

# On the Optimal Approximation of Queries Using Tractable Propositional Languages

Robert Fink and Dan Olteanu  
Oxford University Computing Laboratory  
{robert.fink,dan.olteanu}@comlab.ox.ac.uk

## ABSTRACT

This paper investigates the problem of approximating conjunctive queries without self-joins on probabilistic databases by lower and upper bounds that can be computed more efficiently. We study this problem via an indirection: Given a propositional formula  $\Phi$ , find formulas in a more restricted language that are greatest lower bound and least upper bound, respectively, of  $\Phi$ . We study bounds in the languages of read-once formulas, where every variable occurs at most once, and of read-once formulas in disjunctive normal form.

We show equivalences of syntactic and model-theoretic characterisations of optimal bounds for unate formulas, and present algorithms that can enumerate them with polynomial delay. Such bounds can be computed by queries expressed using first-order queries extended with transitive closure and a special choice construct.

Besides probabilistic databases, these results can also benefit the problem of approximate query evaluation in relational databases, since the bounds expressed by queries can be computed in polynomial combined complexity.

## Categories and Subject Descriptors

H.2.4 [Database Management]: Systems—*Query Processing*

## General Terms

Algorithms, Languages, Theory

## Keywords

Probabilistic database, Query approximation

## 1. INTRODUCTION

Approximate query evaluation is preferred over exact query evaluation in cases when the exact evaluation is too expensive and quick approximate answers are favoured over delayed exact answers. This paper investigates the problem of query approximation through the glasses of propositional

formula approximation. Formulas and queries are intimately connected in the following way: Given an input database where every tuple is annotated by a distinct variable, each tuple  $t$  in the query answer is annotated by a formula over the input tuples that contributed to  $t$ .

We seek bounds expressible in restricted classes of propositional formulas that are model-based, optimal, efficiently computable, and allow for efficient probability computation, in case they are defined over random variables. Bounds are defined semantically based on the set of models (or satisfying assignments) of formulas. The set of models of a formula is included in the set of models of any of its upper bounds and includes the set of models of any of its lower bounds. Lower and upper bounds are optimal with respect to a language if there are no other lower and upper bounds, respectively, in that language whose sets of models are closer to the set of models of the formula to approximate.

In this paper, bounds are expressed in the languages of read-once formulas or of its intersection with the language of formulas in disjunctive normal form (DNF). Read-once formulas are propositional formulas where every variable occurs at most once, that is, they are in *1-Occurrence Form*. We denote this class by IOF. The DNF restriction of IOF is denoted here by *independent DNF*, or iDNF for short.

The proposed approximation framework is relevant to database scenarios that can benefit from optimal and efficient approximation of query evaluation. We discuss three such scenarios: approximate probability computation in probabilistic databases, approximation of provenance information and explanations in provenance databases, and approximate query evaluation in relational databases.

## Scenario 1: Probabilistic Databases

Our main motivation for this work comes from probabilistic databases. The problem of query evaluation in probabilistic databases has received tremendous attention in recent years, a fairly recent survey was compiled by Dalvi, Ré, and Suciu [2]. This problem is known to be #P-hard already for simple conjunctive queries and restricted probabilistic databases. One particularly promising evaluation strategy for hard queries is to consider approximate evaluation.

Consider so-called tuple-independent probabilistic databases, where each tuple is associated with a distinct Boolean random variable. By query evaluation techniques reminiscent of conditioned tables [14], propositional formulas called *lineage* can be computed along with the query answers and may subsequently be used to compute answer probabilities. Figure 1 shows the lineage of three queries evaluated on a probabilistic database. The lineages of positive queries are

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

ICDT 2011, March 21–23, 2011, Uppsala, Sweden.

Copyright 2011 ACM 978-1-4503-0529-7/11/0003 ...\$10.00

		S			T	
R	A		A	B		B
		y <sub>1</sub>	1	1	z <sub>1</sub>	1
x <sub>1</sub>	1	y <sub>2</sub>	1	2	z <sub>2</sub>	2
x <sub>2</sub>	2	y <sub>3</sub>	2	1	z <sub>3</sub>	3
x <sub>3</sub>	3	y <sub>4</sub>	2	2	z <sub>4</sub>	4
		y <sub>5</sub>	3	3		
		y <sub>6</sub>	3	4		

Query  $Q_1:-R(A), S(A, B), T(B)$  has lineage

$$\Phi_1 = x_1 y_1 z_1 \vee x_1 y_2 z_2 \vee x_2 y_3 z_1 \vee x_2 y_4 z_2 \vee x_3 y_5 z_3 \vee x_3 y_6 z_4.$$

Query  $Q_2:-R(A), S(A', B), T(B), A \leq A'$  has lineage

$$\Phi_2 = x_1 y_1 z_1 \vee x_1 y_2 z_2 \vee x_1 y_3 z_1 \vee x_1 y_4 z_2 \vee x_1 y_5 z_3 \vee x_1 y_6 z_4 \vee x_2 y_3 z_1 \vee x_2 y_4 z_2 \vee x_2 y_5 z_3 \vee x_2 y_6 z_4 \vee x_3 y_5 z_3 \vee x_3 y_6 z_4.$$

Query  $Q_3:-R(A), S(A, B)$  has lineage

$$\begin{aligned} \Phi_3 &= x_1 y_1 \vee x_1 y_2 \vee x_2 y_3 \vee x_2 y_4 \vee x_3 y_5 \vee x_3 y_6 \\ &= x_1(y_1 \vee y_2) \vee x_2(y_3 \vee y_4) \vee x_3(y_5 \vee y_6). \end{aligned}$$

**Figure 1: Relations  $R, S, T$  annotated with variables from disjoint sets  $V_1 = \{x_1, \dots, x_3\}$ ,  $V_2 = \{y_1, \dots, y_6\}$ ,  $V_3 = \{z_1, \dots, z_4\}$ , respectively. Lineages  $\Phi_1$  to  $\Phi_3$  of four queries  $Q_1$  to  $Q_3$  evaluated on this database.**

unate DNF formulas. A key contribution of this work is an approach that can efficiently derive optimal lower bounds  $\Phi_L$  and upper bounds  $\Phi_U$  in tractable languages such as iDNF or IOF for any unate DNF formula  $\Phi$ . Moreover, we can lift this result to the query level by showing how lineage bounds can be computed by queries.

Let  $\mathcal{M}$  and  $P$  be functions that map formulas to their sets of models and to their probabilities, respectively. Model-based bounds imply probability-based bounds for any probability distribution over the Boolean random variables:

PROPOSITION 1. *Let  $\Phi$ ,  $\Phi_L$ , and  $\Phi_U$  be formulas. Then,*

$$\mathcal{M}(\Phi_L) \subseteq \mathcal{M}(\Phi) \subseteq \mathcal{M}(\Phi_U) \Rightarrow P(\Phi_L) \leq P(\Phi) \leq P(\Phi_U).$$

We devise algorithms that can enumerate optimal bounds efficiently. Such bounds cannot lead, however, to approximations with arbitrary precision, since exact probability computation is  $\#P$ -hard already for unate bipartite DNF formulas. Efficiently computable bounds can be nevertheless very useful in conjunction with known procedures that incrementally decompose the input formula  $\Phi$  into subformulas and compute bounds on  $P(\Phi)$  based on bounds of these subformulas [19, 7]. This is repeated until the desired precision is reached. Efficient computation of good probability bounds for intermediate decomposition steps is essential to the effectiveness of such techniques.

The closest in spirit to our approach is a technique developed independently by Gatterbauer, Abhay, and Suciu [9]. Their technique computes upper bounds for probabilities of conjunctive queries without self-joins. These bounds are not model-based, and in particular not optimal. Although the query lineage is used to compute the upper bounds, the interpretation of this lineage is non-standard: Each literal is interpreted as being a unique variable.

## Scenario 2: Provenance Databases

Our work is applicable to the approximation of provenance information in the answers to queries on databases with annotations. Similarly to the probabilistic case, query lineage encodes symbolically all possible explanations for the existence of a tuple in the query answer in terms of combina-

tions of the input tuples. Lineage bounds can then represent coarser, more compact, and efficiently computable explanations of the query answer. Lower bounds are *correct but possibly incomplete* explanations, while upper bounds are *complete but potentially incorrect* explanations in the following sense: Every explanation for a lower bound is necessarily an explanation for the answer; however, there may exist explanations for the answer that are not explanations for the lower bound. The situation for upper bounds is symmetric.

## Scenario 3: Relational Databases

Existing work on approximate query answering in relational databases considers the use of synopses (histograms, join synopses, and wavelets) to speed-up the evaluation of aggregate queries with the goal of quickly reporting the leading digits of the answers. A survey of these synopsis-based techniques is given by Garofalakis and Gibbons [8]. Their focus is not on deriving optimal bounds within a given language.

Given a query  $Q$  in a language  $\mathcal{QL}$ , the approximation problem we consider in the relational setting is to compute two queries  $Q_L$  and  $Q_U$  in a language  $\mathcal{QL}'$  such that  $Q_L(D) \subseteq Q(D) \subseteq Q_U(D)$  for any relational database  $D$  and the computational complexity for  $\mathcal{QL}'$  is lower than for  $\mathcal{QL}$ .

When applied to relational databases, our approximation approach comes with the following twist. We can derive queries  $Q_L$  and  $Q_U$  such that their answers can be represented compactly as a IOF or iDNF expression over tuples. Such answer representations can be computed in polynomial combined complexity and allow for the enumeration of (possibly exponentially many) answers with polynomial delay.

## Contributions

The main contributions of this paper are as follows:

- We introduce a query approximation framework based on novel results on model-based optimal approximation of unate formulas. We investigate two languages to express lower and upper bounds: read-once formulas and their DNF retractions.
- We give syntactic characterisations of optimal lower and upper bounds in iDNF and IOF, as well as algorithms that can enumerate such bounds for unate DNF formulas with polynomial delay.
- We show that the computation of model-based bounds on the lineage of conjunctive queries without self-joins is expressible using queries that can be evaluated with polynomial combined complexity using iDNF or IOF representations of the query answers.

Due to lack of space, the proofs of the statements made in this paper are given in an extended technical report [6].

## 2. PRELIMINARIES

We use the language of propositional formulas over a finite set of Boolean variables with its standard interpretation.

**Syntax of Formulas.** We assume the standard syntactic notions of variables, constants 1 (*true*) and 0 (*false*), literals, and (propositional) formulas constructed by combining literals or constants using the logical connectives  $\vee$  and  $\wedge$ . In this paper, we consider unate formulas, i.e. formulas in which all literals are positive. A clause is a conjunction of literals. A formula in disjunctive normal form (DNF) is a disjunction of clauses. At the syntactic level, we see clauses

as sets of literals and DNF formulas as sets of clauses. This explains set notations such as  $\subseteq$  on clauses and formulas. Consequently,  $\Phi \subseteq \Psi$  means that  $\Phi$  consists of a subset of the clauses or literals of  $\Psi$ , if  $\Phi$  and  $\Psi$  are DNF formulas or clauses, respectively. The constants 0 and 1 are special sets such that 1 is contained in any set and 0 contains any set. Two clauses are (syntactically) *independent* if their variable sets are disjoint. We use  $\text{vars}(\Phi)$  to denote the set of variables of a formula  $\Phi$ .

*Definition 1.* Let  $(V_1, \dots, V_k)$  be a tuple of disjoint sets of variables. A formula  $\Phi$  is *k-partite* over  $(V_1, \dots, V_k)$  if each clause in  $\Phi$  has the form  $v_1 \wedge \dots \wedge v_k$ , where  $v_i \in V_i$  for all  $1 \leq i \leq k$ . The *projection*  $\Phi_{i,j}$  of  $\Phi$  is a formula obtained from  $\Phi$  by removing all literals not in  $V_i$  or  $V_j$ .  $\square$

*Example 1.* Figure 1 shows 3-partite unate DNF formulas  $\Phi_1$  and  $\Phi_2$  over the tuple of disjoint variable sets  $(V_1, V_2, V_3)$ . Formula  $\Phi_3$  is a 2-partite over  $(V_1, V_2)$ .  $\square$

The graph of a formula  $\Phi$ , denoted by  $G(\Phi)$ , has one node for each distinct variable in  $\Phi$ , and one edge  $(x, y)$  for each pair of literals  $x, y$  that appear together in a clause of  $\Phi$ . For *k-partite* formulas, we use a more specific representation.

*Definition 2.* Let  $\Phi$  be a *k-partite* formula over  $V_1, \dots, V_k$ . The bipartite graph  $B_{i,j} = (V_i, V_j, E_{i,j})$  of projection  $\Phi_{i,j}$  of  $\Phi$  has an edge  $(x, y)$  for each clause  $xy \in \Phi_{i,j}$ . The *set of projection graphs* of  $\Phi$  is  $\mathcal{B}(\Phi) = \bigcup \{B_{i,j} \mid 1 \leq i < j \leq k\}$ .  $\square$

For a given *k-partite* formula with  $n$  variables and  $m$  clauses, its graph has  $n$  nodes and  $O(m \cdot k^2)$  edges. Figure 2 depicts the graphs of the formulas from Figure 1.

**Semantics of Formulas.** We assume the standard semantic notion of satisfying assignments, also called models, of a formula. Let  $\mathcal{M}$  be a function mapping formulas to their sets of models. Given two formulas  $\Phi$  and  $\Psi$  over the same variables,  $\Phi$  is *semantically* contained in  $\Psi$ , denoted by  $\Phi \models \Psi$ , if  $\mathcal{M}(\Phi) \subseteq \mathcal{M}(\Psi)$ . We also say that  $\Phi$  *implies*  $\Psi$ . Equivalence is two-way containment.

**Linking Semantics and Syntax.** Given two clauses  $\varphi$  and  $\psi$ , it holds that  $\varphi \models \psi$  if and only if  $\varphi \supseteq \psi$ . For DNF formulas, we will use the following key result:

LEMMA 1. [22] Let  $\Phi$  and  $\Psi$  be two DNF formulas. Then

$$\Phi \models \Psi \quad \text{iff} \quad \forall \varphi \in \Phi \exists \psi \in \Psi : \varphi \models \psi. \square$$

A unate DNF formula  $\Phi$  is *reducible* if there exist clauses  $\varphi, \psi \in \Phi$  that satisfy  $\varphi \models \psi$ , and *irreducible* otherwise. In the former case,  $\varphi$  is redundant in  $\Phi$  and can be removed.

### Queries on Annotated and Probabilistic Databases.

We consider conjunctive queries without self-joins expressed using Datalog syntax. By default, equi-joins are the only type of joins allowed in such queries. Figure 1 shows three Boolean queries. Query  $Q_1$  has two query variables  $A$  and  $B$  that express joins between  $R, S$ , and  $T$ . Query  $Q_2$  uses an inequality join  $A \leq A'$ .

An annotated database is a relational database, where each tuple is annotated by a propositional formula, sometimes called *lineage*. If all formulas in an annotated database are pairwise independent, then the database is called *tuple-independent*. Such annotated databases have been extensively studied in the context of provenance in databases, incomplete databases [14], probabilistic databases, and in a unifying framework by Green et al. [12].

The evaluation of queries on annotated databases follows standard query evaluation in relational databases, with the addition that the annotations of the input tuples are propagated to the answer tuples they contributed to. The annotation of a distinct tuple created (a) by a join of two input tuples is the conjunction of the annotations of these input tuples, and (b) by a projection or union is the disjunction of the annotations of the input tuples. By construction, the lineage of a conjunctive query without self-joins and with  $k$  relation symbols is thus a unate *k-partite* DNF formula defined over variable sets  $(V_1, \dots, V_k)$ , with one variable set  $V_i$  per relation symbol.

In probabilistic databases, the lineage is defined over random variables and can be used to compute the probabilities of tuples in the answers to queries: Given a tuple  $t$  in the answer to query  $Q$  with lineage  $\Phi$ , the probability of  $t$  is equal to the probability  $P(\Phi)$  that a random valuation of the variables in  $\Phi$  satisfies  $\Phi$ . Computing  $P(\Phi)$  is a #P-hard problem for conjunctive queries, though there are known classes of queries for which the problem is in polynomial time, e.g. [3]. We next recall such a tractable class.

For a query variable  $A$  in a query  $Q$ , we denote by  $s(A)$  the set of relation symbols in  $Q$  that mention  $A$ .

*Definition 3.* A conjunctive query  $Q$  is *hierarchical* if for any two non-head query variables  $A$  and  $B$ , either  $s(A) \cap s(B) = \emptyset$ , or  $s(A) \subseteq s(B)$ , or  $s(A) \supset s(B)$ .  $\square$

In the class of conjunctive queries without self-joins, the hierarchical queries are the only tractable queries [3].

*Example 2.* Consider the queries from Figure 1. Query  $Q_1$  is not hierarchical since  $s(A) = \{R, S\}$ ,  $s(B) = \{S, T\}$ . Query  $Q_3$  is hierarchical, since  $s(B) \subset s(A)$  and  $A$  and  $B$  are the only query variables. Query  $Q_2$  has an inequality join; there are extensions of the hierarchical property to cope with inequality joins [18].  $\square$

**Assumptions and conventions.** Unless explicitly stated otherwise, *we assume formulas to be unate and irreducible* by default. For arbitrary formulas  $\Phi$ , set-based notions like  $\varphi \in \Phi$  or  $x \in \varphi$  are thus with respect to the unique irreducible model-equivalent DNF formula of  $\Phi$ . If a unate formula is given in DNF, computing the irreducible equivalent DNF can be done in polynomial time. Turning formulas into DNF representations has exponential overhead in worst-case.

Lowercase Latin letters denote variable symbols, capital Latin letters denote query variables, and uppercase (lowercase) Greek letters denote formulas (clauses).

## 3. FORMULA APPROXIMATION

The core problem investigated in this paper is: Given a formula  $\Phi$  from a propositional language, find formulas from a more restrictive language that approximate  $\Phi$ .

*Definition 4.* Let  $\mathcal{L}'$  and  $\mathcal{L}$  be languages of propositional formulas such that  $\mathcal{L}' \subset \mathcal{L}$  and  $\Phi \in \mathcal{L}$ . Formulas  $\Phi_L, \Phi_U \in \mathcal{L}'$  are *lower and upper bounds* for  $\Phi$  with respect to  $\mathcal{L}'$ , if  $\mathcal{M}(\Phi_L) \subseteq \mathcal{M}(\Phi)$  and  $\mathcal{M}(\Phi) \subseteq \mathcal{M}(\Phi_U)$  respectively.

If in addition there are no formulas  $\Phi'_L, \Phi'_U \in \mathcal{L}'$  such that  $\mathcal{M}(\Phi_L) \subset \mathcal{M}(\Phi'_L) \subseteq \mathcal{M}(\Phi)$  and  $\mathcal{M}(\Phi) \subseteq \mathcal{M}(\Phi'_U) \subset \mathcal{M}(\Phi_U)$ , then  $\Phi_L$  is a *greatest lower bound (GLB)* and  $\Phi_U$  is a *least upper bound (LUB)* for  $\Phi$  with respect to  $\mathcal{L}'$ .  $\square$

The notions GLB and LUB provide an intuitive semantic characterisation of *optimal* bounds. We consider languages  $\mathcal{L}'$  such that (i) we can efficiently find optimal bounds in  $\mathcal{L}'$ , and (ii)  $\mathcal{L}'$  allows for efficient key computational tasks such as probability computation. The first language  $\mathcal{L}'$ , called *one-occurrence form* (IOF), is that of formulas in which every variable occurs at most once, and of formulas that are equivalent to formulas in which every variable occurs at most once. The second language  $\mathcal{L}'$ , called *independent DNF* (iDNF), is the intersection of IOF and DNF.

*Example 3.* The formula  $x_1y_1 \vee x_1y_2$  is in DNF, but not in iDNF, since the two clauses have the variable  $x_1$  in common. This formula and its equivalent  $x_1(y_1 \vee y_2)$  are in IOF. Formula  $\Phi_3$  from Figure 1 is in IOF.  $\square$

The IOF language has a long history and many names, such as read-once functions [10], fanout-free functions, or non-repeating trees [21], and also many applications including logic synthesis and circuit design [21] and probabilistic databases [17]. Problems such as satisfiability, model counting, and probability computation are hard for general propositional formulas are tractable for IOFs. The IOF property is decidable in polynomial time for DNF formulas. Furthermore, the equivalent (syntactic) IOF can be found in polynomial time and is unique up to commutativity of the binary connectives [21]. The IOF language is particularly relevant in the context of probabilistic databases, since the lineage of any tractable relational algebra query without repeating symbols on tuple-independent probabilistic databases is in IOF and its probability can be computed using relational query plans [17, 7].

Languages beyond IOF lose tractability. It is NP-hard to decide if a given unate formula admits an equivalent formula where every variable occurs at most twice (read-twice formulas) [5]; this is still open in case of input DNF formulas. In case of formulas where every variable occurs at most four times, model counting (and hence probability computation) is already #P-hard [26]. A different direction is to consider complete languages, i.e. languages that can represent any propositional formula either directly or in a formalism that preserves its models, and that still allow for efficient probability computation. A prime example of such languages is the family of Binary Decision Diagrams, including OBDDs, FBDDs, and d-DNNFs [4]. Such languages allow for several equivalent representations of the same input formula but with an exponential gap between their sizes. Finding a minimal representation for an input formula is NP-hard [16].

*Remark 1.* There is a solid body of work on approximating various computational tasks on propositional formulas. Exact and approximate approaches to model counting, which are also relevant to probabilistic databases, are surveyed by Gomes et al. [11]. These approaches however do not consider the approximation of propositional formulas by formulas in a more restrictive language, as done here. The closest in spirit to our work is that by Selman on approximating CNF theories by model-based lower and upper bound conjunctions of Horn clauses [23]. The languages iDNF and IOF are incomparable to conjunctions of Horn clauses.  $\square$

Before we discuss optimal iDNF and IOF approximations, we give a syntactic characterisation of (not necessarily optimal) lower and upper bound formulas for arbitrary unate formulas: *All* bounds for a given unate formula  $\Phi$  can be defined by simple *syntactic* manipulations of  $\Phi$ .

PROPOSITION 2. Let  $\Phi$ ,  $\Phi_L$ , and  $\Phi_U$  be unate formulas.

- $\Phi_L$  is a lower bound of  $\Phi$  if and only if  $\Phi_L$  can be obtained by removing clauses from  $\Phi$  or by adding literals to its clauses. Removing all clauses results in the lower bound 0.
- $\Phi_U$  is an upper bound of  $\Phi$  if and only if  $\Phi_U$  can be obtained by adding clauses to  $\Phi$  or by removing literals from its clauses. Removing all literals from a clause results in the upper bound 1.  $\square$

These syntactic manipulations necessarily lead to lower and upper bounds. Indeed, by removing clauses from a formula we reduce its set of models. By adding variables to clauses, we further constrain their satisfiability and hence reduce the set of models. In both cases, one obtains formulas that are lower bounds. The inverse manipulations lead to upper bounds. More importantly, Proposition 2 states that all bounds can be gained in this way, a fact that follows from Lemma 1. Section 4 shows that it is not necessary to add literals to clauses in order to find *optimal* iDNF lower bounds, while removing variables can be required in order to get to optimal iDNF upper bounds. Section 5 shows how to obtain optimal IOF bounds that can be constructed by solely removing or adding clauses.

*Example 4.* Neither formula  $\Phi_1$  nor  $\Phi_2$  from Figure 1 are in iDNF. The following iDNF formulas are lower and upper bounds for  $\Phi_1$  (several others are possible):

$$\Phi_L = x_1y_1z_1 \vee x_2y_4z_2 \vee x_3y_5z_3, \quad \Phi_U = z_1 \vee x_1y_2 \vee x_2y_4z_2 \vee x_3.$$

Indeed,  $\Phi_L$  is a lower bound since it is a subset of  $\Phi$ . Each clause in  $\Phi$  implies at least one clause of  $\Phi_U$ , hence  $\Phi_U$  is an upper bound. A closer inspection reveals that each clause in  $\Phi_U$  can be obtained by dropping literals from clauses in  $\Phi$ . As we will show in Section 4, the bounds  $\Phi_L$  and  $\Phi_U$  are already optimal for iDNF. They are, however, not optimal for IOF, since they can be improved as follows:

$$\Phi'_L = x_1y_1z_1 \vee x_2y_4z_2 \vee x_3(y_5z_3 \vee y_6z_4).$$

$$\Phi'_U = (y_1 \vee y_3)z_1 \vee (x_1y_2 \vee x_2y_4)z_2 \vee x_3(y_5z_3 \vee y_6z_4). \square$$

## 4. OPTIMAL IDNF APPROXIMATIONS

We characterise syntactically the optimal iDNF bounds for unate formulas, and give algorithms that find such bounds.

### 4.1 iDNF Greatest Lower Bounds

Our main tool used to find optimal iDNF lower bounds is the syntactic notion of maximal lower bounds.

*Definition 5.* Let  $\Phi$  be a unate formula. An iDNF formula  $\Phi_L$  is called a *maximal lower bound (MLB)* for  $\Phi$  if it satisfies the following conditions:

1. (*Lower bound*)  $\Phi_L$  contains a subset of the clauses of  $\Phi$ :  $\Phi_L \subseteq \Phi$ ;
2. (*Maximality*)  $\Phi_L$  cannot be further extended: There is no clause  $\varphi \in \Phi$  such that  $\text{vars}(\varphi) \cap \text{vars}(\Phi_L) = \emptyset$ .  $\square$

The first criterion ensures that  $\Phi_L$  is indeed an iDNF lower bound. The maximality property enforces that it is a *maximal* lower bound. An MLB for a formula  $\Phi$  is thus a maximal set of pairwise independent clauses of  $\Phi$ .

*Example 5.* The MLBs for formula  $\Phi_1$  from Figure 1 are:

$$x_1y_1z_1 \vee x_2y_4z_2 \vee x_3y_5z_3, \quad x_1y_1z_1 \vee x_2y_4z_2 \vee x_3y_6z_4 \\ x_1y_2z_2 \vee x_2y_3z_1 \vee x_3y_5z_3, \quad x_1y_2z_2 \vee x_2y_3z_1 \vee x_3y_6z_4 \quad \square$$

With respect to the iDNF language, the maximal lower bounds correspond precisely to the greatest lower bounds. This is particularly important, since MLB is a syntactic notion defined in terms of sets of clauses, and GLB is a semantic notion defined in terms of sets of models.

**THEOREM 1.** *Let  $\Phi$  be a unate formula. An iDNF formula  $\Phi_L$  is a maximal lower bound for  $\Phi$  if and only if  $\Phi_L$  is a greatest lower bound for  $\Phi$ .*

The implication  $\text{GLB} \Rightarrow \text{MLB}$  holds since if any of the two MLB conditions fail, then we cannot obtain a GLB. For the other direction, it can be shown that no strict upper bound  $\Phi'_L$  for  $\Phi_L$  is a lower bound for  $\Phi$ . The latter case has two sub-cases: There is a clause  $\varphi \in \Phi$  that is in  $\Phi'_L$  and either is or is not in  $\Phi_L$ . In the first sub-case,  $\varphi$  necessarily shares variables with one clause in  $\varphi_L \in \Phi_L$ , and, according to Lemma 1, it must be contained in  $\varphi_L$ . But then,  $\varphi$  must be  $\varphi_L$ , otherwise  $\Phi_L$  can be further extended and hence is no MLB. The second sub-case follows similarly.

The syntactic definition of MLBs suggests an algorithm for enumerating all GLBs by recursively constructing all maximal subsets of non-conflicting clauses of  $\Phi$ . This algorithm may need time exponential in the number of clauses of  $\Phi$ . The following proposition shows that this is necessarily so:

**PROPOSITION 3.** *The unate DNF formula*

$$\Phi = (x_1y_1 \vee x_1y_2) \vee \dots \vee (x_ny_{2n-1} \vee x_ny_{2n})$$

*has  $2n$  clauses and  $2^n$  iDNF GLBs. Any formula containing either  $x_iy_{2i-1}$  or  $x_iy_{2i}$  for all  $i \leq n$  is an iDNF GLB of  $\Phi$ .  $\square$*

Since there can be exponentially many GLBs that are by definition incomparable with respect to their model sets, it may be desirable to find “the best” GLB for a formula according to different criteria. Possible rankings of GLBs are on the number of clauses or, an arguably more useful criterion, on their probabilities, provided they are defined over random variables. Additionally, it can be useful to have an enumeration of all GLBs with polynomial delay [15], which means that the time before finding the first GLB, as well as the time between every two consecutive GLBs, is polynomial only in the input size, and not in the number of GLBs. We obtain results for ranking and enumerating GLBs by exploiting the following correspondence between MLBs and independent sets in graphs.

**Definition 6.** Let  $\Phi$  be a DNF formula. The *clause-dependency graph* of  $\Phi$  is a graph  $G = (\Phi, E)$ , where nodes are clauses of  $\Phi$ , and edges represent pairs of dependent clauses:

$$E = \{(\varphi, \psi) \mid \varphi, \psi \in \Phi \text{ and } \text{vars}(\varphi) \cap \text{vars}(\psi) \neq \emptyset\}. \square$$

**LEMMA 2.** *Let  $\Phi$  be a unate DNF formula. A formula  $\Phi_L \subseteq \Phi$  is an iDNF greatest lower bound for  $\Phi$  if and only if the clause-dependency graph of  $\Phi_L$  is a maximal independent set in the clause-dependency graph of  $\Phi$ .*

A similar equivalence can be established between maximum weighted independent sets and iDNF greatest lower bounds with maximum probability. Therefore, known results for independent sets, such as enumeration of maximal independent sets with polynomial delay [25, 15], immediately carry over to the case of iDNF greatest lower bounds.

**MUB (Unate formula  $\Phi$ , iDNF upper bound  $\Phi_U$ )**  
**outputs all iDNF MUBs of  $\Phi$**

$\Phi_0 \leftarrow \{\varphi_u \in \Phi_U : \varphi_u \text{ has no critical witness in } \Phi\}$   
mergeCandidates  $\leftarrow \{(\varphi_u, w, \varphi'_u) \mid w \in \Phi; \varphi_u, \varphi'_u \in \Phi_0;$   
 $w \text{ is witness of } \varphi_u \text{ and } \varphi'_u\}$

**if  $\Phi_0 \neq \emptyset$  then**

    // Remove clauses without crit. witness  
    **foreach**  $\varphi_u \in \Phi_0$  **do** MUB( $\Phi, \Phi_U \setminus \{\varphi_u\}$ )  
    // Merge clauses without critical witness  
    **foreach**  $(\varphi_u, w, \varphi'_u) \in \text{mergeCandidates}$  **do**  
        MUB( $\Phi, (\Phi_U \setminus \{\varphi_u, \varphi'_u\}) \cup \{\{\varphi_u \cup \varphi'_u\}\}$ )

**else**

$V \leftarrow \text{vars}(\Phi) \setminus \text{vars}(\Phi_U)$   
    Extend  $\Phi_U$  by variables in  $V$  in every possible way  
    without altering the witness relations  
    **Output**  $\Phi_U$

**Algorithm 1: Finding all MUBs of a unate formula  $\Phi$ . Initial call MUB( $\Phi, \bigcup_{v \in \text{vars}(\Phi)} \{v\}$ ).**

**COROLLARY 1** (LEMMA 2, [25, 15]). *The set of all iDNF greatest lower bounds for a unate DNF formula can be enumerated with polynomial delay.  $\square$*

Since finding a maximum independent set is NP-hard, we cannot efficiently enumerate the maximal independent sets, or equivalently the iDNF GLBs, in decreasing order of their size or probabilities.

**PROPOSITION 4.** *Enumerating the iDNF greatest lower bounds of a unate DNF formula in decreasing order of the number of their clauses or their probabilities is NP-hard.  $\square$*

Although we cannot efficiently obtain the iDNF GLB with maximum probability, a constant-factor approximation can be obtained for input unate  $k$ -partite DNF formulas.

**PROPOSITION 5.** *Let  $\Phi$  be a  $k$ -partite unate DNF formula. There exists a polynomial time algorithm that constructs an iDNF greatest lower bound  $\Phi_L$  for  $\Phi$  such that  $P(\Phi_L^{\text{opt}}) \leq k \cdot P(\Phi_L)$ , where  $\Phi_L^{\text{opt}}$  is the iDNF greatest lower bound for  $\Phi$  with the highest probability amongst all of  $\Phi$ 's iDNF greatest lower bounds.  $\square$*

The algorithm constructs  $\Phi_L$  by iterating over  $\Phi$ 's clauses in decreasing order of their probabilities and greedily selecting a maximal set of pairwise independent clauses.  $\Phi_L$  thus contains the first clause  $\varphi_1$  in the order, then the next clause independent of  $\varphi_1$ , and so on. In the construction of an upper bound for  $P(\Phi_L^{\text{opt}})$ , we use the fact that the maximum number of pairwise independent clauses between two clauses of  $\Phi_L$  in the descending order of clauses is smaller than  $k+1$ .

## 4.2 iDNF Least Upper Bounds

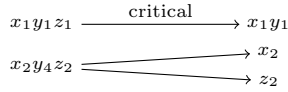
Similar to the case of iDNF greatest lower bounds, we give a syntactic characterisation of iDNF least upper bounds for unate DNF formulas. We introduce a few necessary notions.

**Definition 7.** Let  $\Phi$  and  $\Psi$  be formulas. A clause  $\varphi \in \Phi$  is a *witness* for a clause  $\psi \in \Psi$  if  $\varphi \models \psi$ . We also say that  $\psi$  has the witness  $\varphi$ . If in addition there is no clause  $\psi' \in \Psi$  with  $\psi \neq \psi'$  and  $\varphi \models \psi'$ , then  $\varphi$  is a *critical witness* for  $\psi$ . We then also say that  $\psi$  has the *critical witness*  $\varphi$ .  $\square$

*Example 6.* Consider formula  $\Phi_1$  from Figure 1 and the formula  $\Psi = x_1y_1 \vee x_2 \vee z_2 \vee x_1y_6$ . Clause  $x_1y_1 \in \Psi$  has the critical witness  $x_1y_1z_1 \in \Phi_1$ , clause  $x_2y_4z_2$  is a non-critical witness for  $x_2$  and  $z_2 \in \Psi$ , and finally clause  $x_1y_6 \in \Psi$  has no witness in  $\Phi_1$ .  $\square$

*Definition 8.* Let  $\Phi$  and  $\Psi$  be formulas. The bipartite witness graph  $\Gamma_{\Phi, \Psi} = (V, V', E)$  from  $\Phi$  to  $\Psi$  is defined as follows:  $V$  is the set of clauses of  $\Phi$ ,  $V'$  is the set of clauses of  $\Psi$ , and there is a directed edge  $(\varphi, \psi) \in E$  from a clause  $\varphi \in \Phi$  to a clause  $\psi \in \Psi$  if  $\varphi \models \psi$ .  $\square$

A witness graph  $\Gamma_{\Phi, \Psi}$  shows which clauses of  $\Phi$  imply which clauses of  $\Psi$ . The set of clauses with incoming edges thus represent an upper bound for the set of clauses with outgoing edges. A node  $\varphi \in V$  is then a witness for all nodes  $\psi \in V'$  with  $(\varphi, \psi) \in E$ , and it is a critical witness for a node  $\psi \in V'$  if  $(\varphi, \psi) \in E$  and there is no other node in  $V$  that is a witness for  $\psi$ . The following figure shows the connected nodes in the witness graph for formulas  $\Phi_1, \Psi$  from Example 6.



We are now ready to present our main syntactic tool for finding upper bounds.

*Definition 9.* Let  $\Phi$  be a unate formula. An iDNF formula  $\Phi_U$  is a *minimal upper bound (MUB)* for  $\Phi$  if it satisfies the following conditions:

1. (*Upper bound*) Every  $\varphi \in \Phi$  is a witness for some clauses in  $\Phi_U$ ;
2. (*Maximality*) There is no clause  $\varphi_u \in \Phi_U$  that can be extended by a variable from  $vars(\Phi)$  while preserving condition (1) and keeping  $\Phi_U$  in iDNF;
3. (*Irreducibility*) Every clause  $\varphi_u \in \Phi_U$  has some critical witnesses in  $\Phi$ .  $\square$

Following Lemma 1, the first condition ensures that  $\Phi \models \Phi_U$ . The second condition ensures that we cannot obtain an iDNF refinement of  $\Phi_U$  which is still an upper bound for  $\Phi$  by making clauses in  $\Phi_U$  more specific (thus narrowing the set of models of  $\Phi_U$ ). The third condition states that all clauses in  $\Phi_U$  are necessary, and dropping any of them would lead to clauses in  $\Phi$  violating Lemma 1.

This syntactic characterisation of upper bounds precisely matches the semantic notion of least upper bounds:

**THEOREM 2.** *Let  $\Phi$  be a unate formula. An iDNF formula  $\Phi_U$  is a minimal upper bound for  $\Phi$  if and only if  $\Phi_U$  is a least upper bound for  $\Phi$ .*  $\square$

The implication LUB  $\Rightarrow$  MUB holds since, if any of the three MUB conditions fail, we cannot obtain an LUB. The other direction uses the fact that any difference between an MUB  $\Phi_U$  for  $\Phi$  and a potentially better upper bound  $\Phi'_U$  consists of having a clause  $\varphi'_u \in \Phi'_U$  strictly containing (syntactically) a clause  $\varphi \in \Phi_U$ . A case analysis on a variable  $x \in \varphi'_u$  and  $x \notin \varphi_u$  shows that  $\Phi_U$  does not satisfy the syntactic characterisation of an MUB: Either the irreducibility condition fails, or a clause in  $\Phi_U$  may be expanded by  $x$  to obtain a smaller bound.

PolyMUB (Unate formula  $\Phi$ , iDNF upper bound  $\Phi_U$ )

**outputs** iDNF MUBs of  $\Phi$

$\Phi_0 \leftarrow \{\varphi_u \in \Phi_U : \varphi_u \text{ has no critical witness in } \Phi\}$

$mergeCandidates \leftarrow \{(\varphi_u, w, \varphi'_u) \mid w \in \Phi; \varphi_u, \varphi'_u \in \Phi_0; w \text{ is witness of } \varphi_u \text{ and } \varphi'_u\}$

**if**  $\Phi_0 \neq \emptyset$  **then**

**foreach**  $(\varphi_u, w, \varphi'_u) \in mergeCandidates$  **do**

**foreach**  $v \in \varphi_u$  **do**

$A_{v,v'} \leftarrow 1; A_{v',v} \leftarrow 1$

**if**  $\forall v \in vars(\varphi_u), \forall v' \in vars(\varphi'_u) : A_{v,v'} = 0$

**then**  $MUB(\Phi, (\Phi_U \setminus \{\varphi_u, \varphi'_u\}) \cup \{\varphi_u \cup \varphi'_u\})$

**else**

$\perp$  **Output**  $\Phi_U$

**Algorithm 2: Enumerating a subset of all MUBs of a formula  $\Phi$  with polynomial delay.** Initial call **PolyMUB**( $\Phi, \bigcup_{v \in vars(\Phi)} \{\{v\}\}$ ). The algorithm assumes the existence of a *global* variable  $A$  such that for each pair of variables  $v, v' \in vars(\Phi)$ ,  $A_{v,v'}$  is a boolean flag. Initially,  $A_{v,v'} = 0$  for all  $v, v'$ .

Algorithm 1 constructs MUBs for a formula  $\Phi$  by starting with a pessimistic upper bound  $\Phi_U = \bigcup_{v \in vars(\Phi)} \{\{v\}\}$  and refining it recursively: In every recursion step, witness graphs are constructed for the current upper bound  $\Phi_U$  and each tighter bound representing a refinement of it obtained by merging or removing clauses in  $\Phi_U$ . The set *mergeCandidates* contains pairs of clauses of  $\Phi_U$  together with their common witness that may be merged without violating the *upper bound* condition. Every bound produced by the algorithm is indeed an MUB: The *upper bound* criterion is trivially satisfied by the initial pessimistic bound and remains satisfied, because only non-critical clauses are merged or removed. The base case of the recursion is reached whenever all clauses in  $\Phi_U$  have at least one critical witness (*irreducibility*), i.e.,  $\Phi_0 = \emptyset$  in the algorithm. Every bound is then expanded to satisfy the *maximality* property.

*Example 7.* The formula  $\Phi_{1,u}^1 = x_1y_1 \vee y_2 \vee x_2y_3z_1 \vee y_4z_2 \vee x_3y_5z_3 \vee y_6z_4$  is an MUB for  $\Phi_1$  from Figure 1. Every clause in  $\Phi_{1,u}^1$  has exactly one witness in  $\Phi_1$  which is also critical. No clause can be extended further, since  $vars(\Phi_{1,u}^1) = vars(\Phi_1)$ . The formula  $\Phi_{1,u}^2 = x_1 \vee x_2 \vee x_3y_5z_3 \vee y_6z_4$  is also an MUB; here clauses  $x_1$  and  $x_2$  both have two critical witnesses. It can be checked that  $\Phi_{1,u}^2$  cannot be extended by any variable from  $vars(\Phi_1)$ .  $\square$

As in the case of greatest lower bounds, we can enumerate all iDNF least upper bounds of a given unate formula  $\Phi$ :

**PROPOSITION 6.** *Algorithm 1 computes all iDNF least upper bounds of a unate formula.*  $\square$

Soundness is sketched above. For completeness, consider an arbitrary MUB  $\Phi_U$  and for illustration purposes, imagine Algorithm 1 to be non-deterministic, i.e. instead of iterating over possibilities, it chooses one randomly. It is straightforward to see that there exists a computation branch that produces  $\Phi_U$ : First, remove all variables that are not in  $\Phi_U$ , then merge clauses. These merging and removing steps necessarily have to occur intertwined. Consider formula  $ax \vee bx \vee ay \vee by$ . The LUB  $a \vee b$  cannot be obtained if

1OF-LUB ( $k$ -partite unate formula  $\Phi$ )

**outputs** 1OF LUBs of  $\Phi$  that are  $k$ -partite

$\mathcal{B} \leftarrow$  set of projections graphs.

**Complete** connected components of each projection graph in  $\mathcal{B}$ .

$\mathcal{M} \leftarrow$  subset of all projection graphs  $\mathcal{B}$  that contain a component that is not aligned with components of a different projection graph in  $\mathcal{B}$ .

**if**  $\mathcal{M} = \emptyset$  **then Output** the 1OF represented by  $\mathcal{B}$ .

**foreach** projection graph  $B \in \mathcal{M}$  **do**

$\mathcal{B}' \leftarrow$  copy of  $\mathcal{B}$ .

**while** there exist two projection graphs in  $\mathcal{B}'$  that have non-aligned components  $C, C'$  **do**

        // Wlog assume that the

        // projection graph of  $C$  is not  $B$ .

**Align**  $C$  to  $C'$  in  $\mathcal{B}'$ .

**Complete** connected components of the projection graph of  $C$  in  $\mathcal{B}'$ .

**Output** the 1OF represented by  $\mathcal{B}'$ .

**Complete** projection graph  $B$  in  $\mathcal{B}'$

**foreach** Connected component  $C$  in  $B$  **do**

    Replace component  $C$  with  $(V, W, V \times W)$  in  $\mathcal{B}'$ .

**Align** component  $C$  to component  $C'$  in  $\mathcal{B}'$

**if**  $C$  and  $C'$  are not aligned **then**

    A variable set  $V$  of  $C = (V, W, E)$  is incomparable with a variable set  $V'$  of  $C' = (V', W', E')$ , i.e.  $V \cap V' \neq \emptyset$  and  $V \not\subseteq V'$  and  $V \not\supseteq V'$ .  
     $D \leftarrow (V - V') \times W'$ .

    Add edges  $D$  to projection graph of  $C'$  in  $\mathcal{B}'$ .

**Algorithm 3: Finding 1OF LUBs of a  $k$ -partite unate formula  $\Phi$ .**

clauses are merged first, and the LUB  $ax \vee b \vee y$  cannot be obtained if clauses are removed first.

The intertwined occurrence of merging, removing, and adding clauses makes the design of a polynomial-delay algorithm for enumerating MUBs difficult, since different computation branches may not necessarily lead to different MUBs.

*Example 8.* Consider the formula  $ab \vee ax \vee xy$ . Starting with the usual pessimistic bound  $a \vee b \vee x \vee y$ , the MUB  $ab \vee x$  is found in a computation branch that first merges  $a$  and  $b$  and then removes  $y$ . A second branch might first remove  $a$  and  $y$ , and then expand clause  $b$  by  $a$ . A third branch will delete  $y$  and then merge  $a$  and  $b$ . All these branches lead to the same MUB  $ab \vee x$ .  $\square$

Polynomial delay can be obtained, however, if we only enumerate MUBs that have all input variables:

**PROPOSITION 7.** *Given a unate DNF formula  $\Phi$ , the set of all iDNF least upper bounds  $\Phi_U$  for  $\Phi$  such that  $\text{vars}(\Phi) = \text{vars}(\Phi_U)$  can be enumerated with polynomial delay.*  $\square$

Algorithm 2 can be used to this effect. Soundness and completeness with respect to this restricted class of MUBs follow equivalently to Algorithm 1. The generated MUBs contain all variables of the input formula, since no variable ever gets dropped from the initial upper bound that contains all variables. It remains to be shown that the algorithm never returns two identical MUBs due to alternative

merging orders that may lead to the same MUB and that the variables cannot be first deleted and later added again, as it can happen in Algorithm 1. The symmetric matrix  $A$  records the history of clause merges:  $A_{v,v'} = 1$  if and only if an MUB has already been discovered in which variables  $v, v'$  appear in the same clause. The condition  $A_{v,v'} = 0$  before merging clauses thus ensures that any new MUB does not have a clause with  $v$  and  $v'$ . Furthermore, it is guaranteed that every new MUB contains a clause in which  $v$  or  $v'$  are paired with some other variable. Since all MUBs are iDNF formulas, they are thus pairwise incomparable. Since all computation branches are bounded in depth by  $|\text{vars}(\Phi)|$ , return distinct MUBs, and are independent and may thus be removed from the stack after returning, and since matrix  $A$  has size  $|\text{vars}(\Phi)|^2$ , the algorithm needs space polynomial in the size of the input formula.

*Example 9.* The MUB  $\Phi_{1,u}^2$  from Example 7 has fewer variables than  $\Phi_1$  and can thus not be found by Algorithm 2. Consider the iDNF formula  $abc \vee xyz$  which is also its only MUB. The algorithm indeed finds this bound, for example by merging steps  $ab \vee c \vee x \vee y \vee z$ ,  $abc \vee x \vee y \vee z$ ,  $abc \vee xy \vee z$ ,  $abc \vee xyz$  and sets  $A_{a,b} = A_{b,c} = A_{a,c} = A_{x,y} = A_{y,z} = A_{x,z} = 1$ . Consequently, upon returning from the recursion to the top-level invocation of PolyMUB with  $\Phi_U = a \vee b \vee c \vee x \vee y \vee z$ , the for-loop will not invoke another recursive PolyMUB call, e.g., for  $a \vee bc \vee x \vee y \vee z$ , since  $A_{b,c} = 1$ . There is thus exactly one computation branch for one MUB.  $\square$

The time complexity of Algorithms 1 and 2 can be exponential in the size of the input, and this is necessarily so.

**PROPOSITION 8.** *The unate DNF formula*

$$\Phi = (x_1y_1 \vee x_1y_2) \vee \dots \vee (x_ny_{2n-1} \vee x_ny_{2n})$$

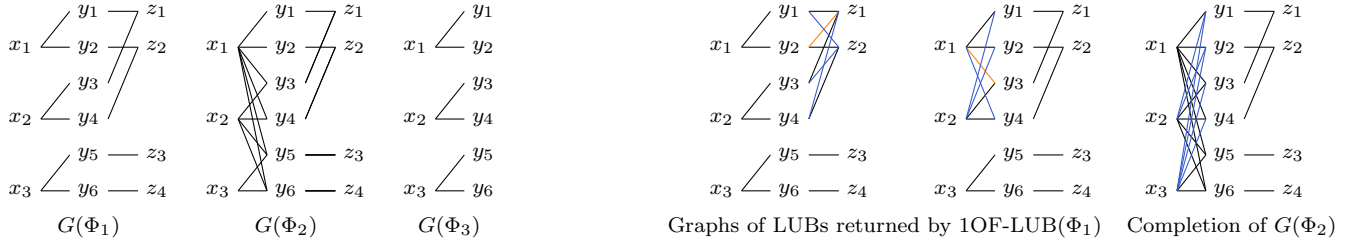
*has  $2n$  clauses and  $3^n$  iDNF least upper bounds. Any formula consisting of either  $x_i$ , or  $x_iy_{2i-1} \vee y_{2i}$ , or  $y_{2i-1} \vee x_iy_{2i}$  for all  $1 \leq i \leq n$  is an iDNF least upper bound of  $\Phi$ .*  $\square$

## 5. OPTIMAL 1OF APPROXIMATIONS

In this section we show how to derive optimal 1OF bounds for unate  $k$ -partite formulas representing lineage of conjunctive queries without self-joins.

Recognizing 1OF formulas can be done using the following characterisation [13]: A unate formula  $\Phi$  is in 1OF if and only if  $\Phi$  is *normal* and  $G(\Phi)$  is  $P_4$ -free. A graph is  $P_4$ -free if it has no induced subgraph that is isomorphic to  $P_4$ , i.e. the graph  $\bullet - \bullet - \bullet - \bullet$ . A formula is normal if it can be unambiguously recovered from its graph. For instance, in the graphs of Figure 2 (with edges  $(x_i, z_j)$  added),  $G(\Phi_2)$  has the  $P_4$  subgraph  $x_1 - z_1 - x_2 - z_3$ , and  $G(\Phi_1)$  has the  $P_4$  subgraph  $x_1 - z_1 - x_2 - y_4$ . The formulas  $\Phi_1$  and  $\Phi_2$  are thus not in 1OF. For unate DNF formulas,  $P_4$ -freeness and subsequent normality checks can be performed in polynomial time in their size [10]. The graph of a  $k$ -partite formula is always normal, if it is  $P_4$ -free [24].

A 1OF formula can be recovered from its graph  $G$  in polynomial time [13, 10]. We next sketch the implementation of a function  $\rho$  that maps graphs to 1OF formulas. If  $G$  is a set of connected components  $G_1, \dots, G_n$ , then  $\rho(G) = (\rho(G_1) \vee \dots \vee \rho(G_n))$ . Otherwise, we compute the complement of  $G$  and check again for connected components  $G_1, \dots, G_m$ . Then,  $\rho(G) = (\rho(G_1) \wedge \dots \wedge \rho(G_m))$ . We recurse with smaller graphs and check for connected components in



**Figure 2: Graphs of  $\Phi_1, \Phi_2, \Phi_3$  from Figure 1. Edges  $(x_i, z_j)$  are not shown to avoid clutter, but can be inferred here since  $x_i$  is connected to  $y_i$ 's that are connected to  $z_j$ 's. The fourth and fifth graphs represent LUBs returned by Algorithm 3; orange and blue edges are added due to the alignment and completion steps, respectively. The sixth graph depicts the completion of  $G(\Phi_2)$ .**

these graphs and their complements. Before any complementation, the formula is obtained as the disjunction of the subformulas for the connected components. With each complementation, we alternate conjunctions and disjunctions. In case of a node, we return its literal. If  $G$  represents an IOF formula, then this procedure necessarily finds it.

*Example 10.* Consider the graph  $G$  that is the third from left in Figure 2. It has three components:  $C_1, C_2, C_3$ . Then,  $\rho(G) = \rho(C_1) \vee \rho(C_2) \vee \rho(C_3)$ . We next illustrate for  $C_1$ , the same applies to the other two components. We complement it, and obtain two components:  $C_{11}$  consisting of one node  $x_1$  and  $C_{12}$  consisting of two connected nodes  $y_1$  and  $y_2$ . Then,  $\rho(C_1) = x_1 \wedge \rho(C_{12})$ . We complement  $C_{12}$  and obtain two components, one with the node  $y_1$  and the other with the node  $y_2$ . Hence,  $\rho(C_{12}) = y_1 \vee y_2$ . The IOF formula represented by  $G(\Phi_3)$  is  $x_1(y_1 \vee y_2) \vee x_2(y_3 \vee y_4) \vee x_3(y_5 \vee y_6)$ .

For the rightmost graph in Figure 2, we complement it and obtain two connected components: One component consists of a clique of three nodes  $x_1, x_2, x_3$ , and the other consists of all the remaining nodes. The obtained formula is thus a conjunction of the formulas of these two components. After several decomposition steps, we obtain

$$(x_1 \vee x_2 \vee x_3)[z_1(y_1 \vee y_3) \vee z_2(y_2 \vee y_4) \vee z_3y_5 \vee z_4y_6]. \square$$

We next give an alternative graph-based characterisation of  $k$ -partite IOF formulas. This is used in the next sections to create optimal IOF bounds and is based on two graph notions: alignment and completeness. A bipartite graph  $(X, Y, E)$  is complete if  $E = X \times Y$ .

*Definition 10.* Two bipartite graphs  $G = (V_1, V_2, E)$  and  $G' = (V'_1, V'_2, E')$  are *aligned* if for all pairs  $(V_i, V'_j)$ ,  $1 \leq i, j \leq 2$ , it holds that  $V_i \cap V'_j = \emptyset$  or  $V_i \subseteq V'_j$  or  $V'_j \subseteq V_i$ .  $\square$

Our characterisation of  $k$ -partite IOFs is as follows:

**LEMMA 3.** *Let  $\mathcal{B}$  be the set of projection graphs of a unate  $k$ -partite formula. The set of connected components of the bipartite graphs in  $\mathcal{B}$  are complete and pairwise aligned if and only if the formula represented by  $\mathcal{B}$  is in IOF.*  $\square$

Since a  $k$ -partite formula  $\Phi$  is in general not normal, its graph may represent a different formula. The next lemma shows that the represented formula is a  $k$ -partite upper bound of  $\Phi$  and that there is no smaller normal  $k$ -partite upper bound. This implies that the graph only represents clauses that are necessarily in every  $k$ -partite IOF upper bound.

**LEMMA 4.** *Let  $\Phi$  be a  $k$ -partite formula,  $\mathcal{B}$  its set of projection graphs, and  $\Phi_u$  the formula represented by  $\mathcal{B}$ . Then*

$\Phi_u$  is  $k$ -partite,  $\mathcal{M}(\Phi) \subseteq \mathcal{M}(\Phi_u)$ , and there is no normal  $k$ -partite formula  $\Phi'_u$  such that  $\mathcal{M}(\Phi) \subseteq \mathcal{M}(\Phi'_u) \subset \mathcal{M}(\Phi_u)$ .

The following lemma states that in order to find optimal IOF bounds for a unate  $k$ -partite formula  $\Phi$ , it is sufficient to remove clauses from  $\Phi$  or add clauses to  $\Phi$ .

**LEMMA 5.** *Let  $\Psi$  be a GLB (LUB) for a unate  $k$ -partite formula  $\Phi$  with respect to  $k$ -partite IOFs. Then  $\Psi$  is GLB (LUB) for  $\Phi$  with respect to general IOFs.*  $\square$

Both Lemmata 3 and 5 are key ingredients for finding optimal IOF bounds for unate  $k$ -partite formulas.

## 5.1 IOF Least Upper Bounds

Algorithm 3 enumerates optimal  $k$ -partite upper bounds using an approach that closely follows Lemma 3. We first *complete* all connected components of bipartite projection graphs of the input formula. Then, we *align* components from different projection graphs that have overlapping, but incomparable sets of nodes, by introducing new edges across these components. Whereas there is one possible completion for these graphs, there may be two possible alignments for any pair of components.

While the completion step implements the canonical notion of a complete bipartite graph, the alignment step needs further explanation. If two components  $C = (V, W, E)$  and  $C' = (V', W', E')$  are non-aligned, then at least one pair of variable sets  $(V, V'), (V, W'), (W, V'), (W, W')$  — assume wlog it is  $(V, V')$  — must be incomparable, i.e. satisfy  $V \cap V' \neq \emptyset$  and  $V \not\subseteq V'$  and  $V \not\supseteq V'$ . We align the components  $C, C'$  by enlarging  $C'$  such that  $V \subseteq V'$ . This is achieved by adding edges  $(V - V') \times W'$  to the projection graph of  $C'$  as shown in the *Align* procedure in Algorithm 3. All these edges are necessary for alignment; further edges could be added but they would violate the optimality of our construction of least upper bounds.

*Example 11.* The graph  $G(\Phi_1)$  from Figure 2 has three bipartite projection graphs  $\mathcal{B}(\Phi) = \{B_{xy}, B_{xz}, B_{yz}\}$ . Consider the sets of variables  $V_1 = \{y_1, y_2\}$ ,  $V_2 = \{y_3, y_4\}$ ,  $V_3 = \{y_1, y_3\}$ , and  $V_4 = \{y_2, y_4\}$ . The first two sets and the last two sets belong to different components of the projections graphs  $B_{xy}$  and  $B_{yz}$ , respectively. Any of  $V_1$  and  $V_2$  is incomparable with any of  $V_3$  and  $V_4$ , that is, they overlap yet none is included in the other. This means that their components are not aligned. Diagrams four and five in Figure 2 show two possible alignment and completion steps for  $G(\Phi_1)$ . The edge  $y_1 - z_2$  (orange) in the fourth diagram



aligns node sets  $V_1$  and  $V_3$  by extending the component containing nodes  $V_3 = \{y_1, y_3\}$  in  $B_{yz}$  by node  $y_2$ , leading to  $V_1 \subseteq \{y_1, y_2, y_3\}$ . The edges  $y_1 - z_2$ ,  $y_3 - z_2$ , and  $y_4 - z_2$  (blue) are due to the subsequent completion step in the projection graph  $B_{yz}$ .

The projection graph  $B_{xy}$  in the graph  $G(\Phi_2)$  has one component only. This component is not complete, since, for instance,  $(x_3, y_1) \notin B_{xy}$ . The rightmost diagram in Figure 2 shows  $G(\Phi_2)$  after completion. The projection graphs in the latter graph are already aligned: Indeed, the variable set  $\{y_1, \dots, y_6\}$  of the large component includes the variable sets  $\{y_1, y_3\}$ ,  $\{y_2, y_4\}$ , and  $\{y_5, y_6\}$  of the other components.  $\square$

Both steps, completion and alignment, may introduce new clauses that are necessary to obtain a 1OF formula. They detect  $P_4$  paths and add edges between their nodes that correspond to new clauses of arity  $k$ . When these steps are done, the graph is  $P_4$ -free and defines a 1OF formula.

We are now ready to state the main result of this section.

**THEOREM 3.** *Algorithm 3 computes  $k$ -partite 1OF least upper bounds of a unate  $k$ -partite formula.*  $\square$

Soundness follows by Lemmata ?? and 5 and the fact that by completion and alignment no unnecessary edges are added to the graph (in other words, no unnecessary clause is added to the input formula). In particular, note that according to Lemma 5, bounds computed by Algorithm 3 are optimal even with respect to the 1OF language, and not only to the language of  $k$ -partite 1OFs.

*Example 12.* The aligned and completed graphs of  $G(\Phi_1)$  given in Figure 2 correspond to the 1OF formulas:

$$\begin{aligned} 1: & (x_1 \vee x_2)[z_1(y_1 \vee y_3) \vee z_2(y_2 \vee y_4)] \vee x_3(y_5 z_3 \vee y_6 z_4) \\ 2: & [x_1(y_1 \vee y_2) \vee x_2(y_3 \vee y_4)](z_1 \vee z_2) \vee x_3(y_5 z_3 \vee y_6 z_4) \end{aligned}$$

These two formulas are the only 3-partite LUBs for  $\Phi_1$ .

The completion of  $G(\Phi_2)$  given in Figure 2 corresponds to the only 1OF 3-partite LUB formula for  $\Phi_2$ :

$$\Phi_U^2 = (x_1 \vee x_2 \vee x_3)[z_1(y_1 \vee y_3) \vee z_2(y_2 \vee y_4) \vee z_3 y_5 \vee z_4 y_6]$$

Example 10 shows how to compute the 1OF formula of the graph obtained after completing  $G(\Phi_2)$ .  $\square$

**THEOREM 4.** *Algorithm 3 enumerates 1OF least upper bounds for a  $k$ -partite unate DNF formula with polynomial delay.*  $\square$

The graph of a  $k$ -partite unate DNF formula can be constructed in time polynomial in  $k$ , in the size and in the number of variables of the formula, and has  $O(k^2)$  bipartite subgraphs. Each iteration of the *foreach*-loop outputs a distinct LUB, since every iteration fixes a projection graph  $B$  that is never aligned; two LUBs produced by different runs are thus incomparable. It remains to be shown that any *foreach*-iteration takes at most polynomial time. The number of alignment and completion steps — hence the number of *while*-iterations — required is bounded by the maximal number of edges in the graph,  $O(n^2)$ . Each alignment and completion step in a component of a projection graph costs at most quadratic time in the number of its nodes.

**Enumeration of all  $k$ -partite 1OF LUBs.** Algorithm 3 does not enumerate *all*  $k$ -partite 1OF LUBs of a given formula. This is due to the restriction that for each projection graph  $B$  that has a non-aligned component only one LUB

1OF-GLB ( $k$ -partite unate formula  $\Phi$ )  
**outputs** all 1OF GLBs of  $\Phi$  that are  $k$ -partite

```

Γ ← ∅ // Γ is the set of all GLBs seen so far
foreach  $n$  from  $|\Phi|$  downto 1 do
  foreach  $\Phi' \subseteq \Phi$  such that  $|\Phi'| = n$  do
    if  $\Phi'$  is in 1OF and  $\exists \Phi_0 \in \Gamma$  such that  $\Phi' \subseteq \Phi_0$ 
    then
      Output  $\rho(\Phi')$ 
      if  $\Phi' = \Phi$  then
        return
       $\Gamma \leftarrow \Gamma \cup \{\Phi'\}$ 

```

**Algorithm 4: Finding 1OF GLBs of a  $k$ -partite unate formula.**

is constructed. By removing the *foreach*-statement from the algorithm and implementing the choice in the *while*-statement to be non-deterministic, this restriction can be lifted to obtain a non-deterministic algorithm that finds all  $k$ -partite 1OF LUBs. To obtain a deterministic version of this algorithm, we iterate and branch over all such choices. Since different branches may find the same LUB several times, this algorithm does not enumerate LUBs with polynomial delay. Indeed, this extended algorithm has exponential time complexity, since (i) we can choose one of two possible alignments for each pair of non-aligned projection graphs, and (ii) the number of alignments of a  $k$ -partite graph can be exponential in the size of the graph.

**PROPOSITION 9.** *The graph of the 3-partite formula  $\Phi_n^3 = \bigvee_{i=1}^n (x_i y_{4i-3} z_{4i-3} \vee x_i y_{4i-2} z_{4i-2} \vee x_{i+1} y_{4i-1} z_{4i-3} \vee x_i y_{4i} z_{4i-2})$  has size  $4n$  and  $2^n$  alignments.*  $\square$

The graph in Proposition 9 consists of  $n$  unconnected copies of the subgraph represented by the upper part  $U$  of the graph  $G(\Phi_1)$  in Figure 2. There are exponentially many alignments since  $U$  has two alignments and the  $n$  copies of  $U$  are unconnected.

*Remark 2.* Algorithm 3 can also be used to compute  $k$ -partite 1OF upper bounds for formulas that represent lineage of conjunctive queries *with self-joins*. Those formulas are similar to  $k$ -partite ones, except that they are not over disjoint variable sets  $V_1, \dots, V_k$ , but some  $V_i, V_j$  can be identical. To obtain bounds in this case, first drop a minimal number of variables from clauses such that each clause contains at most one variable from each distinct variable set. Subsequently, Algorithm 3 can be applied to this smaller formula. An important direction of future work is to study the case of optimal bounds for such formulas as well as unrestricted unate formulas.  $\square$

## 5.2 1OF Greatest Lower Bounds

For a given  $k$ -partite unate formula  $\Phi$ , Algorithm 4 computes  $k$ -partite 1OF greatest lower bounds. It iterates over all subformulas  $\Phi'$  representing subsets of  $\Phi$  starting with the largest subformulas first. If  $\Phi'$  is a 1OF formula and is not included in any bounds previously found by the algorithm, then  $\Phi'$  is reported as a greatest lower bound.

**PROPOSITION 10.** *Algorithm 4 computes all  $k$ -partite 1OF greatest lower bounds of a unate  $k$ -partite formula.*  $\square$

Since it only checks for subsets of the input formula  $\Phi$ , the algorithm reports lower bounds that have arity  $k$ , which is sufficient according to Lemma 5. The 1OF lower bounds computed by the algorithm are optimal, since we check that there is no larger 1OF lower bound. We do this by iterating over the larger subformulas before the smaller subformulas, and also over  $\Phi$ , and memorizing all GLBs encountered so far (in  $\Gamma$ ). The algorithm is complete since it iterates over the entire space of possible  $k$ -partite lower bounds.

*Example 13.* Formula  $\Phi_1$  in Figure 1 has four 1OF GLBs:

- 1 :  $(x_1y_1 \vee x_2y_3)z_1 \vee x_3(y_5z_3 \vee y_6z_4)$
- 2 :  $x_1y_1z_1 \vee x_2y_4z_2 \vee x_3(y_5z_3 \vee y_6z_4)$
- 3 :  $(x_1y_2 \vee x_2y_4)z_2 \vee x_3(y_5z_3 \vee y_6z_4)$ .
- 4 :  $x_1y_2z_2 \vee x_2y_3z_1 \vee x_3(y_5z_3 \vee y_6z_4)$ .  $\square$

Algorithm 4 needs exponential time to report all  $k$ -partite 1OF GLBs of a unate DNF formula. This is partly due to the exponential number of such lower bounds. The problem of enumerating 1OF greatest lower bounds with polynomial delay remains open.

PROPOSITION 11. *The unate bipartite DNF formula*

$$\Phi = (x_1y_1 \vee x_1y_2 \vee x_2y_2) \vee \dots \vee (x_{2n-1}y_{2n-1} \vee x_{2n-1}y_{2n} \vee x_{2n}y_{2n})$$

has  $3n$  clauses and  $3^n$  1OF bipartite greatest lower bounds.

Any formula consisting of either  $x_{2i-1}(y_{2i-1} \vee y_{2i})$ , or  $(x_{2i-1} \vee x_{2i})y_{2i}$ , or  $(x_{2i-1}y_{2i-1} \vee x_{2i}y_{2i})$  for all  $1 \leq i \leq n$  is a 1OF bipartite greatest lower bound of  $\Phi$ .  $\square$

## 6. APPROXIMATION BY QUERIES

In this section we show how to derive queries that compute bounds for the lineage of conjunctive queries without self-joins on databases in which each tuple is annotated with a fresh Boolean variable. In case of random variables, such databases are tuple-independent probabilistic databases. Bounds for the lineage are expressed in the languages of iDNF or 1OF, as discussed in the previous sections.

Our motivation behind this study is twofold. Firstly, bounds that are expressible by queries can be computed in polynomial time using relational database engines. This is essential for probabilistic databases, where query evaluation is  $\#P$ -hard in general since it involves probability computation of arbitrary propositional formulas. Secondly, we would like to understand whether queries with polynomial data complexity on classical relational databases can be approximated by bound queries that have polynomial combined complexity and have answers representable in the iDNF or 1OF propositional languages.

### 6.1 Expressing Bound Queries

The main challenge in expressing bounds using queries is the encoding (i) of the completion and alignment steps for 1OF upper bounds, and (ii) of the choice of a maximal subset of clauses that is still in iDNF or 1OF for lower bounds. Deriving queries that express iDNF upper bound computation is subject to future work.

We first discuss the case of 1OF bounds. For a given conjunctive query  $Q$ , its lineage is already in 1OF if  $Q$  is hierarchical and has no self-joins [17]. By Definition 3, the hierarchical property is violated by relation symbols with several query variables that also appear in other distinct relation symbols. For instance, relation  $S$  in queries  $Q_1, Q_2$ ,

Query-UB (Conjunctive Query without Self-Joins  $Q$ )

---

**outputs** upper bound queries for  $Q$

**foreach** non-head query variables  $A, B$  such that not  $(s(A) \cap s(B) = \emptyset$  or  $s(A) \subseteq s(B)$  or  $s(A) \supset s(B))$  **do**

**foreach** relation  $S \in (s(A) \cap s(B))$  **do**

//Step 1: complete  $S$  on columns  $(A, B)$   
replace  $S$  by  $\mathbf{trans}_Q(S, A, B)$  in  $Q$

//Step 2: align  $S$  on columns  $(A, B)$   
 $\alpha_A(A, B, \bar{X}) :- S(A, B', \bar{X}), S^d(A', B, \bar{X}), S^d(A', B', \bar{X})$   
 $\alpha_B(A, B, \bar{X}) :- S(A', B, \bar{X}), S^d(A, B', \bar{X}), S^d(A', B', \bar{X})$

**choose** one of the following:

1. replace  $S$  by  $\alpha_A$  in  $Q$
2. replace  $S$  by  $\alpha_B$  in  $Q$

**output**  $Q$

**Algorithm 5: Deriving upper bound queries.**

and  $Q_3$  of Figure 1 violates the hierarchical property. In case of  $Q_1$ , the relation  $S$  can be constructed such that it pairs arbitrary tuples from relations  $R$  and  $T$ , and in particular arbitrary annotations from the two relations. The lineage of such queries cannot thus be in 1OF in general, though it may be in 1OF for particular extensions or reductions of such violating relations. The key insight is that one particular case of extensions or reductions of violating relations is by means of completion and alignment of their content.

Algorithm 5 extends  $Q$  with completion and alignment steps necessary for computing 1OF upper bounds, as detailed in Section 5.1. It first detects the binary projections of violating relations  $S$  that have incomplete and/or misaligned connected components. This detection is based on the hierarchical property of conjunctive queries as given in Definition 3, which corresponds precisely to the alignment check used by Algorithm 3.

The completion step requires transitive closure on the bipartite graph  $G_{AB}$  represented by the projection  $S(A, B)$ : For each tuple  $t$  from a relation joined with  $S$  via an equality or inequality condition on  $A$  and that has an  $A$ -value  $a$ , we introduce new tuples in  $S$ . These tuples allow for combinations of  $t$  with  $B$ -values that are annotated by their variables in  $S$  and that appear with  $A$ -values from the same connected component of  $G_{AB}$  with value  $a$ . We denote this closure by  $\mathbf{trans}_Q(S, A, B)$  in Algorithm 5.

*Example 14.* Consider the query  $Q :- R(A), S(A, B), T(B)$ , where the relations are defined as follows:

$R$	$A$	$S$	$A$	$B$	$T$	$B$	$S'$	$A$	$B$
$x_1$	$a_1$	$y_1$	$a_1$	$b_1$	$z_1$	$b_1$	$y_4$	$a_1$	$b_1$
$x_2$	$a_2$	$y_2$	$a_1$	$b_2$	$z_2$	$b_2$	$y_5$	$a_1$	$b_4$
		$y_3$	$a_1$	$b_3$	$z_3$	$b_3$	$y_1$	$a_2$	$b_1$
		$y_4$	$a_2$	$b_1$	$z_4$	$b_4$	$y_2$	$a_2$	$b_2$
		$y_5$	$a_2$	$b_4$			$y_3$	$a_2$	$b_3$

The completion of the violating relation  $S$  leads to the addition of the tuples  $S'$  to  $S$ . The lineage of the upper bound query  $Q^u :- R(A), \mathbf{trans}_Q(S, A, B), T(B)$  is the 1OF formula  $(x_1 \vee x_2)[(y_1 \vee y_4)z_1 \vee y_2z_2 \vee y_3z_3 \vee y_5z_4]$ . This is indeed an upper bound of the original lineage  $x_1(y_1z_1 \vee y_2z_2 \vee y_3z_3) \vee x_2(y_4z_1 \vee y_5z_4)$ , which is not in 1OF.  $\square$

After completion, the alignment step can add tuples to  $S$  so that all variables associated with either  $A$  or  $B$  are now

<p><b>Query-LB (Conjunctive Query without Self-Joins <math>Q</math>)</b>  <b>outputs</b> lower bound queries for <math>Q</math></p> <p><b>choose</b> an order <math>&lt;_Q</math> on the query variables of <math>Q</math>  <b>foreach</b> <i>non-head</i> query variables <math>A, B</math> such that not <math>(s(A) \cap s(B) = \emptyset</math> or <math>s(A) \subseteq s(B)</math> or <math>s(A) \supset s(B))</math> <b>do</b></p> <div style="margin-left: 20px;"> <p>1. In case of 1OF bounds:  <b>if</b> <math>A &lt;_Q B</math> <b>then</b>            <math>\lfloor f \leftarrow</math> functional dependency <math>(A \rightarrow B)</math>  <b>else</b>            <math>\lfloor f \leftarrow</math> functional dependency <math>(B \rightarrow A)</math></p> <p>2. In case of iDNF bounds:  <math>f \leftarrow</math> functional dependencies <math>(A \rightarrow B, B \rightarrow A)</math>  <b>foreach</b> relation <math>S</math> in <math>s(A) \cap s(B)</math> <b>do</b>            <math>\lfloor</math> replace <math>S</math> by <math>repair_f(S)</math> in <math>Q</math></p> </div> <p><b>output</b> <math>Q</math></p>
--

**Algorithm 6: Deriving lower bound queries.**

also associated with the other  $B$  or  $A$  values, respectively, in the same connected component  $G_{AB}$ . We denote alignment by  $\alpha_v$  in the algorithm;  $S^d$  is relation  $S$  without variables and  $v$  denotes our choice of aligning on either  $A$  or  $B$ . The alignment step can be encoded as a conjunctive query.

*Example 15.* In query  $Q_1$  from Figure 1,  $S$  is the only violating relation. This relation is already complete. There are two possible bound queries:  $Q_1^{u_1} :- R(A), \alpha_A(A, B), T(B)$  and  $Q_1^{u_2} :- R(A), \alpha_B(A, B), T(B)$ , where  $\alpha$  is defined as given, where in addition  $S$  is replaced by  $\mathbf{trans}_{Q_1}(S, A, B)$ . After the alignment  $\alpha_A$ , the following tuples are added to  $S$ :  $y_1(1, 2), y_2(1, 1), y_3(2, 2), y_4(2, 1), y_5(3, 4)$ , and  $y_6(3, 3)$ . With this extended relation  $S$ , the query would produce the lineage  $[x_1(y_1 \vee y_2) \vee x_2(y_3 \vee y_4)](z_1 \vee z_2) \vee x_3(y_5 \vee y_6)(z_3 \vee z_4)$ . This lineage is an 1OF upper bound of the original lineage.  $\square$

**THEOREM 5.** *Given a conjunctive query without self-joins  $Q$ , Algorithm 5 constructs queries  $Q_U$  with a transitive closure construct such that:*

- The size of  $Q_U$  is at most quadratic in the size of  $Q$ ,
- For every relational database  $D$ ,  $Q(D) \subseteq Q_U(D)$ ,
- For every annotated database  $D$ , the lineage of  $Q_U(D)$  is a  $|Q|$ -partite 1OF formula, and
- For every probabilistic database  $D$  and tuple  $t \in Q(D)$ ,  $P(t \in Q(D)) \leq P(t \in Q_U(D))$  and, given  $Q_U(D)$ ,  $P(t \in Q_U(D))$  can be computed in polynomial time.  $\square$

This result draws on the connection between Algorithms 5 and 3, since the former only provides query encodings of the completion and alignment steps of the latter. It can be further shown that the completion step is redundant when all joins are key-foreign key joins. If in addition these joins build a hierarchy based on inclusion of the query variables involved in each join, then the query is hierarchical and the alignment step is also not needed, since then the lineage is guaranteed to be in 1OF. Alternative upper bound queries without the transitive closure construct can be constructed by dropping one of the join conditions for each violating relation in the query.

Theorem 5 also holds for input queries with inequality joins ( $<$ ). In case of queries with self-joins, we can obtain upper bounds by first dropping all but one occurrence of each relation symbol in the query, and then use Algorithm 5.

To express lower bound computation using queries, we need to resolve incompleteness and misalignment of graph components of many-to-many relations by dropping tuples and hence clauses from the lineage to be approximated. This can be achieved by non-deterministically choosing maximal functional bipartite subgraphs of many-to-many binary relations, and as such, it requires a choice operator, whose expressiveness is beyond first order [20]. A similar operator has been proposed in the context of uncertain databases for constructing all possible (maximal) *repairs* of a relation under given functional dependencies [1]. Algorithm 6 uses a repair operator to choose a maximal subset of a relation that satisfies a given set of functional dependencies. An alternative approach to the choice operator is to extend the language of first-order queries with a new binary order  $<$  operator, and consider an arbitrary linear ordering of the universe. Linearly ordered finite domains provide a first-order definable choice function, namely the one that always picks  $<$ -least elements [20].

**THEOREM 6.** *Given a conjunctive query without self-joins  $Q$ , Algorithm 6 constructs queries  $Q_L$  with a repair construct such that:*

- The size of  $Q_L$  is at most quadratic in the size of  $Q$ ,
- For every relational database  $D$ ,  $Q_L(D) \subseteq Q(D)$ ,
- For every annotated database  $D$ , the lineage of  $Q_L(D)$  is a  $|Q|$ -partite 1OF or iDNF formula, and
- For every probabilistic database  $D$  and tuple  $t \in Q(D)$ ,  $P(t \in Q_L(D)) \leq P(t \in Q(D))$  and, given  $Q_L(D)$ ,  $P(t \in Q_L(D))$  can be computed in polynomial time.  $\square$

By enforcing a functional dependency  $A \rightarrow B$  on each violating relation  $S$  in the query, we turn the bipartite graph of  $S(A, B)$  into the graph of a function, and drop clauses from the lineage of the original query that forbid factorisation in 1OF. In the iDNF case, we enforce a two-way dependency  $A \leftrightarrow B$  on each violating relation  $S$ , hence each variable annotating a tuple in  $S$  can only occur in one clause in the lower bound lineage. In the 1OF case, the constructed query turns the database into a chain of one-to-many joined relations. All connected components of the projection graphs of the lineage graph are then necessarily complete and aligned, hence the lineage is in 1OF.

*Example 16.* We enforce any of the functional dependencies  $A \rightarrow B$  or  $B \rightarrow A$  on relation  $S$  in the query  $Q_1$  from Figure 1. The number of maximal functional subgraphs of  $S$  is eight in the first case and four in the second case. In the former case, the 1OF lower bound query would be  $Q_1^{l_1} :- R(A), repair_{A \rightarrow B}(S(A, B)), T(B)$ . If we choose the first tuples for each  $A$ -value in  $S$ , we obtain the lineage  $(x_1 y_1 \vee x_2 y_3) z_1 \vee x_3 y_5 z_3$ .

The query  $Q_1^{l_2} :- R(A), repair_{A \leftrightarrow B}(S(A, B)), T(B)$  is the iDNF lower bound query for  $Q_1$ . If we choose the first, fourth, and sixth tuples in  $S$ , we obtain the lineage  $x_1 y_1 z_1 \vee x_2 y_4 z_2 \vee x_3 y_5 z_3$ . The iDNF bound query is optimal since its bound is optimal. The 1OF bound query is not optimal, Example 13 lists the 1OF greatest lower bounds for  $Q_1$ .  $\square$

As exemplified above, the 1OF lower bound queries constructed by Algorithm 6 are usually not optimal. To obtain optimality, we would need a more involved version of the repair operator that can enforce any of the two functional dependencies independently for each of the connected components of the bipartite graph of  $S(A, B)$ .

## 6.2 iDNF and IOF Representation Systems

Since iDNF and IOF formulas representing the lineage of bound queries have sizes linear in the number of their variables, a natural question is whether bound queries can be evaluated more efficiently than the queries they approximate. The initial answer is “no”, since there are bound queries, such as the product of  $n$  relations, that produce answers of size exponential in  $n$ . The structure of answers to bound queries, however, follows the regular structure of IOF or iDNF languages. This motivates the use of these languages as succinct representations of query answers.

*Definition 11.* Let  $Q$  be a IOF bound query,  $D$  be an annotated database, and  $\Phi$  be the IOF lineage of the query answer  $Q(D)$ . The *IOF representation* of  $Q(D)$  is  $(\Phi_e, \mathcal{S})$ , where  $\mathcal{S}$  is the schema of the query answer and  $\Phi_e$  is  $\Phi$  extended by inserting at each variable  $x$  occurring in  $\Phi$  the input tuple  $t$  annotated with  $x$  together with the schema of  $t$ . The case of iDNF is defined analogously.  $\square$

*Example 17.* Consider the lineage of the IOF upper bound query  $Q_1^{u1}$  on database  $D$  from Example 15:  $[x_1(y_1 \vee y_2) \vee x_2(y_3 \vee y_4)](z_1 \vee z_2) \vee x_3(y_5 \vee y_6)(z_3 \vee z_4)$ . The IOF representation of  $Q_1^{u1}(D)$  is  $(\Phi_e, \mathcal{S})$  with  $\mathcal{S} = \emptyset$  and

$$\begin{aligned} \Phi_e = & \left[ x_1\langle 1 \rangle (y_1\langle 1, 1 \rangle \vee y_2\langle 1, 2 \rangle) \vee \right. \\ & \left. x_2\langle 2 \rangle (y_3\langle 2, 1 \rangle \vee y_4\langle 2, 2 \rangle) \right] (z_1\langle 1 \rangle \vee z_2\langle 2 \rangle) \\ & \vee \left[ x_3\langle 3 \rangle (y_5\langle 3, 3 \rangle \vee y_6\langle 3, 4 \rangle) \right] (z_3\langle 3 \rangle \vee z_4\langle 4 \rangle). \square \end{aligned}$$

This representation is closed under queries with projection, selection, and product operators, but not under queries with joins and projection. In contrast to classical DNF representations of query answers (equivalent to a list of answer tuples), the IOF representation is more succinct and enables lower complexity of query evaluation for IOF bound queries.

**THEOREM 7.** *Given a bound query  $Q$  as constructed by Algorithms 5 and 6 and an annotated database  $D$ , the IOF or iDNF representation of  $Q(D)$  needs space  $O(|D|)$  and can be computed in polynomial combined complexity.*

*Given an IOF or iDNF representation of  $Q(D)$ , the answers to  $Q(D)$  can be enumerated with polynomial delay.*  $\square$

## 7. CONCLUSION AND FUTURE WORK

This paper introduces a novel approach to query approximation in probabilistic and relational databases via two tractable propositional languages. We show equivalences between semantic and syntactic characterisations of optimal bounds; while the former provide an intuitive understanding of optimality, the latter are more easily accessible to efficient algorithmic treatment. We propose two methods for the computation of optimal bounds, one based on the syntax of lineage formulas, and one based on query rewriting.

We see several promising directions for future work. We plan to integrate the approximation framework into our query engine for probabilistic databases called SPROUT. The application of our results to the relational case suggests a novel paradigm to the representation and computation of query answers. A promising direction is the investigation of queries whose answers are  $k$ -readable, that is, answers that can be represented as formulas where each input tuple occurs at most  $k$  times, where  $k$  is a polynomial in the sizes of the

database and the query. The problem of enumerating IOF greatest lower bounds with polynomial delay is still open.

**Acknowledgments.** This research was funded by the FP7 ERC grant agreement FOX number FP7-ICT-233599.

## 8. REFERENCES

- [1] L. Antova, C. Koch, and D. Olteanu. “From Complete to Incomplete Information and Back”. In *SIGMOD*, 2007.
- [2] N. Dalvi, C. Ré, and D. Suciu. Probabilistic databases: diamonds in the dirt. *Commun. ACM*, 52(7), 2009.
- [3] N. Dalvi and D. Suciu. “Efficient Query Evaluation on Probabilistic Databases”. *VLDB Journal*, 16(4), 2007.
- [4] A. Darwiche and P. Marquis. “A knowledge compilation map”. *Journal of AI Research*, 17, 2002.
- [5] K. Elbassioni, K. Makino, and I. Rauf. On the Readability of Monotone Boolean Formulae. In *COCOON*, 2009.
- [6] R. Fink and D. Olteanu. On the Optimal Approximation of Queries Using Tractable Propositional Languages. [Extended version, available on author’s homepage]. 2011.
- [7] R. Fink, D. Olteanu, and S. Rath. “Providing Support for Full Relational Algebra Queries in Probabilistic Databases”. In *ICDE*, 2011.
- [8] M. Garofalakis and P. Gibbon. Approximate query processing: Taming the terabytes. In *VLDB*, 2001.
- [9] W. Gatterbauer, A. K. Jha, and D. Suciu. Dissociation and propagation for efficient query evaluation over probabilistic databases. TR UW-CSE-10-04-01, U. Washington, 2010.
- [10] M. Golumbic, A. Mintza, and U. Rotics. “Read-Once Functions Revisited and the Readability Number of a Boolean Function”. In *Int. Colloq. on Graph Theory*, 2005.
- [11] C. P. Gomes, A. Sabharwal, and B. Selman. *Handbook of Satisfiability*, chapter Model Counting. IOS Press, 2009.
- [12] T. J. Green, G. Karvounarakis, and V. Tannen. Provenance semirings. In *PODS*, 2007.
- [13] V. Gurvich. Criteria for repetition-freeness of functions in the algebra of logic. In *Soviet math. dokl.*, 43(3), 1991.
- [14] T. Imielinski and W. Lipski. “Incomplete information in relational databases”. *Journal of ACM*, 31(4), 1984.
- [15] D. Johnson, C. Papadimitriou, and M. Yannakakis. On generating all maximal independent sets. *Inf. Process. Lett.*, 1988.
- [16] C. Meinel and T. Theobald. *Algorithms and Data Structures in VLSI Design*. Springer-Verlag, 1998.
- [17] D. Olteanu and J. Huang. “Using OBDDs for Efficient Query Evaluation on Probabilistic Databases”. In *SUM*, 2008.
- [18] D. Olteanu and J. Huang. “Secondary-Storage Confidence Computation for Conjunctive Queries with Inequalities”. In *SIGMOD*, 2009.
- [19] D. Olteanu, J. Huang, and C. Koch. “Approximate Confidence Computation for Probabilistic Databases”. In *ICDE*, 2010.
- [20] M. Otto. Epsilon Logic is More Expressive than First-Order Logic over Finite Structures. *J. Symb. Log.*, 65(4), 2000.
- [21] J. Pe’er and R. Y. Pinter. Minimal decomposition of boolean functions using non-repeating literal trees. In *IFIP Workshop on Logic and Architecture Synthesis*, 1995.
- [22] Y. Sagiv and M. Yannakakis. Equivalences among relational expressions with the union and difference operators. *J. of the ACM*, 27(4), 1980.
- [23] B. Selman. Knowledge compilation and theory approximation. *J. of the ACM*, 43(2), 1996.
- [24] P. Sen, A. Deshpande, and L. Getoor. Read-once functions and query evaluation in probabilistic databases. In *VLDB*, 2010.
- [25] S. Tsukiyama, M. Ide, H. Ariyoshi, and I. Shirakawa. A New Algorithm for Generating All the Maximal Independent Sets. *SIAM J. Comput.*, 6(3), 1977.
- [26] S. Vadhan. The Complexity of Counting in Sparse, Regular, and Planar Graphs. *SIAM J. Comput.*, 32(2), 2001.

## APPENDIX

### Proof of Proposition 2

Let  $\Phi$  be a unate formula.

**Lower bounds** . Let  $\Phi_L$  be a formula obtainable by removing clauses from  $\Phi$  or adding literals to its clauses. Let  $\varphi_L$  be a clause in  $\Phi_L$ . Then there exists a clause  $\varphi$  in  $\Phi$  such that  $\varphi_L \supseteq \varphi$ . Hence  $\Phi_L$  is a lower bound for  $\Phi$  by virtue of Definition 4 and Lemma 1.

**Upper bounds** . Let  $\Phi_U$  be a formula obtainable by adding clauses to  $\Phi$  or removing literals from its clauses. Let  $\varphi$  be a clause in  $\Phi$ . Then there exists a clause  $\varphi_U$  in  $\Phi_U$  such that  $\varphi \supseteq \varphi_U$ . Hence  $\Phi_U$  is an upper bound for  $\Phi$  by virtue of Definition 4 and Lemma 1.

### Proof of Theorem 1

**MLB  $\Rightarrow$  GLB**. Let  $\Phi$  be an irreducible unate DNF formula and  $\Phi_L$  an MLB for  $\Phi$ . Assume  $\Phi_L$  is no GLB for  $\Phi$ . Then there exists an iDNF formula  $\Phi'_L$  such that

$$\mathcal{M}(\Phi_L) \subset \mathcal{M}(\Phi'_L) \subseteq \mathcal{M}(\Phi)$$

and the following properties hold:

- (i)  $\Phi_L \subseteq \Phi$  (every clause in  $\Phi_L$  is also a clause in  $\Phi$ )
- (ii) Every clause  $\varphi$  in  $\Phi$  and not in  $\Phi_L$  contains a conflicting variable, i.e. a variable  $x \in \varphi$  such that there exists a clause  $\varphi_L \in \Phi_L$  with  $x \in \varphi_L$
- (iii)  $\forall \varphi'_L \in \Phi'_L : \exists \varphi \in \Phi : \varphi'_L \supseteq \varphi$  ( $\Phi'_L \models \Phi$  and Lemma 1)
- (iv)  $\forall \varphi_L \in \Phi_L : \exists \varphi'_L \in \Phi'_L : \varphi_L \supseteq \varphi'_L$   
( $\Phi_L \models \Phi'_L$  and Lemma 1)
- (v) No two clauses  $\varphi, \bar{\varphi} \in \Phi$  satisfy  $\varphi \models \bar{\varphi}$   
( $\Phi$  is irreducible)

Properties (i) and (ii) follow from  $\Phi_L$  being an MLB for  $\Phi$ , (iii) – (iv) from  $\mathcal{M}(\Phi_L) \subset \mathcal{M}(\Phi'_L) \subseteq \mathcal{M}(\Phi)$ . We prove that  $\Phi_L$  is equivalent to  $\Phi'_L$  by case differentiation:

1.  $1 \in \Phi$ . Then since 1 has no conflicting variables with any other clause also  $1 \in \Phi_L$  and thus  $\mathcal{M}(\Phi_L) = \mathcal{M}(\Phi'_L) = \mathcal{M}(\Phi)$  which is a contradiction to the assumption that  $\Phi_L$  is no GLB for  $\Phi$ .
2.  $1 \notin \Phi$ . Let  $\varphi'_L \in \Phi'_L$  be any clause in  $\Phi'_L$  and  $\varphi \in \Phi$  a clause such that  $\varphi'_L \supseteq \varphi$  according to (iii). With respect to property (i), we have the following two cases:
  - (a)  $\varphi \in \Phi_L$ . Let  $\bar{\varphi}'_L \in \Phi'_L$  be a clause with  $\varphi \supseteq \bar{\varphi}'_L$  according to (iv). Together with  $\varphi'_L \subseteq \varphi$  from above we have  $\bar{\varphi}'_L \supseteq \varphi \supseteq \varphi'_L$ , and since none of  $\varphi, \varphi'_L, \bar{\varphi}'_L$  can be the empty clause (see case (a)),  $\varphi'_L$  and  $\bar{\varphi}'_L$  must share at least one variable. Since  $\Phi'_L \in$  iDNF, it follows that  $\varphi'_L$  and  $\bar{\varphi}'_L$  are the same clause. Thus:
$$\varphi \supseteq \bar{\varphi}'_L = \varphi'_L \supseteq \varphi \Rightarrow \varphi = \varphi'_L.$$

It follows that  $\Phi_L = \Phi'_L$  which is a contradiction to the assumption  $\mathcal{M}(\Phi_L) \subset \mathcal{M}(\Phi'_L)$ .

- (b)  $\varphi \notin \Phi_L$ . According to (ii), there is a variable  $x \in \varphi$  such that there exists a clause  $\varphi_L \in \Phi_L$  with  $x \in \varphi_L$ ; furthermore,  $\varphi_L \in \Phi$  due to (i). From  $\varphi'_L \supseteq \varphi$  it follows  $x \in \varphi'_L$ . From (iv) it follows that there exist a clause  $\bar{\varphi}'_L \in \Phi'_L$  such that  $\varphi_L \supseteq \bar{\varphi}'_L$ . We distinguish two cases:

- i.  $x \in \bar{\varphi}'_L$ . It follows  $\bar{\varphi}'_L = \varphi'_L$ , because  $\Phi'_L \in$  iDNF. From  $\varphi_L \supseteq \bar{\varphi}'_L = \varphi'_L \supseteq \varphi$  it follows that  $\varphi_L \supseteq \varphi$ . Since  $\varphi_L, \varphi \in \Phi$ , this is a contradiction to the assumption that  $\Phi$  is irreducible, property (v).
- ii.  $x \notin \bar{\varphi}'_L$ . Then according to (iii) there is a  $\bar{\varphi} \in \Phi$  with  $\bar{\varphi}'_L \supseteq \bar{\varphi}$  and thus  $x \notin \bar{\varphi}$  which in turn implies  $\bar{\varphi} \neq \varphi$  because  $x \in \varphi$ . Transitivity of  $\supseteq$  implies  $\varphi_L \supseteq \bar{\varphi}$  which is a contradiction to the assumption that  $\Phi$  is irreducible, property (v).

**GLB  $\Rightarrow$  MLB**. Assume that a GLB  $\Phi_L$  for  $\Phi$  is no MLB for  $\Phi$ . Then at least one of the two MLB properties in Definition 5 is unsatisfied. We show that in either case  $\Phi_L$  is no GLB for  $\Phi$ .

1. Assume  $\Phi_L \not\subseteq \Phi$ . Then  $\Phi_L$  contains a clause  $\varphi_L$  that is not in  $\Phi$ .
  - (a) If there is a clause  $\varphi \in \Phi$  such that  $\varphi_L \supseteq \varphi$ , then  $\varphi_L \supset \varphi$ . Let  $x$  be a variable that occurs in  $\varphi_L$  but not in  $\varphi$ . Then the iDNF formula obtained from  $\Phi_L$  by removing the variable  $x$  from  $\varphi$  has strictly more models than  $\Phi_L$  and thus  $\Phi_L$  is no GLB for  $\Phi$ .
  - (b) If there is no such clause, then  $\Phi_L \not\models \Phi$  following Lemma 1 and thus  $\Phi_L$  is no lower bound and in particular no greatest lower bound for  $\Phi$ .
2. If there is a clause  $\varphi \in \Phi$  such that  $\Phi_L \cup \{\varphi\} \in$  iDNF, then  $\mathcal{M}(\Phi_L) \subset \mathcal{M}(\Phi_L \cup \{\varphi\})$  because none of the variables in  $\varphi$  occur in  $\Phi_L$  and thus  $\Phi_L$  is no GLB for  $\Phi$ .

### Proof of Theorem 2

**MUB  $\Rightarrow$  LUB**. Assume  $\Phi_U$  is a minimal upper bound for  $\Phi$  but not an LUB for  $\Phi$ . Then there exists a better iDNF upper bound  $\Phi'_U$  for  $\Phi$  that satisfies

$$\mathcal{M}(\Phi) \subseteq \mathcal{M}(\Phi'_U) \subset \mathcal{M}(\Phi_U).$$

Using Lemma 1, the second inclusion unfolds to:

$$\begin{aligned} \mathcal{M}(\Phi'_U) \subset \mathcal{M}(\Phi_U) \\ \text{iff } \mathcal{M}(\Phi'_U) \subseteq \mathcal{M}(\Phi_U) \text{ and not } \mathcal{M}(\Phi_U) \subseteq \mathcal{M}(\Phi'_U) \\ \text{iff } \forall \varphi'_U \in \Phi'_U \exists \varphi_U \in \Phi_U. \varphi'_U \supseteq \varphi_U \\ \text{and } \exists \varphi_U \in \Phi_U \forall \varphi'_U \in \Phi'_U. \neg(\varphi_U \supseteq \varphi'_U). \end{aligned}$$

By collecting the assumptions and unfolding the definitions:

- (i)  $\forall \varphi \in \Phi : \exists \varphi'_U \in \Phi'_U : \varphi \supseteq \varphi'_U$  ( $\Phi \models \Phi'_U$ )
- (ii)  $\forall \varphi \in \Phi : \exists \varphi_U \in \Phi_U : \varphi \supseteq \varphi_U$  ( $\Phi \models \Phi_U$ )
- (iii)  $\forall \varphi'_U \in \Phi'_U : \exists \varphi_U \in \Phi_U : \varphi'_U \supseteq \varphi_U$  ( $\Phi'_U \models \Phi_U$ )
- (iv)  $\exists \varphi_U \in \Phi_U : \forall \varphi'_U \in \Phi'_U : \neg(\varphi_U \supseteq \varphi'_U)$  ( $\Phi_U \not\models \Phi'_U$ )
- (v) There is no clause  $\varphi_U \in \Phi_U$  that can be extended by a variable from  $vars(\Phi)$  and the resulting formula is still in iDNF and implied by  $\Phi$
- (vi) Every  $\varphi \in \Phi$  is a witness for at least one clause  $\varphi_U \in \Phi_U$ , i.e.  $\Phi \models \Phi_U$  according to Lemma 1

- (vii) Every clause  $\varphi_U \in \Phi_U$  has some critical witnesses in  $\Phi$ . Sentences (i)–(iv) are due to the assumption that  $\Phi_U$  is upper bound but no LUB for  $\Phi$  and sentences (v)–(vii) are the syntactical characterisation of the MUB property of  $\Phi_U$ .

Let  $\varphi_U \in \Phi_U$  be as in (iv). Let  $\varphi'_U \in \Phi'_U$  be a clause in  $\Phi'_U$  that shares at least one variable with  $\varphi_U$ . If no such

clause exists, then by transitivity of  $\supseteq$  and  $\Phi_U, \Phi'_U \in \text{iDNF}$ ,  $\Phi_U$  cannot have a witness in  $\Phi$  which is a contradiction to (vii).

Then, since  $\Phi_U, \Phi'_U \in \text{iDNF}$ ,  $\varphi_U$  can be the only clause in  $\Phi_U$  that satisfies sentence (iii) and together with  $\neg(\varphi_U \supseteq \varphi'_U)$  from (iv) we can conclude  $\varphi'_U \supset \varphi_U$ . Let  $x$  be a variable that occurs in  $\varphi'_U$ , but not in  $\varphi_U$ . We show a contradiction to the above sentences (i) – (vii) by case differentiation:

1. If  $x \notin \Phi_U$ : According to (vii),  $\varphi_U$  has a critical witness  $w \in \Phi$ . Consider Figure 3a.
  - (a) If  $w \models \varphi'_U$ : Then  $\varphi_U$  can be extended by  $x$  and the resulting formula is still in iDNF and implied by  $\Phi$ . This is a contradiction to (v).
  - (b) If  $w \not\models \varphi'_U$ : Then according to (i) there is a  $\bar{\varphi}'_U \in \Phi'_U$  such that  $w \models \bar{\varphi}'_U$ . Since  $\Phi'_U \in \text{iDNF}$ ,  $\varphi'_U$  and  $\bar{\varphi}'_U$  share no variables; due to (iii) there must be a  $\bar{\varphi}_U \in \Phi_U$  which is implied by  $\bar{\varphi}'_U$  and is different from  $\varphi_U$  due to the iDNF property of  $\Phi_U$ . From the transitivity of the implication relation it follows  $w \models \bar{\varphi}_U$  which is a contradiction to the assumption that  $w$  is a critical witness for  $\varphi_U$ .
2. If  $x \in \Phi_U$ : Let  $\bar{\varphi}_U \in \Phi_U$  such that  $x \in \bar{\varphi}_U$ , and  $w \in \Phi$  be a critical witness for  $\bar{\varphi}_U$ . Consider Figure 3b.
  - (a) If  $w \models \varphi'_U$ : As above, transitivity of  $\models$  implies  $w \models \varphi_U$  which is a contradiction to the assumption that  $w$  is a critical witness for  $\bar{\varphi}_U$ .
  - (b) If  $w \not\models \varphi'_U$ : Then according to (i) there is a  $\bar{\varphi}'_U \in \Phi'_U$  such that  $w \models \bar{\varphi}'_U$ . Since  $x \in \varphi'_U$  and  $\Phi'_U \in \text{iDNF}$ ,  $\bar{\varphi}'_U$  does not contain the variable  $x$ ; due to (iii) there must be a  $\hat{\varphi}_U \in \Phi_U$  which is implied by  $\bar{\varphi}'_U$  and does thus not contain the variable  $x$  and is hence different from  $\bar{\varphi}_U$ . Transitivity of  $\models$  implies  $w \models \hat{\varphi}_U$  which is a contradiction to the assumption that  $w$  is a critical witness for  $\bar{\varphi}_U$ .

This completes the proof for the direction  $\text{MUB} \Rightarrow \text{LUB}$ .

**LUB  $\Rightarrow$  MUB.** Assume  $\Phi_U$  is no MUB for  $\Phi$ . We need to show that whenever any of the three conditions in Definition 9 does not hold, then  $\Phi_U$  is no LUB for  $\Phi$ .

1. If there is a clause in  $\Phi$  which is not a witness of clauses in  $\Phi_U$ , then  $\Phi \not\models \Phi_U$  and thus  $\Phi_U$  is no upper bound for  $\Phi$  and in particular no least upper bound.
2. Let  $x$  be a variable such that a clause  $\varphi_U \in \Phi_U$  can be extended by  $x$  and the resulting formula is in iDNF and implied by  $\Phi$ . It is clear that  $x$  does not occur in  $\Phi_U$ , because otherwise it would not be possible to extend the formula without violating the iDNF property. Then the assignment with  $x \leftarrow \text{false}$  and all other variables *true* is a model of  $\Phi_U$  but no model of the extended formula  $\Phi_U^{\text{ext}}$ . Thus  $\mathcal{M}(\Phi_U^{\text{ext}}) \subset \mathcal{M}(\Phi_U)$  and  $\Phi_U$  is no LUB for  $\Phi$ .
3. Let  $\varphi_U \in \Phi_U$  be a clause without a critical witness in  $\Phi$ . Then removing  $\varphi_U$  from  $\Phi_U$  creates a formula  $\bar{\Phi}_U$  with the following properties:
  - Since  $\Phi_U$  is a iDNF formula, removing  $\varphi_U$  creates a formula with strictly fewer models than  $\Phi_U$ .
  - Since all witnesses of  $\varphi_U$  are non-critical, they imply other clauses in  $\Phi_U$  as well. Hence  $\bar{\Phi}_U$  is still implied by  $\Phi$ .

$\bar{\Phi}_U$  is thus a better upper bound for  $\Phi$  and  $\Phi_U$  is no LUB.

## Proof of Lemma 2 and Proposition 5

We first show the close correspondence between finding greatest lower iDNF bounds for a formula and finding maximal independent sets in its clause-dependency graph, as formalised by Lemma 2.

**PROOF OF LEMMA 2.** Let  $\Phi$  be a unate DNF formula and  $G = (\Phi, E)$  its clause-dependency graph. Then we have the following equivalences:

$$\Phi_L \subseteq \Phi \text{ is an iDNF GLB for } \Phi$$

iff All clauses in  $\Phi_L$  are pairwise independent and there is no clause in  $\Phi$  that is independent from  $\Phi_L$  [By Definition 6]

iff No two clauses in  $\Phi_L$  share a variable and all clauses in  $\Phi - \Phi_L$  share a variable with some clause in  $\Phi_L$

iff No nodes from  $\Phi_L$  are connected in  $G$  and every node in  $\Phi - \Phi_L$  is connected to some node in  $\Phi_L$  [By Definition 6]

iff  $\Phi_L$  is a maximal independent set in  $G$  [Definition of maximal independent set]

□

As the following Proposition shows, the correspondence holds of Lemma 2 can be extended to the case in which one is interested in finding the greatest lower bound with the largest probability. Such a bound represents the maximal independent set with the largest sum of node weights.

**PROPOSITION 12.** *Let  $\Phi$  be a DNF formula over a set of Boolean random variables  $x_1, \dots, x_n$  with probabilities  $p(x_i)$ , and  $G = (\Phi, E)$  its clause-dependency graph. Define node weights for nodes  $\varphi \in \Phi$  in graph  $G$  as follows:*

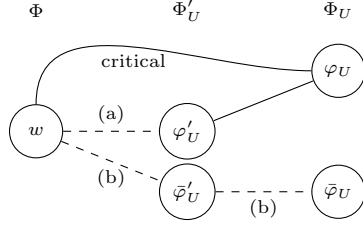
$$w(\varphi) = -\ln\left(1 - \prod_{x \in \Phi} p(x)\right). \quad (1)$$

*A formula  $\Phi_u$  is an iDNF greatest lower bound of  $\Phi$  with maximum probability among all iDNF greatest lower bounds of  $\Phi$ , if and only if  $\Phi_u$  is an independent set with maximum node weight*

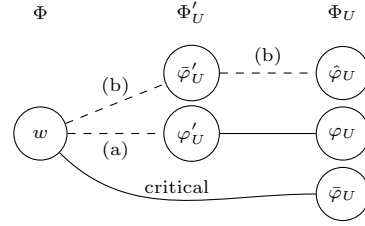
$$W(\Phi) = \sum_{\varphi \in \Phi} w(\varphi)$$

*in the node-weighted graph  $G$ .*

**PROOF.** We show the proposition by proving that the total orders on formula probabilities and weighted independent sets in  $G$  are compatible under  $w(\cdot)$ . Let  $\Phi_L$  and  $\Phi'_L$  two iDNF GLBs of  $\Phi$ .



(a) Witness graph illustration of the witnesses in  $\mathcal{M}(\Phi) \subseteq \mathcal{M}(\Phi'_U) \subset \mathcal{M}(\Phi_U)$  for the case  $x \notin \Phi_U$ . In case 1(b), the assumption that  $w$  is a critical witness for  $\varphi_U$  is disproved by clause  $\bar{\varphi}_U$  which is also implied by  $w$  due to the transitivity of the implication relation.



(b) Witness graph illustration of the witnesses in  $\mathcal{M}(\Phi) \subseteq \mathcal{M}(\Phi'_U) \subset \mathcal{M}(\Phi_U)$  for the case  $x \in \Phi_U$ . In cases 2(a) and 2(b), the assumption that  $w$  is a critical witness for  $\bar{\varphi}_U$  is disproved by clause  $\varphi_U$  or  $\bar{\varphi}_U$ , which are also implied by  $w$  due to the transitivity of the implication relation.

**Figure 3: Witness graphs used in the proof of Theorem 2.**

$$\begin{aligned}
& Pr(\Phi_L) \geq Pr(\Phi'_L) \\
\Leftrightarrow & 1 - \prod_{\varphi \in \Phi_L} (1 - \prod_{x \in \varphi} p(x)) \geq 1 - \prod_{\varphi \in \Phi'_L} (1 - \prod_{x \in \varphi} p(x)) \\
\Leftrightarrow & \prod_{\varphi \in \Phi_L} (1 - \prod_{x \in \varphi} p(x)) \leq \prod_{\varphi \in \Phi'_L} (1 - \prod_{x \in \varphi} p(x)) \\
\Leftrightarrow & \sum_{\varphi \in \Phi_L} \ln(1 - \prod_{x \in \varphi} p(x)) \leq \sum_{\varphi \in \Phi'_L} \ln(1 - \prod_{x \in \varphi} p(x)) \\
\Leftrightarrow & \sum_{\varphi \in \Phi_L} -\ln(1 - \prod_{x \in \varphi} p(x)) \geq \sum_{\varphi \in \Phi'_L} -\ln(1 - \prod_{x \in \varphi} p(x)) \\
\Leftrightarrow & \sum_{\varphi \in \Phi_L} w(\varphi) \geq \sum_{\varphi \in \Phi'_L} w(\varphi) \\
\Leftrightarrow & W(\Phi_L) \geq W(\Phi'_L)
\end{aligned}$$

It follows that the GLB with the largest probability is the independent set in the clause-dependency graph with maximum weight.  $\square$

We need the following helpers in order to show Proposition 5.

PROPOSITION 13. *Let  $k, k_1, \dots, k_m, h \in \mathbb{N}^+$  such that*

$$\begin{aligned}
& \sum_{i=1}^l k_i \leq lk \quad \text{for all } l = 1, \dots, m-1, \text{ and} \\
& \sum_{i=1}^m k_i \leq mk - h.
\end{aligned}$$

Then  $\tilde{k}_1, \dots, \tilde{k}_m \in \mathbb{N}^+$  with

$$\tilde{k}_i = \begin{cases} k_i & i < m \\ k_i + h & i = m \end{cases}$$

satisfy

$$\sum_{i=1}^l \tilde{k}_i \leq lk \quad \text{for all } l = 1, \dots, m.$$

PROOF. This is clear for  $l < m$ . For  $l = m$ :

$$\begin{aligned}
\sum_{i=1}^m \tilde{k}_i &= \left( \sum_{i=1}^{m-1} \tilde{k}_i \right) + \tilde{k}_m \\
&= \left( \sum_{i=1}^{m-1} k_i \right) + k_m + h \\
&= \left( \sum_{i=1}^m k_i \right) + h \\
&\leq mk - h + h \\
&= mk
\end{aligned}$$

$\square$

PROPOSITION 14. *Let  $k, k_1, \dots, k_m \in \mathbb{N}^+$  such that*

$$\sum_{i=1}^l k_i \leq lk \quad \text{for all } l = 1, \dots, m.$$

Let  $w_1, \dots, w_m \in \mathbb{R}^+$  such that  $w_1 \geq w_2 \geq \dots \geq w_m$ . Then

$$\sum_{i=1}^m kw_i \geq \sum_{i=1}^m k_i w_i.$$

PROOF. We prove the proposition by induction on  $m$ . For  $m = 1$  we have  $kp_1 \geq k_1p_1$  iff  $k \geq k_1$ . Now let the statement be true for  $m - 1$  and let  $k, k_1, \dots, k_m \in \mathbb{N}^+$  such that

$$\sum_{i=1}^l k_i \leq lk \quad \text{for all } l = 1, \dots, m.$$

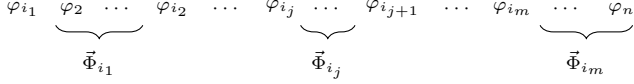
Let  $h = k_m - k$ . We distinguish two cases:

1.  $h \leq k$ : Then via the induction hypothesis the following inequalities hold:

$$\sum_{i=1}^m k_i w_i \leq \sum_{i=1}^{m-1} kw_i + k_m w_m \leq \sum_{i=1}^{m-1} kw_i + kw_m \leq \sum_{i=1}^m kw_i.$$

2.  $h > k$ : First note that

$$\sum_{i=1}^{m-1} k_i = \sum_{i=1}^m k_i - k_m \leq mk - k_m = (m-1)k - h.$$



**Figure 4: Notation used in the proof of Proposition 5.**

Hence by Proposition 13 there exist  $\tilde{k}_1, \dots, \tilde{k}_{m-1}$  with

$$\tilde{k}_i = \begin{cases} k_i & i < m-1 \\ k_i + h & i = m-1 \end{cases}$$

such that

$$\sum_{i=1}^l \tilde{k}_i \leq lk \text{ for all } l = 1, \dots, m-1.$$

By the induction hypothesis and due to the decreasing order  $w_n \geq \dots \geq w_1$  the following inequalities hold:

$$\begin{aligned} \sum_{i=1}^m k_i w_i &= \sum_{i=1}^{m-2} k_i w_i + k_{m-1} w_{m-1} + k_m w_m \\ &\leq \sum_{i=1}^{m-2} k_i w_i + (k_{m-1} + h) w_{m-1} + k_m w_m \\ &\leq \sum_{i=1}^{m-1} \tilde{k}_i w_i + k_m w_m \\ &\leq \sum_{i=1}^{m-1} k w_i + k w_m \\ &\leq \sum_{i=1}^m k w_i \end{aligned}$$

□

**PROPOSITION 15.** *Let  $k \in \mathbb{N}^+$  and  $x \in [0, 1)$ . Then*

$$\frac{1 - x^k}{1 - x} \leq k.$$

**PROOF.** We proof by induction over  $k$ . For  $k = 1$  we have  $1 = 1$ . For the inductive step, assume  $1 - x^k \leq k(1 - x)$ . We need to show  $1 - x^{k+1} \leq (k+1)(1 - x) = k(1 - x) + (1 - x)$ .

$$\begin{aligned} 1 - x^{k+1} &= x(1 - x^k) + (1 - x) \\ &\leq xk(1 - x) + (1 - x) \\ &\leq k(1 - x) + (1 - x) \end{aligned}$$

The first inequality is due to the induction hypothesis, the second due to  $0 \leq x \leq 1$ . □

**PROOF OF PROPOSITION 5.** We prove that algorithm 7 satisfies Proposition 5.

Let  $\Phi = \varphi_1 \vee \dots \vee \varphi_n$  a  $k$ -partite formula as in the proposition, such that the clauses are ordered by descending probability, i.e.  $p(\varphi_i) \geq p(\varphi_j)$  iff  $i \leq j$ . Let  $\Phi_L$  the GLB returned by the algorithm. Clearly, the algorithm runs in time  $\mathcal{O}(|\Phi| \log(|\Phi|))$ . It remains to be shown that – given that  $\Phi_L^{\text{opt}}$  is the GLB for  $\Phi$  with the highest probability –  $\Phi_L$  satisfies the inequality  $P(\Phi_L^{\text{opt}}) \leq k \cdot P(\Phi_L)$ .

**GreedyIDnfGlb ( $k$ -partite unate formula  $\Phi$ )**  
**outputs** An iDNF GLB of  $\Phi$

Sort clauses of  $\Phi$  in descending order  $\varphi_1, \dots, \varphi_n$  of their probabilities  $p(\varphi_1) \geq p(\varphi_2) \geq \dots \geq p(\varphi_n)$ .

$\Phi_L \leftarrow \emptyset$

**for**  $i = 1 \dots n$  **do**

**if**  $\text{vars}(\varphi_i) \cap \text{vars}(\Phi_L) = \emptyset$  **then**  
        [  $\Phi_L \leftarrow \Phi_L \cup \{\varphi_i\}$  ]

**Output**  $\Phi_L$

**Algorithm 7: Algorithm for finding an iDNF maximum lower bound for a  $k$ -partite formula whose probability is at most factor  $k$  smaller than the probability of the best iDNF GLBs.**

We use the notation depicted in Figure 4:  $\Phi_L = \varphi_{i_1} \vee \varphi_{i_2} \vee \dots \vee \varphi_{i_m}$ ; for every  $\varphi_{i_j} \in \Phi_L$  the set of clauses  $\vec{\Phi}_{i_j}$  are the clauses between  $\varphi_{i_j}$  and  $\varphi_{i_{j+1}}$  for  $j < m$ , and the clauses after  $\varphi_{i_j}$  for  $j = m$ :

$$\vec{\Phi}_{i_j} = \begin{cases} \{\varphi_{i_{j+1}}, \dots, \varphi_{i_{j+1}-1}\} & \text{if } j < m, \\ \{\varphi_{i_{j+1}}, \dots, \varphi_n\} & \text{if } j = m. \end{cases}$$

With respect to the order  $\varphi_1, \dots, \varphi_n$ , the existence of the GLB  $\Phi_L$  imposes the following structural constraints on  $\Phi$ :

1. Every clause  $\varphi \in \vec{\Phi}_{i_j}$  shares at least one variable with one of  $\varphi_{i_1}, \dots, \varphi_{i_j}$ . This is due to the greedy selection in the algorithm that constructs  $\Phi_L$ .
2. For each  $l \leq m$  the set of clauses  $\bigcup_{j=1}^l \vec{\Phi}_{i_j}$  contains at most  $l \cdot k$  independent clauses. Assume on contrary that  $\bigcup_{j=1}^l \vec{\Phi}_{i_j}$  contains  $h > lk$  pairwise independent clauses. Since none of these clause is in  $\Phi_L$ , they must each share at least one variable with one of  $\varphi_{i_1}, \dots, \varphi_{i_l}$ ; those  $h$  variables are pairwise different by the assumption that the  $h$  clauses are pairwise independent. This is a contradiction to the fact that the  $l$   $k$ -partite clauses  $\varphi_{i_1}, \dots, \varphi_{i_l}$  only contain  $lk$  variables.

The above structural constraints and the weight-ordered probabilities  $w(\varphi_1) \geq \dots \geq w(\varphi_n)$  imply the conditions of Proposition 14. It follows that  $W(\Phi_L^{\text{opt}}) \leq \sum_{j=1}^m k w(\varphi_{i_j})$ .

Now note that the mapping from clauses probabilities to weights in equation (1) is strictly monotonic and thus the clause order  $\varphi_1, \dots, \varphi_n$  is also decreasing with respect to clause weights  $w(\varphi_i)$ . In the light of Proposition 12 this implies that we can reason about sums of clause weights and formula probabilities equivalently and hence  $P(\Phi_L^{\text{opt}}) \leq 1 - \prod_{j=1}^m (1 - P(\varphi_{i_j}))^k$ . We thus have

$$\begin{aligned} \frac{P(\Phi_L^{\text{opt}})}{P(\Phi_L)} &\leq \frac{1 - \prod_{j=1}^m (1 - P(\varphi_{i_j}))^k}{1 - \prod_{j=1}^m (1 - P(\varphi_{i_j}))} \\ &= \frac{1 - \left[ \prod_{j=1}^m (1 - P(\varphi_{i_j})) \right]^k}{1 - \prod_{j=1}^m (1 - P(\varphi_{i_j}))} \\ &\leq \frac{1 - x^k}{1 - x} \\ &\leq k. \end{aligned}$$



The last inequality is due to Proposition 15 with

$$0 \leq x = (1 - \prod_{j=1}^m P(\varphi_{i_j})) < 1.$$

□

### Proof of Lemmata 3 and 4

We need the following helpers in order to prove Lemma 3. We start with several facts about the connection between  $k$ -partite formulas and their graph representations.

**PROPOSITION 16.** *Let  $\Phi, \Psi$  be  $k$ -partite formulas. If  $\mathcal{M}(\Phi) \subset \mathcal{M}(\Psi)$  then formula  $\Phi$  contains a clause that is not in not a clause in formula  $\Psi$ .*

**PROOF.** According to Lemma 1 the following implications hold:

$$\mathcal{M}(\Phi) \subset \mathcal{M}(\Psi)$$

$$\Rightarrow \forall(\varphi \in \Phi) \exists(\psi \in \Psi) : \varphi \supseteq \psi \text{ and } \exists(\psi \in \Psi) \forall(\varphi \in \Phi) : \psi \not\supseteq \varphi$$

$$\Rightarrow \exists(\psi \in \Psi) \forall(\varphi \in \Phi) : \psi \not\supseteq \varphi$$

$$\Rightarrow \exists(\psi \in \Psi) \forall(\varphi \in \Phi) : \psi \neq \varphi$$

□

**PROPOSITION 17.** *The formula represented by the graph of a  $k$ -partite formula is  $k$ -partite.*

**PROOF.** Let  $\Phi$  be a  $k$ -partite formula. We need to show that all maximal cliques in the graph  $G(\Phi)$  are contain  $k$  nodes.

First, the graph of  $\Phi$  can by definition not contain a clique that contains more than  $k$  nodes, since this would require at least one edge between two nodes from the same variable set and the latter cannot exist since  $\Phi$  is  $k$ -partite.

The graph can also not contain a maximal clique with less than  $k$  nodes: Assume, there is a node that is part of a maximal clique of size  $l$  where  $l < k$ . Then there is a node that is not part of a  $k$ -clique. However, by construction of the graph all nodes are part of at least one  $k$ -clique, which is a contradiction. □

Note that not every  $k$ -partite formula is normal. Consider for example  $a_1b_1c_2 \vee a_2b_2c_2 \vee a_2b_1c_1$ . Its graph represents the formula  $a_1b_1c_2 \vee a_2b_2c_2 \vee a_2b_1c_1 \vee a_2b_1c_2$ . However, as Lemma 4 shows, this graph only represents clauses that are necessarily in every  $k$ -partite 1OF upper bound.

**PROOF OF LEMMA 4.** We denote by  $G(\cdot)$  the function that maps a formula to its graph and by  $F(\cdot)$  the function that maps a graph to its represented formula. Note that  $F(G(\Phi)) = \Phi$  iff  $\Phi$  is normal.

Let  $\Phi$  be a formula,  $G(\Phi)$  its graph and  $\Phi_u = F(G(\Phi))$  the formula formula represented by  $G(\Phi)$ .

It follows from Proposition 17 that  $\Phi_u$  is  $k$ -partite. Since by construction every clause in  $\Phi$  is a  $k$ -clique in  $G(\Phi)$ , it is thus also a maximal clique in  $G(\Phi)$ . Thus every clause in  $\Phi$  is also a clause in  $\Phi_u$  and  $\mathcal{M}(\Phi) \subseteq \mathcal{M}(\Phi_u)$ .

We prove by contradiction that there is no smaller normal  $k$ -partite formula: Assume that there is a normal  $k$ -partite formula  $\Phi'_u$  such that  $\mathcal{M}(\Phi) \subseteq \mathcal{M}(\Phi'_u) \subset \mathcal{M}(\Phi_u)$ . The formulas  $\Phi, \Phi'_u, \Phi_u$  are  $k$ -partite and furthermore  $F(G(\Phi'_u)) = \Phi'_u$  and  $F(G(\Phi_u)) = \Phi_u$  and  $G(\Phi) = G(\Phi_u)$ .

It follows from  $\mathcal{M}(\Phi'_u) \subset \mathcal{M}(\Phi_u)$  that  $\Phi_u$  contains a clause that is not in  $\Phi'_u$  by Proposition 16. Let this clause be  $\alpha = x_1 \wedge \dots \wedge x_k$ . Since  $\alpha \notin \Phi'_u$  one of the edges  $(x_i, x_j)$  is not in  $G(\Phi'_u)$ . [Because otherwise  $x_1, \dots, x_k$  is a clique in  $G(\Phi'_u)$  and thus a clause in  $\Phi'_u$ .] Let  $(x_k, x_l)$  be such a missing edge. It follows that  $\Phi'_u$  contains no clause with variables  $x_k$  and  $x_l$ .

Since  $x_1 \wedge \dots \wedge x_k$  is a clause in  $\Phi_u$  and  $G(\Phi) = G(\Phi_u)$  it follows that all edges  $(x_i, x_j)$  are in  $G(\Phi)$ . In particular, edge  $(x_k, x_l)$  is in  $G(\Phi)$  and thus by construction of the graph the formula  $\Phi$  contains a clause containing both variables  $x_k$  and  $x_l$ . This contradicts the assumption  $\mathcal{M}(\Phi) \subseteq \mathcal{M}(\Phi'_u)$  which implies

$$\forall(\varphi \in \Phi) : \exists(\varphi' \in \Phi'_u) : \varphi \supseteq \varphi',$$

since there is no clause in  $\Phi'_u$  containing both variables  $x_k$  and  $x_l$ . □

**PROPOSITION 18.** *Let  $G(\Phi)$  be the graph of a  $k$ -partite formula over variable sets  $V_1, \dots, V_k$ , and  $H$  a subgraph of  $G$  induced by a subset  $\mathcal{V} \subseteq \{V_1, \dots, V_k\}$  of the variable sets. Every set of nodes in  $H$  that contains at most one node from each variable set in  $\mathcal{V}$  can be extended to a maximal clique of size  $|\mathcal{V}|$  in  $H$ .*

**PROOF.** This follows from Proposition 17. □

**PROPOSITION 19.** *Any induced subgraph of a  $P_4$ -free graph is  $P_4$ -free.*

**PROOF.** This follows directly from the definition of  $P_4$ -freeness. □

**PROPOSITION 20.** *Let  $B = (V_1, V_2, E)$  be a bipartite graph. Every  $P_4$  in  $B$  contains exactly two nodes from  $V_1$  and two nodes from  $V_2$ . Four such nodes  $a, b \in V_1$  and  $c, d \in V_2$  form a  $P_4$  if and only if  $E$  contains exactly three of the following edges:  $(a, c), (a, d), (b, c), (b, d)$ .*

**PROOF.** First note that paths in bipartite graphs visit nodes from  $V_1$  and  $V_2$  in alternating order. A  $P_4$  in a bipartite graph must thus contain exactly two nodes from  $V_1$  and two nodes from  $V_2$ . Let  $a, b \in V_1$  and  $c, d \in V_2$  be four such nodes. Since  $B$  is bipartite,  $(a, b) \notin E$  and  $(c, d) \notin E$ . If  $E$  contains less than three of the edges  $(a, c), (a, d), (b, c), (b, d)$ , then  $a, b, c, d$  do not form a path, hence no  $P_4$ . If  $E$  contains all four of those edges, then  $a, b, c, d$  form no  $P_4$  by definition. Conversely, let  $a, b, c, d$  be a  $P_4$  in  $B$ . Then  $E$  contains exactly three of the edges  $(a, b), (a, c), (a, d), (b, c), (b, d), (c, d)$ . Since  $(a, b) \notin E$  and  $(c, d) \notin E$ ,  $E$  must contain exactly three of the remaining edges  $(a, c), (a, d), (b, c), (b, d)$ . □

The first lemma is a fairly general syntactic characterisation of bi-partite 1OF formulas.

**LEMMA 6.** *A bi-partite graph  $B = (V_1, V_2, E)$  is  $P_4$ -free if and only if all connected components of  $B$  are complete bipartite graphs.*

**PROOF.** “ $\Leftarrow$ ”: Let  $B$  be a bi-partite graph in which all connected components are complete graphs. We show that  $B$  does not contain a  $P_4$ . Firstly, there is trivially no path and hence no  $P_4$  that spans more than one connected component of  $B$ . Let  $C$  be a complete bipartite connected component of  $B$  over disjoint node sets  $W_1 \subseteq V_1$  and  $W_2 \subseteq V_2$ . If  $|W_1| < 2$  or  $|W_2| < 2$  then  $C$  does not contain a path with 4 different nodes, hence is  $P_4$ -free. Otherwise, assume

$C$  contains a  $P_4$ . According to Proposition 20, the  $P_4$  must span four distinct nodes  $a, b \in W_1$  and  $c, d \in W_2$ . Since  $C$  is complete, it contains edges  $(a, c), (a, d), (b, c), (b, d)$  and thus nodes  $a, b, c, d$  are no  $P_4$  according to Proposition 20.

“ $\Rightarrow$ ”: For the other direction, first note that  $P_4$ -freeness of  $B$  implies  $P_4$ -freeness of all connected components of  $B$ . We can thus assume wlog that  $B$  has only one connected component. Let  $S$  be the following statement to be proven: *If a connected bi-partite graph  $(V_1, V_2, E)$  is  $P_4$ -free, then it is complete.* The proof of  $S$  is by induction on the numbers of vertices  $|V_1|$  and  $|V_2|$  of the graph. The base cases  $|V_1| + |V_2| \leq 4$  are as follows: If  $|V_1| < 2$  or  $|V_2| < 2$ , then  $B$  is complete since it is connected. If  $|V_1| = |V_2| = 2$ , then the graph contains at least three edges since it is connected. According to Proposition 20 the graph must contain four edges if it is  $P_4$ -free. Hence the graph is complete.

We first prove  $S$  for fixed  $|V_1| = 2$  by induction over  $|V_2|$ . The base case  $|V_1| = |V_2| = 2$  is as above. For the inductive step, assume that  $S$  holds for  $|V_1| = 2, |V_2| = m$ . Let  $B' = (V_1, V_2, E')$  be a bi-partite graph with  $|V_1| = 2, |V_2| = m + 1$ . Assume that  $B'$  is  $P_4$ -free; it remains to be shown that  $B'$  is complete. Pick a node  $v \in V_2$  and let  $B = (V_1, V_2 \setminus \{v\}, E)$  be the subgraph of  $B'$  induced by  $V_1, V_2 \setminus \{v\}$ . Since  $B$  is an induced subgraph of  $B'$  and  $B'$  is  $P_4$ -free,  $B$  is  $P_4$ -free by Proposition 19, and furthermore complete by induction hypothesis, since  $|V_1| = 2$  and  $|V_2 \setminus \{v\}| = m$ . Since  $B'$  is connected, it must contain an edge between  $v$  and some node from  $V_1 = \{x_1, x_2\}$ , assume wlog it is  $x_1$ . Since  $B$  is complete and an induced subgraph of  $B'$ ,  $B'$  contains edges between  $x_1, x_2$  and a node  $y_1 \in V_2$ , and hence the path  $v - x_1 - y_1 - x_2$ . Since  $B'$  is  $P_4$ -free, it must also contain the edge  $v - x_2$  according to Proposition 20. Thus  $v$  is connected to all vertices of  $V_1$  in  $B'$ , and together with the completeness of  $B$  it follows that  $B'$  is a complete bi-partite graph. Hence  $S$  holds on all bi-partite graphs with  $|V_1| = 2$ .

The proof for arbitrary  $|V_1|$  is similar and by induction over  $|V_1|$ : The base case  $|V_1| = 2, |V_2| = m$  is proven above. For the inductive step, let  $S$  be satisfied for  $|V_1| = n, |V_2| = m$ . Let  $B' = (V_1, V_2, E')$  be a bi-partite graph with  $|V_1| = n + 1, |V_2| = m$ . Assume that  $B'$  is  $P_4$ -free; it remains to be shown that  $B'$  is complete. Pick a node  $v \in V_1$  and let  $B = (V_1 \setminus \{v\}, V_2, E)$  be the subgraph of  $B'$  induced by  $V_1 \setminus \{v\}, V_2$ . Since  $B$  is an induced subgraph of  $B'$  and  $B'$  is  $P_4$ -free,  $B$  is  $P_4$ -free by Proposition 19, and furthermore complete by induction hypothesis, since  $|V_1 \setminus \{v\}| = n$  and  $|V_2| = m$ . Since  $B'$  is connected, it must contain an edge between  $v$  and some node from  $V_2 = \{y_1, \dots, y_m\}$ , assume wlog it is  $y_1$ . Let  $x_1 \in V_1 \setminus \{v\}$ . Since  $B$  is complete and an induced subgraph of  $B'$ ,  $B'$  contains edges between  $x_1$  and all nodes  $y_i \in V_2$ . Hence,  $B'$  contains the paths  $v - y_1 - x_1 - y_i$  for all  $1 \leq i \leq m$ . Since  $B'$  is  $P_4$ -free, it must also contain all edges  $v - y_i$  according to Proposition 20. Thus  $v$  is connected to all vertices  $y_i \in V_2$  in  $B'$ , and together with the completeness of  $B$  it follows that  $B'$  is a complete bi-partite graph. Hence  $S$  holds on all bi-partite graphs.  $\square$

The following second lemma is the central connection between alignment of projection graphs and  $P_4$ -freeness of graphs of  $k$ -partite formulas.

LEMMA 7. *Let  $\mathcal{B}(\Phi)$  be the set of projection graphs of a  $k$ -partite formula  $\Phi$  over variable sets  $V_1, \dots, V_k$ , such that all connected components in all projections  $B_{i,j} \in \mathcal{B}$  are complete bi-partite graphs.*

1. *If there exist projection graphs  $B_{i,j}$  and  $B_{l,m}$  and con-*

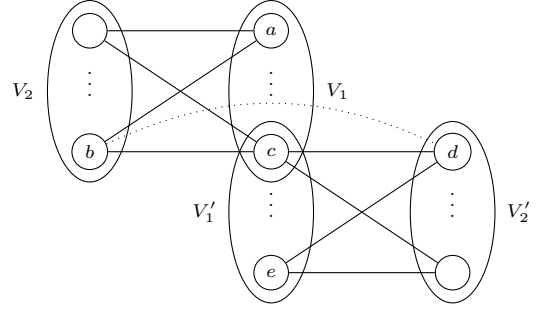


Figure 5: Graphical illustration for the proof of part 1 of Lemma 7

nected components  $C \in B_{i,j}, C' \in B_{l,m}$  in those projection graphs such that  $C$  and  $C'$  are not aligned, then the graph of  $\Phi$  contains a path isomorphic to  $P_4$ .

2. *Conversely, if all such components  $C, C'$  in all projection graphs are aligned, then the graph of  $\Phi$  is  $P_4$ -free.*

PROOF. First, note that a projection graph  $B_{i,j}$  of a projection  $\Phi_{i,j}$  is exactly the subgraph of the graph  $G(\Phi)$  of  $\Phi$  induced by variable sets  $V_i, V_j$ . Edges in  $B_{i,j}$  are thus also edges in  $G(\Phi)$ .

For part 1, let  $C = (V_1, V_2, E)$  and  $C' = (V_1', V_2', E')$  be non-aligned complete components. By definition of alignment, there exists  $1 \leq i, j \leq 2$  such that  $V_i \cap V_j' \neq \emptyset$  and  $V_i \not\subseteq V_j'$  and  $V_i \not\supseteq V_j'$ . Wlog, assume that this holds for  $i = j = 1$ . Then there exists distinct variables  $a, c, e$  with  $c \in V_1$  and  $c \in V_1', a \in V_1$  and  $a \notin V_1'$ , and  $e \notin V_1$  and  $e \in V_1'$ . Additionally,  $C$  and  $C'$  must be components from different projection graphs, since otherwise they would be connected yet not complete. Thus  $V_2 \cap V_2' = \emptyset$ , and there exist variables  $b, d$  with  $b \in V_2$  and  $b \notin V_2'$ , and  $d \notin V_2$  and  $d \in V_2'$ . The situation is depicted in Figure 5. If  $(b, d) \in G(\Phi)$ , then  $a - b - d - e$  is a  $P_4$  in  $G(\Phi)$ , otherwise  $a - b - c - d$  is a  $P_4$  in  $G(\Phi)$ .

For part 2, assume that  $G(\Phi)$  contains a  $P_4$ . We prove the statement by case differentiation on number of variable sets  $V_i$  that are visited by the  $P_4$ .

- There is no path of length four within a single variable set  $V_i$ , since there are no edges between variables/nodes of the same set.
- If the  $P_4$  visits exactly two variable sets  $V_i$  and  $V_j$ , then the projection graph  $B_{i,j}$  contains a component that is not complete by Proposition 20. This is a contradiction to the assumption that all such components are complete bi-partite graphs.
- The  $P_4$  contains nodes from exactly three variable sets  $V_i, V_j, V_k$ . The following construction shows that the assumptions of an existing  $P_4$  and aligned and complete projection graphs is contradictory, compare Figure 6. Assume that  $G$  contains such a  $P_4$  and all components of projection graphs are complete and aligned.

Diagram (a) shows the  $P_4$  in  $G$  that stretches over three variable sets  $\{a_i\}, \{b_i\}, \{c_i\}$ . The path  $b_2 - c_2$  is extendible to a clique of size 3 according to Proposition 18. The third of this clique cannot be  $a_2$ , since then the edge  $a_2 - c_2$  would annul the  $P_4$ , see diagram (b). Hence  $G$  contains an additional node  $a_3$  which forms a clique with  $b_2$  and  $c_2$ , see diagram (c). By the

same reasoning,  $b_1 - c_2$  must be part of a clique of size 3. If the third node is  $a_2$  then the edge  $a_2 - c_2$  annuls the  $P_4$ . If the third node is  $a_3$  – see diagram (d) – then due the assumption of complete bi-partite projection graphs, also  $a_2 - b_2$  would be an edge in  $G$  which again annuls the  $P_4$ . Hence  $G$  contains a node  $a_1$  that is part of the clique  $a_1, b_1, c_2$ , see diagram (e). Analogously,  $G$  must contain a node  $c_1$  to extend  $a_2 - b_1$  to a clique without annulling the  $P_4$ , see diagrams (f) and (g).

Diagram (h) is equivalent to (g) except for graphical reordering and edges between  $b$  and  $c$  nodes left out. Clearly, the components  $C_b = (\{b_1\}, \{a_1, a_2\})$  and  $C_c = (\{c_2\}, \{a_1, a_3\})$  are not aligned. Hence,  $G$  must contain the edge  $b_1 - a_3$  or the edge  $a_2 - c_2$  to achieve alignment between the abovementioned components. In the former case – diagram (i), due to the assumption of complete components also  $a_2 - b_2$  is an edge in  $G$  which annuls the  $P_4$ . In the latter case, edge  $a_2 - c_2$  directly annuls the  $P_4$ , see diagram (j).

This constructions shows that the above assumption is contradictory, i.e.  $G$  cannot be complete and aligned and contain a  $P_4$ .

- The  $P_4$  contains nodes from exactly four variable sets  $V_i, V_j, V_k, V_l$ . Let  $G$  contain the  $P_4$   $a_1 - b_1 - c_1 - d_1$ . Then there according to Proposition 18 there must exist a node  $b_2$  from the same variable set as  $b_1$  to extend the clique  $c_1, d_1$ . Hence  $a_1 - b_1 - c_1 - b_2$  is a  $P_4$  in  $G$  that stretches over three variable sets; the rest follows as above.

□

We are now in position to prove Lemma 3:

PROOF OF LEMMA 3. Let  $G$  be the graph of a unate  $k$ -partite formula  $\Phi$ .  $\Phi$  is 1OF if and only if  $\Phi$  is normal and its graph is  $P_4$ -free [13].  $G$  is normal because all maximal cliques are by construction of size  $k$  and correspond exactly to the clauses of  $\Phi$ . It thus remains to be proven that  $G$  is  $P_4$ -free if and only if the connected components of the bipartite subgraphs in  $G$  are complete and pairwise aligned. We prove the direction ( $\Rightarrow$ ) by contraposition: If there is a connected component of a bipartite subgraph in  $G$  that is not complete, then  $G$  has a  $P_4$  by virtue of Lemma 6. If there exist non-aligned components in any of the projections, then  $G$  has a  $P_4$  according to case 1 of Lemma 7.

The converse direction ( $\Leftarrow$ ) is case 2 of Lemma 7. □

## Proof of Lemma 5

We start with LUBs. The first proposition ensures that there exist  $k$ -partite LUBs for a  $k$ -partite unate formula. In particular, the proposition proves that alternative bounds with clauses smaller than  $k$  are necessarily incomparable to  $k$ -partite ones.

PROPOSITION 21. *Let  $\Phi$  be an unate formula whose shortest clause has arity  $k$ ,  $\Phi_i$  a formula that contains a clause with arity  $i$ ,  $i < k$ . Then  $\Phi_i \not\models \Phi$ .*

PROOF. Let  $\varphi_i \in \Phi_r$  be a clause with arity  $i$ . Assume  $\Phi_i \models \Phi$ . Then according to Lemma 1 there must be a clause  $\varphi \in \Phi$  such that  $\varphi_i \supseteq \varphi$ . This is a contradiction to the assumption that all clauses in  $\Phi$  have at least arity  $k$ . □

The next proposition shows that 1OF bounds comprising clauses with arity larger than  $k$  are necessarily worse than  $k$ -partite bounds.

PROPOSITION 22. *Let  $\Phi$  be a unate  $k$ -partite formula over variable sets  $V_1, \dots, V_k$  and  $\Phi_u^l$  be an 1OF upper bound for  $\Phi$  that contains a clause of arity  $l$ ,  $k < l$ .  $\Phi_u^l$  is no 1OF LUB for  $\Phi$ .*

PROOF.  $\Phi_u^l$  necessarily contains variables that are not from any of the sets  $V_1, \dots, V_k$ ; denote the set of those variables  $X$ . Define the formula  $\Phi_u$  by removing all clauses from  $\Phi_u^l$  that contain a variable from  $X$ . We show that  $\Phi_u$  is a tighter upper 1OF bound than  $\Phi_u^l$ . Non of the removed clauses is implied by a clause in  $\Phi$ ; it follows  $\mathcal{M}(\Phi) \subseteq \mathcal{M}(\Phi_u)$  due to Lemma 1. Furthermore, we have  $\mathcal{M}(\Phi_u) \subset \mathcal{M}(\Phi_u^l)$ , since at least one clause (with arity  $l$ ) was removed and this clause cannot be redundant, because  $\Phi_u^l$  is in 1OF and hence irreducible.

It remains to be shown that  $\Phi_u$  is a 1OF formula. By assumption,  $\Phi_u^l$  is in 1OF, i.e. its graph  $G$  is normal and  $P_4$ -free. Removing all clauses that contain variables from  $X$  means to remove all nodes corresponding to variables in  $X$  from  $G$ . It is clear that this operation removes exactly the maximal cliques from  $G$  that correspond to the removed clauses; all other maximal cliques remain unchanged. The graph of  $\Phi_u$  is thus normal. Furthermore, if a given graph is  $P_4$ -free, then removing nodes from the graph cannot introduce a  $P_4$ . (To see this, consider the contraposition: Let  $K$  be a graph with that contains a  $P_4$   $(a, b, c, d)$ . Adding nodes and edges connecting those nodes among themselves or to nodes in  $K$  does not connect any of the four nodes  $(a, b, c, d)$ ;  $(a, b, c, d)$  is thus still a  $P_4$  in the larger graph.) The graph of  $\Phi_u$  is thus  $P_4$ -free and normal and  $\Phi_u$  is indeed a tighter 1OF upper bound than  $\Phi_u^l$ . □

This completes the proof for LUBs. We now turn to GLBs. Proposition 21 shows that formulas with arity smaller  $k$  cannot be lower bounds for  $k$ -unate formulas. It remains to be shown that bounds that contain clauses with arity larger than  $k$  are not optimal. This is ensured by the following proposition:

PROPOSITION 23. *Let  $\Phi$  be a unate  $k$ -partite formula over variable sets  $V_1, \dots, V_k$  and  $\Phi_l^l$  be an 1OF lower bound for  $\Phi$  that contains a clause of arity  $l$ ,  $k < l$ .  $\Phi_l^l$  is no 1OF GLB for  $\Phi$ .*

PROOF.  $\Phi_l^l$  necessarily contains variables that are not from any of the sets  $V_1, \dots, V_k$ ; denote the set of those variables  $X$ . We first note that, according to Lemma 1, variables from  $X$  can only occur in clauses in  $\Phi_l^l$  that also contain  $k$  variables from  $V_1, \dots, V_k$ .

Define the formula  $\Phi_l$  by removing all variables from set  $X$  from all clauses in  $\Phi_l^l$ . We show that  $\Phi_l$  is a tighter lower 1OF bound than  $\Phi_l^l$ . Let  $\varphi_l^l$  be a clause in  $\Phi_l^l$ . According to Lemma 1,  $\varphi_l^l$  implies a clause in  $\Phi$ . Removing  $X$ -variables from  $\varphi_l^l$  does not change this property, since  $\Phi$  does not contain variables from  $X$ ; it follows  $\mathcal{M}(\Phi_l) \subseteq \mathcal{M}(\Phi)$ . Furthermore, we have  $\mathcal{M}(\Phi_l^l) \subset \mathcal{M}(\Phi_l^l)$ , since at least one variables was removed from some clause with arity  $l$  and this clause cannot be redundant, because  $\Phi_l^l$  is in 1OF and hence irreducible.

It remains to be shown that  $\Phi_l$  is a 1OF formula. This follows completely analogously to the proof of Proposition 22. □

## Comparison with Reference [9]

Reference [9] proposes a technique to derive upper bounds on the probability of conjunctive queries without self-joins.

We show that these bounds are not model-based, and in particular not optimal with respect to IOF bounds (but we conjecture that they are optimal with respect to iDNF bounds). The following arguments require an understanding of the technique proposed in [9].

We give two counter-examples, one for a query that actually yields IOF lineage and one for which the completion-based approach can be better, depending on the probability distribution of the random variables.

*Example 18.* [IOF lineage] Consider the database from Figure 7a. The query  $Q \leftarrow R(x), S(x, y), T(y)$  has lineage  $(x_1 \vee x_2)(y_1 \vee y_2)(z_1 \vee z_2) \in \text{IOF}$  with probability

$$\begin{aligned} & [1 - (1 - P(x_1))(1 - P(x_2))][1 - (1 - P(y_1))(1 - P(y_2))] \\ & [1 - (1 - P(z_1))(1 - P(z_2))] \end{aligned}$$

The dissociation  $R(x), S(x, y), T(x, y)$  proposed in [9] yields the probability

$$1 - \prod_{i=1,2} \left[ 1 - P(x_i) \left( 1 - \prod_{m,n=1,2} (1 - P(y_m)P(z_n)) \right) \right]$$

which is larger than the exact probability in all but pathological cases.  $\square$

*Example 19.* [Comparison of dissociation and alignment] Consider the database from Figure 7b. The query

$$Q \leftarrow R(x), S(x, y), T(y)$$

has lineage  $x_1(y_1 z_1 \vee y_2 z_2) \vee x_2 y_3 z_1$  that is not IOF. The dissociation method allows for two minimal safe dissociations,

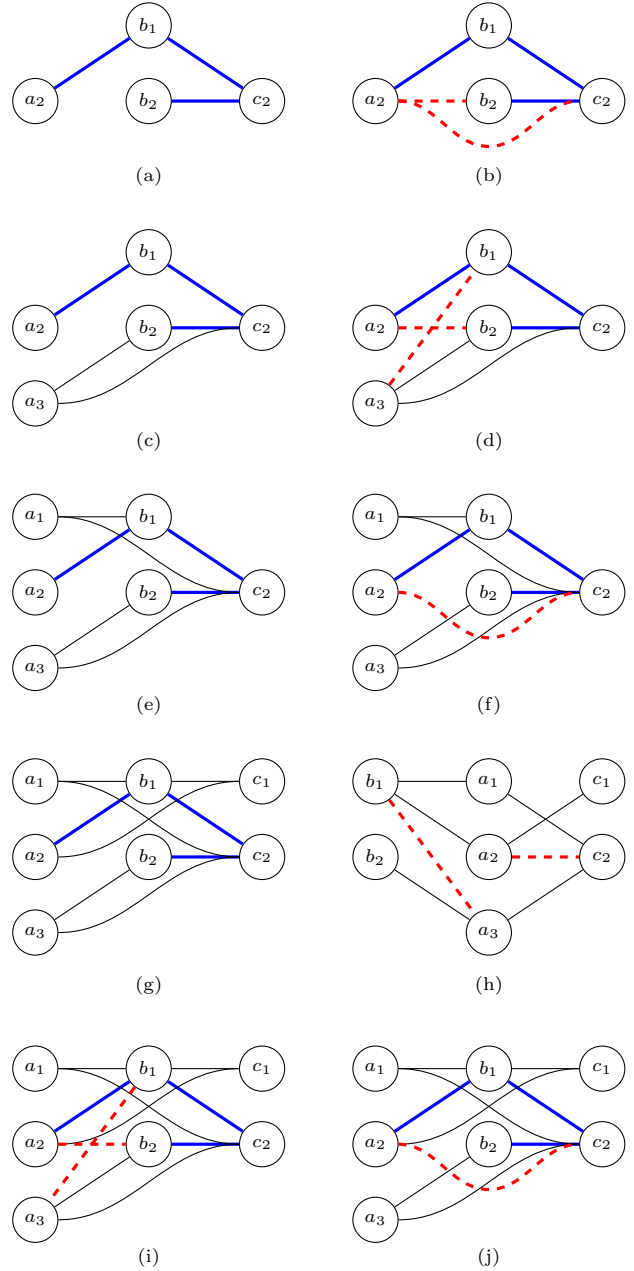
$$R(x), S(x, y), T(x, y) \quad (\text{DissociationRight}),$$

$$R(x, y), S(x, y), T(z) \quad (\text{DissociationLeft}).$$

Similarly, the alignment method allows to align the graph towards the left (AlignLeft) or the Right (AlignRight). The probabilities according to these four methods over the database from Figure 7b are given in the equation below.

$$\begin{aligned} P(\text{AlignRight}) &= (1 - (1 - p(z_1))(1 - p(z_2))) \cdot \\ & (1 - (1 - p(x_1)(1 - (1 - p(y_1))(1 - p(y_2))))(1 - p(x_2)p(y_3))) \\ P(\text{AlignLeft}) &= (1 - (1 - p(x_1))(1 - p(x_2))) \cdot \\ & (1 - (1 - p(z_1)(1 - (1 - p(y_1))(1 - p(y_3))))(1 - p(x_2)p(y_2))) \\ P(\text{DissociateRight}) &= 1 - (1 - p(x_1)(1 - (1 - \\ & p(y_1)p(z_1))(1 - p(y_2)p(z_2))))(1 - p(x_2)p(y_3)p(z_1)) \\ P(\text{DissociateLeft}) &= 1 - (1 - p(z_1)(1 - (1 - \\ & p(x_1)p(y_1))(1 - p(x_2)p(y_3))))(1 - p(x_1)p(y_2)p(z_2)) \end{aligned}$$

It is straightforward to check that, depending on the particular probability distribution of the variables in the database, none of the four upper is *the best*. For example,  $\text{AlignRight} > \text{DissociateRight}$  for  $P(x_1) = \frac{1}{2}, P(x_2) = \frac{1}{2}, P(y_1) = \frac{1}{4}, P(y_2) = \frac{1}{4}, P(y_3) = \frac{1}{4}, P(z_1) = \frac{1}{2}, P(z_2) = \frac{1}{2}$ , but  $\text{AlignRight} < \text{DissociateRight}$  for  $P(x_1) = \frac{511}{512}, P(x_2) = \frac{255}{256}, P(y_1) = \frac{7}{8}, P(y_2) = \frac{1}{4}, P(y_3) = \frac{63}{64}, P(z_1) = \frac{1}{2}, P(z_2) = \frac{1}{2}$ .  $\square$



**Figure 6: Graphical illustration for the proof part 2 of Lemma 7**

R	A	S	A	B	T	B
$x_1$	1	$y_1$	1	1	$z_1$	1
$x_2$	1	$y_2$	1	1	$z_2$	1

(a) Example 18

R	A	S	A	B	T	B
$x_1$	1	$y_1$	1	1	$z_1$	1
$x_2$	2	$y_2$	1	2	$z_2$	2
		$y_3$	2	1		

(b) Example 19

**Figure 7: Tuple-independent databases  $(R, S, T)$  used in Examples 18 and 19.**