the fp-combined complexity and in P in the data complexity, then MPD for OMQ<sub>X</sub> is in NP in the fp-combined complexity.*

Under ba-combined complexity assumptions, we observe an increase in complexity, which intuitively comes from the fact that the query  $(Q_{\perp}, \Sigma^+)$  is not fixed anymore.

**Theorem 16.** *MPD for OMQ<sub>NC</sub> is  $\Sigma_2^P$ -hard in ba-combined complexity, even for  $Q = \top$ .*

We obtain a matching upper bound from Theorem 13 if  $C = \text{NP}$ . Most of the remaining hardness results follow from the complexity of classical OMQ entailment, since an OMQ  $(Q, \Sigma)$  is entailed by a consistent database  $\mathcal{D}$  iff the most probable database w.r.t.  $\{\{t : 1\} \mid t \in \mathcal{D}\}$  has probability 1. In combination with Theorem 13, this yields tight complexity results for large deterministic classes like PSPACE or EXP. This leaves open only the case of OMQ<sub>A</sub>, for which entailment is NEXP-complete in the (ba-)combined complexity, which yields an upper bound of  $\text{NP}^{\text{NEXP}} = \text{P}^{\text{NEXP}} = \text{P}^{\text{NE}}$  due to Theorem 13 and results from [Hemachandra, 1989]. We show that this bound is tight, using a reduction from a P<sup>NEXP</sup>-complete version of the *tiling problem* [Fürer, 1983].

**Theorem 17.** *MPH for OMQ<sub>A</sub> is P<sup>NE</sup>-hard in the ba-combined complexity.*

## 4.2 Most Probable Hypotheses for OMQs

Similarly to MPD, we have to update the definition of MPH to take into account only consistent worlds.

**Definition 18.** The *most probable hypothesis* for an OMQ  $(Q, \Sigma)$  over a PDB  $\mathcal{P}$  is

$$\arg \max_{\mathcal{H}=(Q, \Sigma)} \sum_{\substack{\mathcal{D} \supseteq \mathcal{H} \\ \text{mods}(\mathcal{D}, \Sigma) \neq \emptyset}} P(\mathcal{D}),$$

where  $\mathcal{H}$  is a set of (non-probabilistic) tuples  $t$  occurring in  $\mathcal{P}$ .

Since we consider only monotone queries (i.e., UCQs), we do not have to include negated tuples in the hypothesis.

To solve MPH, one can guess a hypothesis, and then check whether it entails the query, and whether the probability mass of its consistent extensions exceeds the given threshold. The latter part can be done by a PP Turing machine with an oracle

for OMQ entailment (by normalizing the probability of the worlds such that the threshold is exactly 0.5). The oracle can be used also for the initial entailment check.

**Theorem 19.** *If entailment for  $OMQ_X$  is in a complexity class  $C$ , then MPH for  $OMQ_X$  is in  $NP^{PP^C}$  (under the same complexity assumptions).*

For any  $C \subseteq PH$ , this yields  $NP^{PP^{PH}} = NP^{PP^C} = NP^{PP}$  as an upper bound, due to a result from [Toda, 1989]. Except for PP, all other upper bounds in the lower right part of Table 1 also follow from this observation. For  $C = NEXP$ , the whole  $PP^{NEXP}$  computation can be done by an NEXP oracle, and hence we again obtain  $NP^{NEXP} = P^{NE}$  in this case.

On the other hand, we can transfer most of the lower bounds from MPD by a simple reduction.

**Theorem 20.** *If MPD for  $OMQ_X$  is hard for a complexity class  $C$ , then MPH for  $OMQ_X$  is also hard for  $C$  (under any complexity measure except for the data complexity).*

We now discuss the more interesting cases below PSPACE. For FO-rewritable classes of OMQs, we obtain a data complexity of only PP. This is due to the fact that TGDs can be compiled into the query, and the observation of Theorem 8 that existential queries can be processed using polynomially many steps. Nevertheless, in each step, we have to compute a non-trivial sum over the *consistent* worlds, which can be done in PP. Finally, the polynomially many PP checks can be compiled into a single PP operation [Beigel *et al.*, 1995].

**Theorem 21.** *If entailment for  $OMQ_X$  is in  $AC^0$  in the data complexity, then MPH for  $OMQ_X$  is PP-complete in the data complexity.*

To frame this result, we show  $NP^{PP}$ -hardness both for  $OMQ_{GF}$  (the prototypical non-FO-rewritable class) in the data complexity, and for  $OMQ_{NC}$  in the fp-combined complexity. We use two similar reductions from a problem in the polynomial-time counting hierarchy [Wagner, 1986]. In the first one, we rely on the power of guarded and full TGDs; in the second, we use a non-fixed query and a careful choice of probabilities to simulate the effect of these TGDs.

**Theorem 22.** *MPH is  $NP^{PP}$ -hard for  $OMQ_{GF}$  in the data complexity, and for  $OMQ_{NC}$  in the fp-combined complexity.*

We also observe a dichotomy in the data complexity of MPH, which follows from the dichotomy for probabilistic query evaluation [Dalvi and Suciu, 2012]. Note that the PP-hardness holds only under Turing reductions.

**Lemma 23 (Dichotomy).** *For FO-rewritable classes  $OMQ_X$ , MPH is either in P or PP-hard in the data complexity.*

Our results apply to all decidable Datalog<sup>±</sup> languages from the literature, due to the generic nature of our theorems and the known complexity results for OMQ entailment [Cali *et al.*, 2012; Baget *et al.*, 2011; Krötzsch and Rudolph, 2011; Cali *et al.*, 2013]. For example, *full* sets of TGDs (F) behave like  $OMQ_{GF}$ , and *linear full* (LF) and *acyclic full* (AF) sets of TGDs exhibit the same complexities as  $OMQ_L$ .

## 5 Related Work

Our work is inspired by the *maximal posterior probability* computations in PGMs [Pearl, 1988; Koller and Friedman,

2009]. It is well-known that they are a form of *abduction*, which clearly also applies to the problems studied here. Maximal posterior probability computations are central in PGMs and in statistical relational learning. However, analogous problems have not been studied in depth in the context of PDBs. To our knowledge, the only work in this direction is on most probable databases [Gribkoff *et al.*, 2014], which we use as a starting point.

Probabilistic query evaluation has been investigated in depth in the literature [Imieliski and Lipski, 1984; Suciu *et al.*, 2011]. In fact, it is known to be either in P or #P-hard for UCQs in the data complexity [Dalvi and Suciu, 2012]. Extensions of PDBs with ontological rules, in particular with Datalog<sup>±</sup>, are also well-covered in the literature [Gottlob *et al.*, 2013; Borgwardt *et al.*, 2017]. The focus of these works is on probabilistic query answering.

Most of our data complexity results are covered by the classes NP,  $\Sigma_2^P$ , PP, and  $NP^{PP}$ . Though intractable, they are at the core of many important problems, which motivated a body of work tailored towards scalable algorithms for these classes. There is an immediate connection between MPE and weighted MAX-SAT [Sang *et al.*, 2007]. Similarly, advances in knowledge compilation [Park and Darwiche, 2004b; Pipatsrisawat and Darwiche, 2009] and approximate model counting [Chakraborty *et al.*, 2016; Fremont *et al.*, 2017] are tailored to achieve optimal, scalable algorithms for problems in classes such as PP and  $NP^{PP}$ . Both MPD and MPH can be cast into their propositional variants using the lineage representation of queries, which gives immediate access to such algorithms. However, there is need for future work in this direction, as grounding is not an optimal way of handling these problems; performing inference directly on FO-structures is known to be more efficient.

## 6 Summary and Outlook

We studied two inference tasks for PDBs; namely finding the most probable database and the most probable hypothesis for a given query. The focus of this paper was to determine the precise complexity of these problems, and we provided a detailed analysis for these problems relative to a variety of query languages, ranging from conjunctive queries to ontology-mediated queries. The main focus of future work is on the one hand to extend these results to other classes (such as equality generating dependencies) and on the other hand to obtain even more refined results (such as classification results).

## Acknowledgments

This work was supported by DFG in the CRC 912 (HAEC), GRK 1907 (RoSI), and the project BA 1122/19-1 (GoAsQ), and by the UK EPSRC grants EP/J008346/1, EP/L012138/1, EP/M025268/1, and by The Alan Turing Institute under the EPSRC grant EP/N510129/1.

## References

- [Baader *et al.*, 2003] Franz Baader, Diego Calvanese, Deborah L. McGuinness, Daniele Nardi, and Peter F. Patel-Schneider, editors. *The Description Logic Handbook: Theory, Implementation, and Applications*. Cambridge University Press, 2003.

- [Baget *et al.*, 2011] Jean-François Baget, Marie-Laure Mugnier, Sebastian Rudolph, and Michaël Thomazo. Walking the complexity lines for generalized guarded existential rules. In *Proc. IJCAI*, 2011.
- [Beeri and Vardi, 1981] Catriel Beeri and Moshe Y. Vardi. The implication problem for data dependencies. In *Proc. ICALP*, 1981.
- [Beigel *et al.*, 1995] Richard Beigel, Nick Reingold, and Daniel Spielman. PP is closed under intersection. *JCSS*, 50(2):191–202, 1995.
- [Borgwardt *et al.*, 2017] Stefan Borgwardt, İsmail İlkan Ceylan, and Thomas Lukasiewicz. Ontology-mediated queries for probabilistic databases. In *Proc. AAAI*, 2017.
- [Cali *et al.*, 2012] Andrea Cali, Georg Gottlob, and Thomas Lukasiewicz. A general Datalog-based framework for tractable query answering over ontologies. *JWS*, 14:57–83, 2012.
- [Cali *et al.*, 2013] Andrea Cali, Georg Gottlob, and Michael Kifer. Taming the infinite chase: Query answering under expressive relational constraints. *JAIR*, 48:115–174, 2013.
- [Chakraborty *et al.*, 2016] Supratik Chakraborty, Kuldeep S. Meel, and Moshe Y. Vardi. Algorithmic improvements in approximate counting for probabilistic inference: From linear to logarithmic SAT calls. In *Proc. IJCAI*, 2016.
- [Dalvi and Suciu, 2012] Nilesh Dalvi and Dan Suciu. The dichotomy of probabilistic inference for unions of conjunctive queries. *J. ACM*, 59(6):1–87, 2012.
- [Darwiche, 2009] Adnan Darwiche. *Modeling and Reasoning with Bayesian Networks*. Cambridge University Press, 2009.
- [Dong *et al.*, 2014] Xin Luna Dong, Evgeniy Gabrilovich, Jeremy Heitz, Wilko Horn, Ni Lao, Kevin Patrick Murphy, Thomas Strohmann, Shaohua Sun, and Wei Zhang. Knowledge Vault: A web-scale approach to probabilistic knowledge fusion. In *Proc. SIGKDD*, 2014.
- [Fagin *et al.*, 2005] Ronald Fagin, Phokion G. Kolaitis, Renée J. Miller, and Lucian Popa. Data exchange: Semantics and query answering. *TCS*, 336(1):89–124, 2005.
- [Fremont *et al.*, 2017] Daniel Fremont, Markus Rabe, and Sanjit Seshia. Maximum model counting. In *Proc. AAAI*, 2017.
- [Fürer, 1983] Martin Fürer. The computational complexity of the unconstrained limited domino problem (with implications for logical decision problems). In *Logic and Machines*, 1983.
- [Gill, 1977] John T. Gill. Computational complexity of probabilistic Turing machines. *SIAM J. on Computing*, 6(4):675–695, 1977.
- [Gottlob *et al.*, 2013] Georg Gottlob, Thomas Lukasiewicz, Maria Vanina Martinez, and Gerardo I. Simari. Query answering under probabilistic uncertainty in Datalog+/- ontologies. *AMAI*, 69(1):37–72, 2013.
- [Gribkoff *et al.*, 2014] Eric Gribkoff, Guy Van den Broeck, and Dan Suciu. The most probable database problem. In *Proc. BUDA*, 2014.
- [Hemachandra, 1989] Lane A. Hemachandra. The strong exponential hierarchy collapses. *JCSS*, 39(3):299–322, 1989.
- [Hoffart *et al.*, 2013] Johannes Hoffart, Fabian M. Suchanek, Klaus Berberich, and Gerhard Weikum. YAGO2: A spatially and temporally enhanced knowledge base from Wikipedia. *AIJ*, 194:28–61, 2013.
- [Imieliski and Lipski, 1984] Tomasz Imieliski and Witold Lipski. Incomplete information in relational databases. *J. ACM*, 31(4):761–791, 1984.
- [Jung and Lutz, 2012] Jean Christoph Jung and Carsten Lutz. Ontology-based access to probabilistic data with OWL QL. In *Proc. ISWC*, 2012.
- [Koller and Friedman, 2009] Daphne Koller and Nir Friedman. *Probabilistic Graphical Models: Principles and Techniques*. MIT Press, 2009.
- [Krötzsch and Rudolph, 2011] Markus Krötzsch and Sebastian Rudolph. Extending decidable existential rules by joining acyclicity and guardedness. In *Proc. IJCAI*, 2011.
- [Kwisthout, 2011] Johan Kwisthout. Most probable explanations in Bayesian networks: Complexity and tractability. *IJAR*, 52(9):1452–1469, 2011.
- [Littman, 1999] Michael L. Littman. Initial experiments in stochastic satisfiability. In *Proc. AAAI*, 1999.
- [Mitchell *et al.*, 2015] Tom Mitchell, William Cohen, Estevam Hruschka, Partha Talukdar, Justin Betteridge, Andrew Carlson, Bhavana Dalvi Mishra, Matthew Gardner, Bryan Kisiel, Jayant Krishnamurthy, Ni Lao, Kathryn Mazaitis, Thahir Mohamed, Ndapa Nakashole, Emmanouil Platanios, Alan Ritter, Mehdi Samadi, Burr Settles, Richard Wang, Derry Wijaya, Abhinav Gupta, Xinlei Chen, Abulhair Saparov, Malcolm Greaves, and Joel Welling. Never-ending learning. In *Proc. AAAI*, 2015.
- [Park and Darwiche, 2004a] James D. Park and Adnan Darwiche. Complexity results and approximation strategies for MAP explanations. *JAIR*, 21(1):101–133, 2004.
- [Park and Darwiche, 2004b] James D. Park and Adnan Darwiche. A differential semantics for jointree algorithms. *AIJ*, 156(2):197–216, 2004.
- [Pearl, 1988] Judea Pearl. *Probabilistic Reasoning in Intelligent Systems*. Morgan Kaufmann, 1988.
- [Pipatsrisawat and Darwiche, 2009] Knot Pipatsrisawat and Adnan Darwiche. A new d-DNNF-based bound computation algorithm for functional E-MAJSAT. In *Proc. IJCAI*, 2009.
- [Poggi *et al.*, 2008] Antonella Poggi, Domenico Lembo, Diego Calvanese, Giuseppe De Giacomo, Maurizio Lenzerini, and Riccardo Rosati. Linking data to ontologies. *JDS*, 10:133–173, 2008.
- [Poole, 1993] David Poole. Probabilistic Horn abduction and Bayesian networks. *AIJ*, 64(1):81–129, 1993.
- [Poole, 1997] David Poole. The independent choice logic for modelling multiple agents under uncertainty. *AIJ*, 94(1):7–56, 1997.
- [Richardson and Domingos, 2006] Matthew Richardson and Pedro Domingos. Markov logic networks. *ML*, 62(1):107–136, 2006.
- [Sang *et al.*, 2007] Tian Sang, Paul Beame, and Henry Kautz. A dynamic approach to MPE and weighted MAX-SAT. In *Proc. IJCAI*, 2007.
- [Shimony, 1994] Eyal Solomon Shimony. Finding MAPs for belief networks is NP-hard. *AIJ*, 68(2):399–410, 1994.
- [Shin *et al.*, 2015] Jaeho Shin, Sen Wu, Feiran Wang, Christopher De Sa, Ce Zhang, and Christopher Ré. Incremental knowledge base construction using DeepDive. In *Proc. VLDB*, 2015.
- [Staworko and Chomicki, 2010] Sławomir Staworko and Jan Chomicki. Consistent query answers in the presence of universal constraints. *Inf. Syst.*, 35(1):1–22, 2010.
- [Suciu *et al.*, 2011] Dan Suciu, Dan Olteanu, Christopher Ré, and Christoph Koch. *Probabilistic Databases*. Morgan & Claypool, 2011.
- [Toda, 1989] Seinosuke Toda. On the computational power of PP and  $\oplus P$ . In *Proc. SFCS*, 1989.



- [Torán, 1991] Jacobo Torán. Complexity classes defined by counting quantifiers. *J. ACM*, 38(3):753–774, 1991.
- [Valiant, 1979] Leslie Gabriel Valiant. The complexity of computing the permanent. *TCS*, 8(2):189–201, 1979.
- [Wagner, 1986] Klaus W. Wagner. The complexity of combinatorial problems with succinct input representation. *Acta Inf.*, 23(3):325–356, 1986.

## A Complexity Background

In the scope of this paper, we use the standard assumption that the probability values are rational. Apart from the well known classes such as NP,  $\Sigma_2^P$ ,  $\Sigma_3^P$ , PSPACE, EXP, we refer to the probabilistic complexity class PP [Gill, 1977]. Briefly, PP defines the set of languages recognized by a polynomially time-bounded non-deterministic Turing machine that accepts an input if and only if *more than half* of the computation paths are accepting [Torán, 1991]. Intuitively, the class PP can be seen as the decision counterpart of the counting class #P [Valiant, 1979]. In fact, it is known that  $P^{PP} = \#P$  by Toda’s theorem. Toda’s theorem is actually stronger than this well-known result: it holds that  $PP^{PH} \subseteq P^{PP}$  [Toda, 1989].

Observe that the original dichotomy for probabilistic query evaluation by Dalvi and Suciu (2012) is formulated using the class #P. However, #P-hardness is shown as usual using FP-Turing reductions, which translates to a P versus PP dichotomy *under polynomial-time Turing reductions*, for the associated decision problem. Thus, the dichotomy result from Lemma 23 holds under these assumptions. All our other results, however, hold even under standard many-one reductions.

We also refer to the class  $NP^{PP}$ , which encompasses an important class of problems in AI, combining search and counting problems. The complexity classes relevant to our results relate to standard classes as follows:

$$NP \subseteq \Sigma_2^P, PP, PH \subseteq P^{PP} = \#P \subseteq NP^{PP} \subseteq PSPACE \subseteq EXP \\ EXP \subseteq NEXP \subseteq NP^{NEXP} = P^{NEXP} = P^{NE} \subseteq 2EXP$$

## B Proof of Theorem 7

The claim for MPD follows from Theorem 8.

For MPH, observe first that, since the query is monotone, the databases extending the hypothesis  $\mathcal{H}$  satisfy the query only if  $\mathcal{H}$  (extended with all tuples that have probability 1) is already a match for the query. This means that the hypothesis must be a subset of a ground instance of one of the disjuncts of the UCQ Q. In the data complexity, there are only polynomially many such hypotheses, and their probabilities can be computed in polynomial time, which yields an upper bound of P. In the combined complexity, we can guess such a hypothesis, compare its probability to the threshold in polynomial time, and verify that it satisfies one of the disjuncts of Q.  $\square$

## C Proof of Theorem 8

**Upper bounds** We can assume that the query is of the form  $\exists x \phi$ , where  $\phi$  is a disjunction of conjunctions of atoms and negated atoms. It is easy to see that the most probable database must contain an instance of one of these conjunctions, and from the remaining tuples  $t$  in  $\mathcal{P}$  simply contains

those with  $P(t) > 0.5$ . Hence, for the data complexity, one can enumerate all (polynomially many) such instances and check if one of the resulting databases exceeds the probability threshold. For the combined complexity, we can guess the instance in non-deterministic polynomial time, complete it to a database using the polynomial approach described above, and compare its probability to the threshold in polynomial time.

For the MPH, we can guess a hypothesis  $\mathcal{H}$  in polynomial time, and then check the threshold and verify that, for all databases extending  $\mathcal{H}$  that are induced by the PDB, the query is satisfied. The latter is an ordinary database query evaluation problem, and hence can be done in NP in the combined complexity (in P in the data complexity). Thus, the entailment check can be done using a  $\Pi_2^P$  (coNP) oracle, which puts the overall complexity of this algorithm at  $\Sigma_3^P$  ( $\Sigma_2^P$ ).

**Combined complexity lower bound for MPH** We now show  $\Sigma_3^P$ -hardness. Consider a quantified Boolean formula of the form

$$\Phi = \exists u_1, \dots, u_n \forall v_1, \dots, v_m \exists w_1, \dots, w_k \phi,$$

where  $\phi = \phi_1 \wedge \dots \wedge \phi_l$  is in CNF. Checking validity of such formulas is known to be  $\Sigma_3^P$ -complete. We provide a reduction from this problem. We start by defining a PDB  $\mathcal{P}_\Phi$  as follows:

- For every existentially quantified variable  $u_j$ ,  $\mathcal{P}_\Phi$  contains the tuples  $\langle L(u_j, 0) : 0.5 \rangle$  and  $\langle L(u_j, 1) : 0.5 \rangle$ .
- For every universally quantified variable  $v_j$ ,  $\mathcal{P}_\Phi$  contains the tuple  $\langle S(v_j) : 0.5 \rangle$ .
- For the existentially quantified variables  $w_j$ , we only add the two tuples  $\langle T(1) : 1 \rangle$  and  $\langle T(0) : 0 \rangle$  to  $\mathcal{P}_\Phi$ .

We now construct a query from the given propositional formula:

- For every clause  $\phi_j$ , we construct a disjunction  $\psi_j$  by replacing the propositional variables with appropriate L-atoms, S-atoms, or T-atoms. For instance, for  $\phi_j = u_1 \vee \neg u_2 \vee v_3 \vee \neg w_6$ , we use

$$\psi_j = L(u_1, 1) \vee L(u_2, 0) \vee S(v_3) \vee \neg T(w_6).$$

The atom  $L(u_1, 1)$  indicates that the existentially quantified variable  $u_1$  should be true, and dually for  $L(u_2, 0)$ . The atom  $S(v_3)$  says that the universally quantified variable  $v_3$  should be true, and we negate such atoms if the variable is negated in the clause. Finally,  $\neg T(w_6)$  expresses a similar condition on the remaining existentially quantified variable. Note that here  $w_6$  is a variable, while  $u_1, u_2, v_3$  are constants from  $\mathcal{P}_\Phi$ .

- For all existentially quantified variables  $u_i$ , we additionally define the formulas

$$\psi_{e_i} = (\neg L(u_i, 0) \vee \neg L(u_i, 1)) \wedge (L(u_i, 0) \vee L(u_i, 1)).$$

- We finally construct the existential query

$$Q_\Phi = \exists w_1, \dots, w_n \bigvee_{1 \leq j \leq l} \psi_j \wedge \bigwedge_{1 \leq i \leq n} \psi_{e_i}.$$

Intuitively,  $\mathcal{P}_\Phi$  together with  $\psi_1, \dots, \psi_l$  encodes the satisfaction condition of the formula, while  $\psi_{e_1}, \dots, \psi_{e_n}$  force the hypothesis to contain the L-tuples for the existentially quantified variables  $u_1, \dots, u_n$ .

**Claim**  $\Phi$  is valid iff there exists a hypothesis  $\mathcal{H}$  for  $Q_\Phi$  over  $\mathcal{P}_\Phi$  such that  $P(\mathcal{H}) \geq (0.5)^{2n}$  and  $\mathcal{H} \models Q_\Phi$ .

Assume that  $\Phi$  is valid. Then there exists a valuation  $\nu$  for  $u_1, \dots, u_n$ , such that all extensions  $\tau$  of  $\nu$  to  $v_1, \dots, v_m$  admit an extension  $\iota$  to  $w_1, \dots, w_k$  that satisfies  $\phi$ . We choose the hypothesis  $\mathcal{H}$  as follows: We add  $L(u_j, 1)$  and  $\neg L(u_j, 0)$  to  $\mathcal{H}$  if  $\nu(u_j) = 1$ , and otherwise we add  $\neg L(u_j, 1)$  and  $L(u_j, 0)$  to  $\mathcal{H}$ . Hence, we have  $P(\mathcal{H}) = (0.5)^{2n}$ . Now, for any database  $\mathcal{D}$  that is induced by  $\mathcal{P}_\Phi$  and extends  $\mathcal{H}$ , we must have  $T(1) \in \mathcal{D}$  and  $T(0) \notin \mathcal{D}$ , and for each universally quantified variable  $v_j$ ,  $\mathcal{D}$  fixes a truth value via the tuple  $S(v_j)$ . This defines an extension  $\tau_{\mathcal{D}}$  of  $\nu$  by setting  $\tau_{\mathcal{D}}(v_j) = 1$  iff  $S(v_j) \in \mathcal{D}$ . By assumption, we know that there is an extension  $\iota_{\mathcal{D}}$  of  $\tau_{\mathcal{D}}$  to the remaining variables such that  $\phi$  is satisfied. We can hence satisfy  $Q_\Phi$  in  $\mathcal{D}$  by mapping each  $w_j$  to 1 if  $\iota_{\mathcal{D}}(w_j) = 1$ , and to 0 otherwise. This shows that  $\mathcal{H} \models Q_\Phi$ .

Conversely, suppose that  $P(\mathcal{H}) \geq (0.5)^{2n}$  for some hypothesis  $\mathcal{H}$ , i.e.,  $\mathcal{H} \models Q_\Phi$ . By construction of the query, for every existentially quantified variable  $u_j$ ,  $\mathcal{H}$  must contain the two tuples  $L(u_j, 0)$  and  $L(u_j, 1)$ , one positively and the other negatively, to satisfy the query. This implies that, for  $\mathcal{H}$  to achieve the threshold  $(0.5)^{2n}$ , it must contain exactly two L-tuples for each existentially quantified variable  $u_j$  (and it can contain some deterministic tuples). To show that  $\Phi$  is valid, we now define the partial assignment  $\nu$  such that  $\nu(u_j) = 1$  if  $L(u_j, 1) \in \mathcal{H}$  (thus,  $\neg L(u_j, 0) \in \mathcal{H}$ ), and  $\nu(u_j) = 0$  otherwise. Consider now any extension  $\tau$  of  $\nu$  to  $v_1, \dots, v_m$ , and construct the extension  $\mathcal{D}$  of  $\mathcal{H}$  by adding  $S(v_j)$  iff  $\tau(v_j) = 1$ . We must also add  $T(1)$  to  $\mathcal{D}$ . It is easy to see that  $P(\mathcal{D}) > 0$ . Hence, by assumption there must be an instantiation  $\mu$  of the query variables  $w_1, \dots, w_k$  that satisfies  $Q_\Phi$ . We define the extension  $\iota_\mu$  of  $\tau$  to the propositional variables  $w_1, \dots, w_k$  by setting  $\iota_\mu(w_j) = 1$  iff  $T(\mu(w_j)) \in \mathcal{D}$ . Due to the construction of  $Q_\Phi$ , this extension must satisfy  $\phi$ , and hence we know that  $\Phi$  is valid.  $\square$

## D Proof of Theorem 9

**Hardness** We provide a reduction from the satisfiability of propositional 3CNF formulas. Let  $\phi = \bigwedge_i \phi_i$  be a propositional formula in 3CNF. We define the following query

$$\begin{aligned} Q_{\text{SAT}} := \forall x, y, z ( & L(x) \vee L(y) \vee L(z) \vee R_1(x, y, z) \wedge \\ & (\neg L(x) \vee L(y) \vee L(z) \vee R_2(x, y, z)) \wedge \\ & (\neg L(x) \vee \neg L(y) \vee L(z) \vee R_3(x, y, z)) \wedge \\ & (\neg L(x) \vee \neg L(y) \vee \neg L(z) \vee R_4(x, y, z)) \end{aligned}$$

We define the PDB  $\mathcal{P}_\Phi$  as follows.

- For each propositional variable  $u_j$ , the PDB  $\mathcal{P}_\Phi$  contains the tuple  $\langle L(u_j) : 0.5 \rangle$ .
- The clauses  $\phi_j$  are described with the help of the predicates  $R_1, \dots, R_4$ , each of which corresponds to one type of clause. For example, if we have  $\phi_j = u_1 \vee \neg u_2 \vee \neg u_4$ , we add the tuple  $\langle R_3(u_4, u_2, u_1) : 0 \rangle$  to  $\mathcal{P}_\Phi$ , which enforces via  $Q_{\text{SAT}}$  that either  $\neg L(u_4)$ ,  $\neg L(u_2)$  or  $L(u_1)$  holds. All other tuples  $R_i(u_k, u_l, u_m)$  that do not correspond in such a way to one of the clauses we add with probability 1 to  $\mathcal{P}_\Phi$ .

**Claim** The formula  $\phi$  is satisfiable iff there exists a world  $\mathcal{D}$  induced by  $\mathcal{P}_\Phi$  such that  $P(\mathcal{D}) \geq (0.5)^n$  and  $\mathcal{D} \models Q_{\text{SAT}}$ , where  $n$  is the number of variables appearing in  $\phi$ .

Suppose that  $\phi$  is satisfiable and let  $\nu$  be such a satisfying assignment. We define a world  $\mathcal{D}$  such that it contains all the tuples of the form  $L(u_j)$  iff  $\nu(u_j) \mapsto 1$  in the given assignment. Moreover,  $\mathcal{D}$  contains all the tuples which are assigned the probability 1 in  $\mathcal{P}_\Phi$ . It is easy to see that  $\mathcal{D}$  is one of the worlds induced by  $\mathcal{P}_\Phi$ . Observe further that  $\mathcal{P}_\Phi$  contains  $n$  non-deterministic tuples, each with 0.5 probability. By this argument, the probability of  $\mathcal{D}$  is clearly  $(0.5)^n$ . It only remains to show that  $\mathcal{D} \models Q_{\text{SAT}}$ , which is easy to verify.

For the other direction, let  $\mathcal{D} \models Q_{\text{SAT}}$  and  $P(\mathcal{D}) \geq (0.5)^n$  for some world  $\mathcal{D}$ . We define an assignment  $\nu$  by setting the truth value of  $u_j$  to 1 if  $L(u_j) \in \mathcal{D}$ , and to 0 otherwise. Every world contains exactly one assignment for every variable, by our construction. Thus, the assignment  $\nu$  is well-defined. It is easy to verify that  $\nu \models \phi$ .

**Membership** We can guess a world  $\mathcal{D}$ , verify that it satisfies the query, and compare its probability in polynomial time. Furthermore, only the first step is non-deterministic.  $\square$

## E Proof of Theorem 10

**Hardness** We consider the following query  $Q = Q_{\text{VAL}} \wedge Q'$  where

$$\begin{aligned} Q_{\text{VAL}} := \forall x, y, z ( & (L(x) \wedge L(y) \wedge L(z)) \vee R_1(x, y, z) \wedge \\ & ((\neg L(x) \wedge L(y) \wedge L(z)) \vee R_2(x, y, z)) \wedge \\ & ((\neg L(x) \wedge \neg L(y) \wedge L(z)) \vee R_3(x, y, z)) \wedge \\ & ((\neg L(x) \wedge \neg L(y) \wedge \neg L(z)) \vee R_4(x, y, z)) \end{aligned}$$

and

$$Q' := \forall x (\neg M(x) \wedge L(x)) \vee (M(x) \wedge \neg L(x)) \vee K(x)$$

Consider a quantified Boolean formula of the form  $\Phi = \exists u_1, \dots, u_n \forall v_1, \dots, v_m \phi$ , where  $\phi$  is in 3DNF. Checking validity of such formulas is known to be  $\Sigma_2^P$ -complete. We provide a reduction from this problem. First, we define a PDB  $\mathcal{P}_\Phi$  such that

- for each variable  $u$  that appears in  $\phi$ ,  $\mathcal{P}_\Phi$  contains the tuple  $\langle L(u) : 0.5 \rangle$ .
- for every existentially quantified variable  $u_j$ ,  $\mathcal{P}_\Phi$  contains the tuple  $\langle M(u_j) : 0.5 \rangle$ .
- for every universally quantified variable  $v_j$ ,  $\mathcal{P}_\Phi$  contains the tuple  $\langle K(v_j) : 1 \rangle$ .
- for every conjunction in  $\phi$ ,  $\mathcal{P}_\Phi$  contains an  $R_i$ -tuple with probability 0, e.g., for  $\phi_j = u_1 \wedge v_3 \wedge \neg u_5$  the tuple  $\langle R_2(u_5, v_3, u_1) : 0 \rangle$ .
- $\mathcal{P}_\Phi$  contains the remaining  $R_i$ -tuples with probability 1.

In this construction,  $Q_{\text{VAL}}$  encodes the 3DNF and  $Q'$  helps us to distinguish between the existentially and universally quantified variables through the  $K$ - and  $M$ -tuples.

**Claim**  $\Phi$  is valid iff there exists a hypothesis  $\mathcal{H}$  for  $Q$  over  $\mathcal{P}_\Phi$  such that  $P(\mathcal{H}) \geq (0.5)^{2n}$  and  $\mathcal{H} \models Q$ .

Suppose that  $\Phi$  is valid. Then, there exists a valuation  $\nu$  of  $u_1, \dots, u_n$ , such that all valuations  $\tau$  that extend this partial valuation (by assigning truth values to  $v_1, \dots, v_m$ ) satisfy  $\phi$ . We define a hypothesis  $\mathcal{H}$  depending on  $\nu$  as follows. For all assignments  $u_j \mapsto 1$  in  $\nu$ , we add  $L(u_j)$  to  $\mathcal{H}$ ; if, on the other hand,  $u_j \mapsto 0$  in  $\nu$  we add  $\neg L(u_j)$  to  $\mathcal{H}$ . Moreover, to satisfy the query  $Q$ , for every  $L(u_j) \in \mathcal{H}$ , we add  $\neg M(u_j)$  to  $\mathcal{H}$ , and analogously, for every  $\neg L(u_j) \in \mathcal{H}$ , we add  $M(u_j)$  to  $\mathcal{H}$ . By this construction, there are clearly  $2n$  tuples in  $\mathcal{H}$ , each of which has the probability  $0.5$  in  $\mathcal{P}_\Phi$ . Hence, it holds that  $P(\mathcal{H}) = (0.5)^{2n}$ . Finally, it is sufficient to observe that all databases  $\mathcal{D}$  which extend  $\mathcal{H}$  must satisfy the query  $Q$ , as every such database is in one-to-one correspondence with a valuation  $\tau$  that extends  $\nu$ .

For the other direction, we assume that there exists a hypothesis  $\mathcal{H}$  for  $Q$  over  $\mathcal{P}_\Phi$  such that  $P(\mathcal{H}) \geq (0.5)^{2n}$  and  $\mathcal{H} \models Q$ . This implies that  $\mathcal{H}$  contains at most  $2n$  tuples that have probability  $0.5$  in  $\mathcal{P}_\Phi$  (and possibly some deterministic tuples). Furthermore, since  $\mathcal{H} \models Q'$ ,  $\mathcal{H}$  contains each  $M$ -tuple either positively or negatively, and it also contains the complementary  $L$ -tuple. Since these are already  $2n$  tuples,  $\mathcal{H}$  cannot contain any  $L$ -tuples for the universally quantified variables  $v_j$ . We can thus define a valuation  $\nu$  for  $u_1, \dots, u_n$  simply by setting  $u_j \mapsto 1$  if  $L(u_j) \in \mathcal{H}$  and  $u_j \mapsto 0$  if  $\neg L(u_j) \in \mathcal{H}$ . It is easy to see that the extensions  $\tau$  of  $\nu$  are in one-to-one correspondence with the databases that extend  $\mathcal{H}$ , and that  $\phi$  evaluates to true for all of these assignments.

**Membership** This is an immediate consequence of Theorem 11, which proves  $\Sigma_2^P$  membership even for combined complexity.  $\square$

## F Proof of Theorem 11

**Hardness for MPD** We again consider the validity problem for  $\Phi = \exists u_1, \dots, u_n \forall v_1, \dots, v_m \phi$ , where  $\phi$  is a propositional formula. We use the PDB  $\mathcal{P}_\Phi$  that contains all tuples  $\langle L(u_j) : 0.5 \rangle$  for the existentially quantified variables  $u_j$ , and the two additional tuples  $\langle L(0) : 0 \rangle$  and  $\langle L(1) : 1 \rangle$ . The query is  $Q_\Phi = \forall v_1, \dots, v_m \psi$ , where  $\psi$  is obtained from  $\phi$  by replacing all propositional variables  $u$  by  $L(u)$ . Here, the existentially quantified variables are viewed as constants, and the universally quantified variables are still universally quantified, but now range over the database constants instead of the truth values true and false.

**Claim**  $\Phi$  is valid iff there exists a database  $\mathcal{D}$  induced by  $\mathcal{P}_\Phi$  that satisfies  $Q_\Phi$  such that  $P(\mathcal{D}) \geq (0.5)^n$ .

Assume that  $\Phi$  is valid. Then there is a valuation  $\nu$  for  $u_1, \dots, u_n$  such that all extensions  $\tau$  to  $v_1, \dots, v_m$  satisfy  $\phi$ . We choose the database  $\mathcal{D}$  such that  $L(u_j) \in \mathcal{D}$  iff  $\nu(u_j) = 1$ ,  $L(1) \in \mathcal{D}$ , and  $L(0) \notin \mathcal{D}$ . Hence,  $P(\mathcal{D}) = (0.5)^n$ . To show that  $\mathcal{D} \models Q_\Phi$ , consider an assignment  $\sigma$  for  $v_1, \dots, v_m$  in  $Q_\Phi$ , which assigns a constant to each  $v_j$ . We define the extension  $\tau$  of  $\nu$  to the propositional variables  $v_1, \dots, v_m$  by setting  $\tau(v_j) = 1$  iff  $L(\sigma(v_j)) \in \mathcal{D}$ . Since by assumption  $\tau$  satisfies  $\phi$ , it must be the case that  $\sigma$  satisfies  $\psi$ .

Conversely, let  $\mathcal{D} \models Q_\Phi$  be such that  $P(\mathcal{D}) \geq (0.5)^n$ . We define the valuation  $\nu$  for  $u_1, \dots, u_n$  by setting  $\nu(u_j) = 1$  iff  $L(u_j) \in \mathcal{D}$ . Consider now any extension  $\tau$  of  $\nu$  to  $v_1, \dots, v_m$ . We obtain a corresponding assignment for the variables in  $Q_\Phi$  by setting  $\sigma(v_j) = \tau(v_j)$ . Since  $\mathcal{D}$  must satisfy  $L(1)$  and cannot satisfy  $L(0)$ , we know from  $\mathcal{D} \models Q_\Phi$  that  $\tau$  satisfies  $\phi$ .

**Membership** Consider a non-deterministic Turing machine  $M$  with a (co)NP oracle: Given a PDB  $\mathcal{P}$ , a universal query  $Q$ , and a threshold  $p \in (0, 1]$ , we can decide whether there exists a hypothesis  $\mathcal{H}$  such that  $P(\mathcal{H}) \geq p$  by first guessing a hypothesis  $\mathcal{H}$ , and verifying whether (i)  $P(\mathcal{H}) \geq p$  and (ii) for all databases  $\mathcal{D}$  that extend  $\mathcal{H}$  and are induced by  $\mathcal{P}$ , it holds that  $\mathcal{D} \models Q$ . Verification of (i) can be done in deterministic polynomial time and (ii) can be done in coNP (the complement is equivalent to the existence of an extension  $\mathcal{D}$  and a valuation for the query variables that falsifies  $Q$ ).

We can use the same machine for MPD, except that we require the initial guess to be a full database.  $\square$

## G Proof of Theorem 14

We provide a reduction from the well-known 3-colorability problem: Given an undirected graph  $G = (V, E)$ , decide whether the nodes of  $G$  are 3-colorable. We define the PDB  $\mathcal{P}_G$  as follows. For all edges  $(u, v) \in E$ , we add the tuple  $E(u, v)$  with probability 1, and for all nodes  $u \in V$ , we add the tuples  $V(u, 1), V(u, 2), V(u, 3)$ , each with probability 0.7. In this encoding, the tuples  $V(u, 1), V(u, 2), V(u, 3)$  correspond to different colorings of the same node  $u$ . The conditions for 3-colorability are encoded through a set  $\Sigma$  containing only negative constraints (that do not depend on  $G$ ). First, we ensure that each node is assigned *at most* one color:

$$\begin{aligned} V(x, 1), V(x, 2) &\rightarrow \perp \\ V(x, 1), V(x, 3) &\rightarrow \perp \\ V(x, 2), V(x, 3) &\rightarrow \perp \end{aligned}$$

Similarly, we enforce that the neighboring nodes are not assigned the same color:

$$E(x, y), V(x, c), V(y, c) \rightarrow \perp$$

Finally, we define the query  $Q = \tau$ .

**Claim**  $G$  is 3-colorable iff there is a consistent database  $\mathcal{D}$  with  $P(\mathcal{D}) \geq (0.7 \cdot 0.3 \cdot 0.3)^{|V|}$  ( $\mathcal{D} \models Q$  is trivially fulfilled).

Suppose that there exists a database  $\mathcal{D}$  with a probability of at least  $(0.7 \cdot 0.3 \cdot 0.3)^{|V|}$  that satisfies all NCs in  $\Sigma$ . Then, for every node  $u \in V$ ,  $\mathcal{D}$  must contain exactly one tuple  $V(u, c)$  for some color  $c \in \{1, 2, 3\}$ . Recall that *at most* was ensured by the first three NCs; hence, in order to achieve the given threshold, *at least* one of these tuples must be present. This yields a unique coloring for the nodes. Furthermore, since  $\mathcal{D}$  must contain all tuples corresponding to the edges of  $G$ , and  $\mathcal{D}$  satisfies the last NC, we conclude that  $G$  is 3-colorable.

Suppose that  $G$  is 3-colorable. Then, for a valid coloring, we define a DB  $\mathcal{D}$  that contains all tuples that correspond to the edges, and add all tuples  $V(u, c)$  where  $c$  is the color of  $u$ . It is easy to see that  $\mathcal{D}$  is consistent with  $\Sigma$  and  $P(\mathcal{D}) = (0.7 \cdot 0.3 \cdot 0.3)^{|V|}$ , which concludes the proof.  $\square$

## H Proof of Theorem 16

Consider the formula

$$\Phi = \exists x_1, \dots, x_n \forall y_1, \dots, y_m \phi,$$

where  $\phi = \phi_1 \vee \dots \vee \phi_k$  is a propositional formula in 3DNF. We encode the truth values of the variables  $x_i$  using the tuples  $\langle V(x_i, 0) : 0.9 \rangle$ ,  $\langle V(x_i, 1) : 0.9 \rangle$  in the PDB  $\mathcal{P}_\Phi$ , which are constrained by the NC  $V(x, 0) \wedge V(x, 1) \rightarrow \perp$ . The idea is that the worlds  $\mathcal{D}$  with maximal probability must satisfy exactly one of  $V(x_i, 0)$ ,  $V(x_i, 1)$  for each  $x_i$ , thereby representing a truth value assignment for  $\mathbf{x}$ .

We illustrate the encoding of the conjunctions in  $\Phi$  on the example of  $\phi_j = x_1 \wedge \neg y_4 \wedge \neg x_3$ . We introduce a predicate  $C_j(t, y_4)$  that describes the truth value  $t$  of  $\phi_j$  as a function of the truth value of  $y_4$ , which of course depends on the truth values chosen for  $x_1$  and  $x_3$  via the V-tuples. For example, if  $y_4$  is true, then  $\phi_j$  must be false, which is expressed by the tuple  $C_j(0, 1)$ . As above, we add all tuples  $\langle C_j(t, y_4) : 0.9 \rangle$  with  $t, y_4 \in \{0, 1\}$  to  $\mathcal{P}_\Phi$ . The idea is again that, for a certain value of  $y_4$ , exactly one of the two tuples  $C_j(0, y_4)$ ,  $C_j(1, y_4)$  will be true in the chosen world  $\mathcal{D}$  (not both and not neither).

This behavior is enforced by a series of NCs. For example, if  $x_1$  is true and  $x_3$  is false, then  $C_j(1, 0)$  must be true, i.e.,  $C_j(0, 0)$  must be false:

$$V(x_1, 1) \wedge V(x_3, 0) \wedge C_j(0, 0) \rightarrow \perp$$

Moreover, making  $y_4$  true will always make the conjunction false, regardless of the values of  $x_1$  and  $x_3$ :

$$C_j(1, 1) \rightarrow \perp$$

Finally, if either  $x_1$  is false or  $x_3$  is true, then it is impossible to satisfy  $\phi_j$ , regardless of the value of  $y_4$ :

$$V(x_1, 0) \wedge C_j(1, y_4) \rightarrow \perp$$

$$V(x_3, 1) \wedge C_j(1, y_4) \rightarrow \perp$$

These four NCs together ensure that, in every world  $\mathcal{D}$  of maximal probability (which determines a fixed truth value assignment for  $x_1, \dots, x_n$ ), the satisfied tuples  $C_j(t, y_4)$  determine the truth value of  $\phi_j$  as a function of the truth value of  $y_4$ . For example, if  $\mathcal{D}$  satisfies  $V(x_1, 0)$  and  $V(x_3, 1)$  (and hence neither  $V(x_1, 1)$  nor  $V(x_3, 0)$ ), then it must satisfy also  $C_j(0, 0)$  and  $C_j(0, 1)$  (but neither  $C_j(1, 0)$  nor  $C_j(1, 1)$ ).

In general, the predicate  $C_j$  is of arity  $1 + |\mathbf{y}_j|$ , where  $\mathbf{y}_j$  are the variables among  $y_1, \dots, y_m$  that occur in  $\phi_j$ . This results in tuples of the form  $C_j(t, \mathbf{y}_j)$ . For each clause, four NCs like the ones above enforce that exactly one of  $C_j(0, \mathbf{y}_j)$  and  $C_j(1, \mathbf{y}_j)$  is true, for each valuation of  $\mathbf{y}_j$ , and this describes exactly the behavior of  $\phi_j$ .

Finally, we add the NC

$$\bigwedge_{j=1}^k C_j(0, \mathbf{y}_j) \rightarrow \perp$$

to express that every valuation of the variables  $y_1, \dots, y_m$  must satisfy at least one conjunction  $\phi_j$ . Hence, the existence of a consistent world  $\mathcal{D}$  with  $P(\mathcal{D}) \geq (0.09)^\ell$ , where  $\ell = n + \sum_{j=1}^k 2^{|\mathbf{y}_j|}$ , is equivalent to the validity of  $\Phi$ . Since  $\ell$  is bounded by  $n + 8k$ , the threshold can be written using linearly many bits. Observe that we do not need a query (we can choose  $Q = \top$ ) nor any TGDs, and the maximal arity of the used predicates is 4.  $\square$

## I Proof of Theorem 17

We give a reduction from the following *extended tiling problem*, which is  $P^{\text{NEXP}}$ -complete: Given a triple  $(m, TP_1, TP_2)$  of an integer  $m$  in unary and tiling problems  $TP_1$  and  $TP_2$  for the exponential square  $2^n \times 2^n$ , does, for every initial condition  $w = w_1 \dots w_m$  for the first row,  $TP_1$  have no solution with  $w$ , or does  $TP_2$  have some solution with  $w$ ? Here, the *tiling problem* [Fürer, 1983] is defined as follows: Let  $T$  be a set of square tile types,  $H, V \subseteq T \times T$  be the horizontal and vertical compatibility relations, respectively, and  $n$  be an integer in unary. A  $2^n \times 2^n$  tiling is a function  $f: \{1, \dots, 2^n\} \times \{1, \dots, 2^n\} \rightarrow T$  such that  $(f(i, j), f(i, j+1)) \in H$  and  $(f(i, j), f(i+1, j)) \in V$ , for each  $i$  and  $j$ . An instance of the tiling problem is a tuple  $(T, H, V, n)$ , and the question is whether a  $2^n \times 2^n$  tiling exists.

The construction makes use of the result that any instance  $TP$  of tiling the  $2^n \times 2^n$ -square, given  $n$ , relations  $H$  and  $V$ , and an initial tiling condition  $w = w_1 \dots w_m$ , is reducible to OMQ answering with acyclic TGDs in polynomial time such that the query  $(\text{Tiling}, \Sigma^{TP, |w|})$  is entailed by  $\mathcal{D}^{TP} \cup \mathcal{D}^w$ , where  $\Sigma^{TP, |w|}$  is constructed from  $TP$  and  $|w|$ ,  $\mathcal{D}^{TP}$  from  $TP$ , and  $\mathcal{D}^w = \{\text{Init}_j(w_j) \mid 1 \leq j \leq m\}$ , iff  $TP$  has a solution with  $w$ .

We define the PDB  $\mathcal{P}$  as the set of all tuples  $\langle t : 1 \rangle$  such that  $f \in \mathcal{D}^{TP_1} \cup \mathcal{D}^{TP_2}$ , all tuples  $\langle \text{Init}_j(d) : 0.5 \rangle$  such that  $d \in T$  and  $1 \leq j \leq m$ , and the tuples  $\langle \text{Tiling}_2(0) : 0.5 \rangle$  and  $\langle \text{Tiling}_2(1) : 0.5 \rangle$ . We define the program  $\Sigma$  as the union of  $\Sigma_r^{TP_1, |w|}$  and  $\Sigma_r^{TP_2, |w|}$ , obtained from  $\Sigma^{TP_1, |w|}$  and  $\Sigma^{TP_2, |w|}$  by renaming all derived predicates  $P$  to  $P_1$  and  $P_2$ , respectively, and by replacing  $\text{Tiling}_2$  by  $\text{Tiling}_2(1)$ , together with the NCs  $\text{Tiling}_2(0), \text{Tiling}_2(1) \rightarrow \perp$  and all  $\text{Init}_j(d), \text{Init}_j(d') \rightarrow \perp$  such that  $d, d' \in T$  and  $1 \leq j \leq m$ .

Then, with the UCQ  $Q$  defined as

$$\exists x_1, \dots, x_m \bigwedge_{j=1}^m \text{Init}_j(x_j) \wedge \text{Tiling}_1 \wedge \text{Tiling}_2(0),$$

we obtain that there is a world induced by  $\mathcal{P}$  that satisfies  $(Q, \Sigma)$  iff  $(m, TP_1, TP_2)$  is a no-instance of the extended tiling problem.  $\square$

## J Proof of Theorem 20

Consider an OMQ  $(Q, \Sigma)$  over a PDB  $\mathcal{P}$ , for which we want to find a consistent world of probability  $> q$ . We enforce that the most probable hypothesis needs to make a choice for all tuples in  $\mathcal{P}$ , by replacing each tuple  $\langle P(\mathbf{t}) : p \rangle$  by two new tuples  $\langle P'(\mathbf{t}, 1) : p \rangle$  and  $\langle P'(\mathbf{t}, 0) : 1 - p \rangle$ , where  $P'$  is a fresh predicate that increases the arity of  $P$  by one, and 0 and 1 are fresh constants. We denote the resulting PDB by  $\mathcal{P}'$ . As the query, we use

$$Q' := Q^1 \wedge \bigwedge_{\langle P(\mathbf{t}):p \rangle \in \mathcal{P}} \exists z P'(\mathbf{t}, z),$$

where  $Q^1$  is obtained from  $Q$  by replacing all tuples  $P(\mathbf{x})$  by  $P'(\mathbf{x}, 1)$ . The query  $Q'$  is equivalent to a UCQ that is of size polynomial in the size of  $Q$  and  $\mathcal{P}$ . Finally, in the program  $\Sigma$  we similarly replace each tuple  $P(\mathbf{x})$  by  $P'(\mathbf{x}, 1)$ , and add the NCs  $P'(\mathbf{x}, 1) \wedge P'(\mathbf{x}, 0) \rightarrow \perp$  for each of the new predicates  $P'$ .

The resulting program  $\Sigma'$  still satisfies the constraints of fp/ba-combined complexity, and is in the same class as  $\Sigma$ .

We show that the most probable hypothesis for  $(Q', \Sigma')$  over  $\mathcal{P}'$  exceeds the threshold  $q^2$  iff the most probable database for  $(Q, \Sigma)$  over  $\mathcal{P}$  exceeds  $q$ . Assume that the latter is the case, and the database  $\mathcal{D}$  has a probability of  $d > q$  w.r.t.  $\mathcal{P}$ . Then the hypothesis that contains all tuples  $P'(t, 1)$  for which  $P(t) \in \mathcal{D}$ , and  $P'(t, 0)$  whenever  $P(t) \notin \mathcal{D}$  has the probability  $d^2 > q^2$  over  $\mathcal{P}'$  and satisfies the OMQ  $(Q', \Sigma')$ . Conversely, assume that the hypothesis  $\mathcal{H}$  for  $(Q', \Sigma')$  over  $\mathcal{P}'$  has a probability  $d > q^2$ . Since  $\mathcal{H}$  satisfies  $Q'$  and the new NCs in  $\Sigma'$ , for each tuple  $\langle P(t) : p \rangle \in \mathcal{P}$ , it must contain exactly one of the tuples  $P'(t, 0)$  or  $P'(t, 1)$ , which together contribute a probability of  $(1-p)^2$  or  $p^2$ , respectively, over  $\mathcal{P}'$ . Hence, the sum in Definition 18 collapses to one element, which corresponds exactly to the database  $\mathcal{D}$  obtained by collecting all tuples  $P(t)$  for which  $P'(t, 1)$  is in  $\mathcal{H}$ . This database has a probability of  $\sqrt{d} > q$  over  $\mathcal{P}$ , and must satisfy the original query and NCs.  $\square$

## K Proof of Theorem 21

**Hardness** PP-hardness follows from the complexity of probabilistic query evaluation over PDBs [Suciu *et al.*, 2011] since we can choose  $Q = \top$  and reformulate any UCQ into a set of NCs such that the consistency of a database is equivalent to the non-satisfaction of the UCQ.

**Membership** We consider an OMQ  $(Q, \Sigma)$ , a PDB  $\mathcal{P}$ , and a threshold  $p$ . Since the query is FO-rewritable, it is equivalent to an ordinary UCQ  $Q_\Sigma$  over  $\mathcal{P}$ . Similarly, we can rewrite the UCQ  $Q_\perp$  expressing the non-satisfaction of the NCs into a UCQ  $Q_{\perp, \Sigma}$ . By the observation in Theorem 8, we can enumerate all hypotheses  $\mathcal{H}$ , which are the polynomially many matches for  $Q_\Sigma$  in  $\mathcal{P}$ , and then have to check for each  $\mathcal{H}$  whether the probability of all consistent extensions exceeds  $p$ . The latter part is equivalent to evaluating  $\neg Q_{\perp, \Sigma} \wedge \bigwedge_{t \in \mathcal{H}} t$  over  $\mathcal{P}$ , which can be done by a PP oracle [Suciu *et al.*, 2011]. We accept iff one of these PP checks yields a positive answer. In the terminology of [Beigel *et al.*, 1995], this is a polynomial-time disjunctive reduction of our problem to a PP problem. Since that paper shows that PP is closed under such reductions, we obtain the desired PP upper bound.  $\square$

## L Proof of Theorem 22

**Data complexity for OMQ<sub>GF</sub>** We reduce the following problem from [Wagner, 1986], which uses the *counting quantifier*  $C$ : decide the validity of

$$\Phi = \exists x_1, \dots, x_n C y_1, \dots, y_m \phi,$$

where  $\phi = \phi_1 \wedge \dots \wedge \phi_k$  is a propositional formula in CNF, over the variables  $x_1, \dots, x_n, y_1, \dots, y_m$ . This amounts to checking whether there is a partial assignment for  $x_1, \dots, x_n$  that admits at least  $c$  extensions to  $y_1, \dots, y_m$  that satisfy  $\phi$ .

We can assume without loss of generality that each of the clauses  $\phi_i$  contains exactly three literals: shorter clauses can be padded by copying existing literals, and longer clauses can

be abbreviated using auxiliary variables that are included under the counting quantifier  $C^c$ . Since the values of these variables are uniquely determined by the original variables, this does not change the number of satisfying assignments.

We define the PDB  $\mathcal{P}_\Phi$  that describes the structure of  $\Phi$ :

- For each variable  $v$  occurring in  $\Phi$ ,  $\mathcal{P}_\Phi$  contains the tuples  $\langle V(v, 0) : 0.5 \rangle$  and  $\langle V(v, 1) : 0.5 \rangle$ , where  $v$  is viewed as a constant. These tuples represent the assignments that map  $v$  to *false* and *true*, respectively.
- For each clause  $\phi_j$ , we introduce the tuple  $\langle C(v_1, t_1, v_2, t_2, v_3, t_3) : 1 \rangle$ , where  $t_i$  is 1 if  $v_i$  occurs negatively in  $\phi_j$ , and 0 otherwise (again, all terms are constants). For example, for  $x_3 \vee \neg y_2 \vee x_7$ , we use the tuple  $\langle C(x_3, 0, y_2, 1, x_7, 0) : 1 \rangle$ . This encodes the knowledge about the partial assignments that do not satisfy  $\phi$ .
- We use auxiliary tuples  $\langle A(x_1) : 1 \rangle, \langle S(x_1, x_2) : 1 \rangle, \dots, \langle S(x_{n-1}, x_n) : 1 \rangle$ , and  $\langle L(x_n) : 1 \rangle$  to encode the order on the variables  $x_i$ , and similarly for  $y_j$ :  $\langle B(y_1) : 1 \rangle, \langle S(y_1, y_2) : 1 \rangle, \dots, \langle S(y_{m-1}, y_m) : 1 \rangle$ , and  $\langle L(y_m) : 1 \rangle$ .

We now describe the program  $\Sigma$  used for the reduction. First, we detect whether all variables  $x_i$  ( $1 \leq i \leq n$ ) have a truth assignment (i.e., at least one of the facts  $V(x_i, 0)$  or  $V(x_i, 1)$  is present) by the special nullary predicate  $A$ , using the auxiliary unary predicates  $V$  and  $A$ :

$$\begin{aligned} V(x, t) &\rightarrow V(x) \\ A(x) \wedge V(x) \wedge S(x, x') &\rightarrow A(x'), \\ A(x) \wedge V(x) \wedge L(x) &\rightarrow A, \end{aligned}$$

where  $x, x', t$  are variables. We do the same for the variables  $y_1, \dots, y_m$ :

$$\begin{aligned} B(y) \wedge V(y) \wedge S(y, y') &\rightarrow B(y'), \\ B(y) \wedge V(y) \wedge L(y) &\rightarrow B. \end{aligned}$$

Now, the query  $Q = A$  ensures that only such hypotheses are valid that at least contain a truth assignment for the variables  $x_1, \dots, x_n$ .

Next, we restrict the assignments to satisfy  $\phi$  by using additional NCs in  $\Sigma$ . First, we ensure that there is no “inconsistent” assignment for any variable  $v$ , i.e., only one of the facts  $V(v, 0)$  or  $V(v, 1)$  holds:

$$V(v, 0) \wedge V(v, 1) \rightarrow \perp$$

Furthermore, if all variables  $y_1, \dots, y_m$  have an assignment, then none of the clauses in  $\phi$  can be falsified:

$$\begin{aligned} C(v_1, t_1, v_2, t_2, v_3, t_3) \wedge \\ V(v_1, t_1) \wedge V(v_2, t_2) \wedge V(v_3, t_3) \wedge B \rightarrow \perp, \end{aligned}$$

where  $v_1, t_1, v_2, t_2, v_3, t_3$  are variables.

We show that  $\Phi$  is valid iff there exists a hypothesis  $\mathcal{H}$  that satisfies  $(Q, \Sigma)$  such that all consistent databases that extend  $\mathcal{H}$  sum up to a probability (under  $\mathcal{P}_\Phi$ ) of at least  $p = 0.25^n \cdot 0.25^m (3^m - 2^m + c)$ .

Assume that such a hypothesis  $\mathcal{H}$  exists. Since  $\mathcal{H} \models (Q, \Sigma)$ , we know that for each  $x_i$  ( $1 \leq i \leq n$ ) one of the

tuples  $V(x_i, 0), V(x_i, 1)$  is included in  $\mathcal{H}$ . In each consistent extension of  $\mathcal{H}$ , it must be the case that the complementary facts (representing an inconsistent assignment for  $x_i$ ) are false. In particular, these complementary facts cannot be part of  $\mathcal{H}$  since then its probability would be 0. Hence, we can ignore the factor  $0.25^n$  in the following. There are exactly  $3^m - 2^m$  databases satisfying  $\mathcal{H}$  that represent consistent, but incomplete assignments for the variables  $y_j$ . Since these databases do not entail  $B$ , they are all consistent, and hence counted towards the total sum. The inconsistent assignments for  $y_1, \dots, y_m$  yield inconsistent databases, which leaves us only with the  $2^m$  databases representing proper truth assignments. Those that violate at least one clause of  $\phi$  become inconsistent, and hence there are at least  $c$  such consistent databases iff there are at least  $c$  extensions of the assignment represented by  $\mathcal{H}$  that satisfy  $\phi$ . We conclude these arguments by noting that the probability of each individual choice of tuples  $V(y_j, t_j)$  ( $1 \leq j \leq m$ ) is  $0.25^m$ .

On the other hand, if  $\Phi$  is valid, then we can use the same arguments to construct a hypothesis  $\mathcal{H}$  (representing the assignment for  $x_1, \dots, x_n$ ) that exceeds the given threshold.

**Fp-combined complexity for OMQ $_{\emptyset}$**  As before, we consider a formula

$$\Phi = \exists x_1, \dots, x_n \ C^c y_1, \dots, y_m \ \phi,$$

where  $\phi = \phi_1 \wedge \dots \wedge \phi_k$  is in 3CNF and the PDB  $\mathcal{P}_{\Phi}$  is almost the same as before:

- For each variable  $x_i$ ,  $1 \leq i \leq n$ , we use the tuples  $\langle V(x_i, 0) : 0.5 \rangle$  and  $\langle V(x_i, 1) : 0.5 \rangle$ , and for  $y_j$ ,  $1 \leq j \leq m$ , we use  $\langle V(y_j, 0) : p \rangle$  and  $\langle V(y_j, 1) : p \rangle$ , where  $p$  is a fixed, large probability that we specify later.
- For each clause  $\phi_j$ , we introduce the tuple  $\langle C(v_1, t_1, v_2, t_2, v_3, t_3) : 1 \rangle$  as before.

We use the query

$$Q_{\Phi} = \exists t_1, \dots, t_n \ V(x_1, t_1) \wedge \dots \wedge V(x_n, t_n),$$

which is equivalent to the query  $A$  from above, to enforce that any hypothesis contains a truth assignment for the existentially quantified variables. For the variables  $y_1, \dots, y_m$ , this is handled by a special choice of  $p$ , which ensures that the probabilities of the incomplete assignments sum up to a value that is smaller than the probability of a single complete assignment; hence, we can ignore the incomplete assignments when counting the complete assignments. The program  $\Sigma$  hence only needs to contain the two NCs

$$V(v, 0) \wedge V(v, 1) \rightarrow \perp$$

and

$$C(v_1, t_1, v_2, t_2, v_3, t_3) \wedge \\ V(v_1, t_1) \wedge V(v_2, t_2) \wedge V(v_3, t_3) \rightarrow \perp.$$

To find an appropriate value for  $p$ , consider a fixed hypothesis (which specifies an assignment for the existentially quantified variables), and a single database that contains exactly one of each of the pairs of tuples  $V(y_j, 0), V(y_j, 1)$ , for

each  $y_j$ ,  $1 \leq j \leq m$ . This complete assignment has a probability of  $p^m(1-p)^m$  (if we ignore all other tuples). We now compute the probability mass of all incomplete assignments for  $y_1, \dots, y_m$ . For a fixed number  $k$  of “incomplete variables”, there are  $\binom{m}{k} 2^{m-k}$  such assignments, since we can first choose  $k$  out of  $m$  variables  $y_j$  for which neither tuple  $V(y_j, 0), V(y_j, 1)$  is true, and we have binary choice for each of the remaining  $m-k$  variables. Each such assignment has a probability of  $p^{m-k}(1-p)^{m+k}$ , and hence in total the incomplete assignments have a probability of

$$\begin{aligned} & \sum_{k=1}^m \binom{m}{k} 2^{m-k} p^{m-k} (1-p)^{m+k} \\ &= 2^m p^m (1-p)^m \sum_{k=1}^m \binom{m}{k} \left(\frac{1-p}{2p}\right)^k \\ &= 2^m p^m (1-p)^m \left( \left(1 + \frac{1-p}{2p}\right)^m - 1 \right) \end{aligned}$$

by the binomial theorem. Recall that our goal is to make this number smaller than  $p^m(1-p)^m$ , and hence we need to solve the inequation

$$1 > \left(\frac{p+1}{p}\right)^m - 2^m,$$

which is equivalent to

$$p > \frac{1}{(1+2^m)^{\frac{1}{m}} - 1}.$$

Since the latter term is always smaller than 1, it is possible to choose  $p$  as required, e.g.,  $p = 1 - 2^{-2^m}$ , which has only linearly many digits.

Now we have that  $\Phi$  is valid iff there exists a hypothesis  $\mathcal{H} \models (Q_{\Phi}, \Sigma)$  whose consistent extensions sum up to probability (under  $\mathcal{P}_{\Phi}$ ) of at least  $0.25^m \cdot c \cdot p^m(1-p)^m$ . Indeed, if there exists such an  $\mathcal{H}$ , then it represents a truth assignment for  $x_1, \dots, x_m$ , which accounts for the term  $0.25^n$  (see the proof of Theorem 22). But then to obtain the probability threshold, there must exist at least  $c$  complete assignments for  $y_1, \dots, y_m$ , since the incomplete assignments on their own do not add up to  $p^m(1-p)^m$ . Conversely, if  $\Phi$  is valid, we can use that information to choose a hypothesis that admits at least  $c$  extensions that represent complete truth assignments for  $y_1, \dots, y_m$ .  $\square$

## M Proof of Lemma 23

Recall the proof of Theorem 21. According to [Dalvi and Suciu, 2012], the evaluation problem for the UCQ  $Q_{\perp, \Sigma}$  over  $\mathcal{P}$  is either in P or PP-hard (under Turing reductions). In the former case, MPH can also be decided in deterministic polynomial time. In the latter case, we reduce the evaluation problem for  $Q_{\perp, \Sigma}$  over a PDB  $\mathcal{P}$  to the MPH for  $Q_{\Sigma}$  and  $Q_{\perp, \Sigma}$  over some  $\mathcal{P}' \supseteq \mathcal{P}$ . We introduce an “artificial match” for  $Q_{\Sigma}$  into  $\mathcal{P}'$ , by adding new constants and tuples (with probability 1) that satisfy one disjunct of  $Q_{\Sigma}$ , while taking care that these new tuples do not satisfy  $Q_{\perp, \Sigma}$ . Such tuples must exist if  $Q_{\Sigma}$  is not subsumed by  $Q_{\perp, \Sigma}$ ; otherwise, all hypotheses would trivially have the probability 0 (and hence the MPH would be decidable in polynomial time). In  $\mathcal{P}'$ , the probability of the most probable hypothesis for  $Q_{\Sigma}$  and  $Q_{\perp, \Sigma}$  is the same as the probability of  $Q_{\perp, \Sigma}$  over  $\mathcal{P}$ , and hence deciding the threshold is PP-hard.  $\square$