

# The Containment Problem for *Real* Conjunctive Queries with Inequalities

T.S. Jayram  
IBM Almaden

jayram@almaden.ibm.com

Phokion G. Kolaitis\*  
IBM Almaden

kolaitis@almaden.ibm.com

Erik Vee  
IBM Almaden

vee@almaden.ibm.com

## ABSTRACT

Query containment is a fundamental algorithmic problem in database query processing and optimization. Under set semantics, the query-containment problem for conjunctive queries has long been known to be NP-complete. In real database systems, however, queries are usually evaluated under bag semantics, not set semantics. In particular, SQL queries are evaluated under bag semantics and return multisets as answers, since duplicates are not eliminated unless explicitly requested. The exact complexity of the query-containment problem for conjunctive queries under bag semantics has been an open problem for more than a decade; in fact, it is not even known whether this problem is decidable.

Here, we investigate, under bag semantics, the query-containment problem for conjunctive queries with inequalities. It has been previously shown that, under set semantics, this problem is complete for the second level of the polynomial hierarchy. Our main result asserts that, under bag semantics, the query-containment problem for conjunctive queries with inequalities is undecidable. Actually, we establish the stronger result that this problem is undecidable even if the following two restrictions hold at the same time: (1) the queries use just a single binary relation; and (2) the total number of inequalities is bounded by a certain fixed value. Moreover, the same undecidability results hold under bag-set semantics.

## Categories and Subject Descriptors

H.2.4 [Database Management]: Systems—*query processing, relational databases*; H.2.3 [Database Management]: Languages—*query languages*; F.2.2 [Analysis of Algorithms & Problem Complexity]: Nonnumerical Algorithms and Problems—*computations on discrete structures*

## General Terms

Algorithms, Management, Languages, Theory

## Keywords

conjunctive queries, query containment, inequalities, bag semantics, bag-set semantics, undecidability

\*On leave from UC Santa Cruz

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.

PODS'06, June 26–28, 2006, Chicago, Illinois, USA.  
Copyright 2006 ACM 1-59593-318-2/06/0006 ...\$5.00.

## 1. Introduction and Summary of Results

Query containment is regarded as a fundamental algorithmic problem in database query processing and optimization. This problem asks: given two queries  $Q$  and  $Q'$ , is it true that  $Q(D) \subseteq Q'(D)$ , for every database  $D$ ? Over the years, researchers have investigated in depth the query-containment problem for several different classes of frequently asked database queries. The class of conjunctive queries (that is, select-project-join queries) is arguably the most prominent among these classes. Chandra and Merlin [3] showed that the query-containment problem for conjunctive queries is NP-complete. After this, researchers investigated the worst-case complexity of the query-containment problem for broader classes of database queries that naturally subsume conjunctive queries. Specifically, Sagiv and Yannakakis [13] showed that the query-containment problem for unions of conjunctive queries is  $\Pi_2^P$ -complete, where  $\Pi_2^P$  is the second level of the polynomial hierarchy (NP is the first level of the polynomial hierarchy - see [11]). Klug [7] studied conjunctive queries with comparison predicates  $\neq$ ,  $<$ , and  $\leq$ . He showed that the query-containment problem for conjunctive queries with comparison predicates is in  $\Pi_2^P$ , and conjectured that this upper bound is tight. Klug's conjecture was subsequently confirmed by van der Meyden [14], who proved that the query-containment problem for conjunctive queries with comparison predicates is  $\Pi_2^P$ -complete. As a matter of fact, van der Meyden showed that even the query-containment problem for conjunctive queries with just inequalities ( $\neq$ ) as the only comparison predicate is  $\Pi_2^P$ -complete; this result was further refined in [8].

All aforementioned complexity-theoretic results were obtained under the assumption that queries are evaluated under *set* semantics. This means that the database relations given as inputs to queries are sets (i.e., no duplicate tuples are allowed) and that queries return sets as answers. In real database systems, however, queries are usually evaluated under *bag* semantics, not set semantics: input database relations may be bags (multisets), and queries may return bags as answers. In particular, SQL queries are evaluated under bag semantics, since duplicate tuples are not eliminated unless explicitly specified in the syntax using the `SELECT DISTINCT` construct. In addition to faster response, the reason for not eliminating duplicate tuples in SQL is that the values of aggregate operators, such as `AVG` and `COUNT`, depend on the multiplicities of the tuples in the database relations. In a paper titled "Optimization of *Real* Conjunctive Queries" [4], Chaudhuri and Vardi drew attention to this discrepancy between database theory and practice, and raised the question of whether the known complexity results about conjunctive queries carry over from set semantics to bag semantics.

Chaudhuri and Vardi [4] discovered that, under bag semantics, the query-containment problem for conjunctive queries is  $\Pi_2^P$ -hard, which implies that, in all likelihood, the change from set semantics to bag semantics is accompanied by a jump in complexity. More-

over, they found that the same hardness result holds under *bag-set* semantics, the variant of bag semantics in which the input database relations are sets, but the queries return bags as answers. Chaudhuri and Vardi, however, were not able to pinpoint the exact complexity of this problem. As a matter of fact, even though more than a decade has passed since the publication of [4], it is still not known whether, under bag semantics, the query-containment problem for conjunctive queries is decidable. On the other hand, Ioannidis and Ramakrishnan [6] showed that, under bag semantics, the query-containment problem for unions of conjunctive queries is undecidable.

In this paper, we study, under bag semantics, the query-containment problem for conjunctive queries with inequalities ( $\neq$ ). Our study is motivated by two considerations: first, this problem is of interest *in its own right*, as conjunctive queries with inequalities form the most natural extension of conjunctive queries with comparison predicates; second, tools developed in the course of this study may turn out to be of use in attacking, under bag semantics, the containment problem for conjunctive queries. Our main result is that, under bag semantics, the query-containment problem for conjunctive queries with inequalities is *undecidable*. Note that, in general, the inputs to this problem are conjunctive queries over arbitrary relational schemas and with an arbitrary number of inequalities. Thus, it is natural to ask whether solvability results can be obtained by considering restrictions in which the relations in the schemas have bounded arities or the total number of inequalities in the queries is bounded. We establish the stronger result that, under bag semantics, the query-containment problem for conjunctive queries with inequalities is undecidable, even if the following two restrictions hold at the same time: (1) the queries use just a single binary relation; and (2) the total number of inequalities is bounded by a certain fixed value.

To prove these results, we first show that, under bag semantics, the query-containment problem for conjunctive queries with inequalities is polynomial-time equivalent to the same problem under bag-set semantics; moreover, we exhibit polynomial-time reductions that do not increase the total number of inequalities in the queries. We also show that, under either semantics, the problem is polynomial-time equivalent to the restriction in which the queries use just a single binary relation. After this, the crucial undecidability result is established by showing that Hilbert’s Tenth Problem has a recursive reduction to the query-containment problem for conjunctive queries with inequalities and under bag-set semantics. This reduction is carried out in two stages. In the first stage, we identify a class of special databases, which we call *polynomial encoders*, and construct a family of conjunctive queries (without inequalities) so that the evaluation of polynomials can be simulated by evaluating these queries on polynomial encoders under bag-set semantics. This makes it possible to recursively reduce Hilbert’s Tenth Problem to the query-containment problem for conjunctive queries (without inequalities) and under bag-set semantics, provided the input databases are restricted to be polynomial encoders. In the second stage, we amend appropriately the queries used in the first stage (in particular, we add inequalities), so that the containment holds over arbitrary databases, not just polynomial encoders.

It should be pointed out that Hilbert’s Tenth Problem was also used by Ioannidis and Ramakrishnan [6] in proving that, under bag semantics, the containment problem for unions of conjunctive queries is undecidable. Their reduction, however, is much simpler than ours, as the union operation can be used to easily simulate the addition operation on monomials. In the absence of the union operation, we have to develop rather elaborate combinatorial machinery to simulate the evaluation of polynomials; moreover, our simulation makes an essential use of the presence of inequalities in the queries. It should also be pointed out that in a series of papers, including [1, 2, 12], a group

of researchers studied, under bag semantics, the query-containment problems for conjunctive queries with and without comparison predicates, but they did not settle the decidability question for these problems. The main technical assertion in these papers is a necessary and sufficient condition for query containment under bag semantics, a condition that involves a family of *canonical* databases. Unfortunately, there is a subtle counting error in the proof of the claimed necessary and sufficient condition. In fact, a counterexample to this condition is reported in [16].

Finally, it is interesting to note the similarities and differences between query containment and query equivalence, another fundamental problem in database query optimization. Clearly, query equivalence is always reducible to query containment. Consider the case of conjunctive queries (without inequalities). Under set semantics, both query equivalence and query containment are NP-complete [3]. However, the situation is, in all likelihood, different under bag semantics, since query equivalence has the same complexity as GRAPH ISOMORPHISM (hence, is in NP), while query containment is  $\Pi_2^P$ -hard [4]. In the case of conjunctive queries with inequalities, Nutt, Sagiv, and Shurin [10] showed that query equivalence under bag semantics is in PSPACE (see also [5, 15]). Thus, our undecidability result for query containment under bag semantics shows a provable dramatic difference in complexity between the two problems.

## 2. Basic Concepts and Notation

A *bag* or *multiset* is a collection of objects each of which occurs one or more times in the collection. A *relation*  $R$  of arity  $k$  is a bag of  $k$ -tuples whose elements belong to some underlying fixed domain. If  $R$  is a  $k$ -ary relation and  $(A_1, \dots, A_k)$  is a  $k$ -tuple, then we write  $|R(A_1, \dots, A_k)|$  to denote the *multiplicity* (number of occurrences) of  $(A_1, \dots, A_k)$  in the bag  $R$ . In particular,  $|R(A_1, \dots, A_k)| = 0$  means that  $(A_1, \dots, A_k)$  does not occur in the bag  $R$ .

A *database schema* (or, simply, *schema*) is a set  $\mathcal{S}$  of distinct relation symbols  $\{R_1, \dots, R_t\}$  of fixed arities  $k_1, \dots, k_t$ . A *database instance* (or, simply, a *database*)  $D$  for the schema  $\mathcal{S}$  is a set of relations, also denoted by  $R_1, \dots, R_t$ , of arities  $k_1, \dots, k_t$ . If  $R_i$  is one of the relations of a database  $D$  and  $(A_1, \dots, A_{k_i})$  is a  $k_i$ -tuple, we will write  $|R(A_1, \dots, A_{k_i})|_D$ , to emphasize that  $R_i$  is a relation in the database  $D$ , unless the database is understood from the context in which case we will write  $|R(A_1, \dots, A_{k_i})|$ . If  $|R(A_1, \dots, A_{k_i})|_D > 0$ , we say that  $R(A_1, \dots, A_{k_i})$  is a *fact* of  $D$ .

In what follows, we will upper-case letters to denote elements of the domain, and lower-case letters to denote variables. As usual,  $\neq$  denotes the *built-in* inequality relation with the standard interpretation over any domain.

**DEFINITION 1 (CONJUNCTIVE QUERY WITH INEQUALITIES).** Let  $\mathcal{S} = \{R_1, \dots, R_t\}$  be a schema and  $n$  a non-negative integer. A  $n$ -ary *conjunctive query with inequalities* over  $\mathcal{S}$  is a rule of the form:

$$Q(x_1, \dots, x_n) :- T_1(\vec{z}_1), \dots, T_p(\vec{z}_p), \quad (1)$$

where  $T_i$  is one of the relation symbols in  $\mathcal{S}$  or the  $\neq$  symbol, and each  $\vec{z}_i$  is a tuple of variables in set  $Z = \{x_1, \dots, x_n, y_1, \dots, y_m\}$ .

Each expression  $T_i(\vec{z}_i)$  is a *subgoal* of  $Q$ ; the list of subgoals is the *body* of  $Q$ , while  $Q(x_1, \dots, x_n)$  is the head of  $Q$ . Note that the same subgoal may occur more than once in the body of  $Q$ . We only consider *safe* queries, that is, each variable  $x_j$  in the head must appear in at least one of the subgoals of  $Q$  that involve a relation symbol in  $\mathcal{S}$ . The variables  $x_1, \dots, x_n$  are called *free* and the variables  $y_1, \dots, y_m$  are called *bound*. When  $n = 0$  (i.e.,  $Q$  has no free variables), then we say that  $Q$  is a *bound* (or, *Boolean*) query.  $\square$

An *assignment mapping* (or, simply, an *assignment*) is a mapping  $\tau$  from the set of free and bound variables of a query  $Q$  to the underlying domain of elements. If  $\vec{z} = (z_1, \dots, z_s)$  is a tuple of variables, we will write  $\tau(\vec{z})$  to denote the tuple  $(\tau(z_1), \dots, \tau(z_s))$ . Assignments are used to define the semantics of queries, that is, to define what it means to evaluate a query on a database. Under *set semantics*, relations are assumed to be sets, and evaluating a query essentially amounts to determining whether there is an assignment that satisfies every subgoal of the query. Here, we consider two different kinds of semantics, namely *bag semantics* and *bag-set semantics*. The difference from set semantics is that, instead of asking whether a satisfying assignment exists, we now want to know how many different satisfying assignments are there. Informally, in both bag semantics and bag-set semantics, the query is evaluated over a database algebraically by treating the conjunction as a product and then “marginalizing out” the bound variables  $y_j$  by summing over all choices for the  $y_j$ 's. The difference between bag and bag-set semantics is whether the multiplicities of the tuples in the database relations are taken into account while evaluating the product.

Formally, let  $D$  be a database over a schema  $\mathcal{S} = \{R_1, \dots, R_t\}$  and let  $Q$  be a conjunctive query with inequalities as in (1). The *result of evaluating  $Q$  on  $D$  under bag semantics* is the  $n$ -ary relation  $\text{EVAL}_B(Q, D)$  defined as follows: for every  $n$ -tuple  $(A_1, \dots, A_n)$  of elements from the domain,

$$|\text{EVAL}_B(Q, D)(A_1, \dots, A_n)| = \sum_{\tau} \prod_{i=1}^p |T_i(\tau(\vec{z}_i))|_D,$$

where  $\tau$  ranges over all assignments such that  $\tau(x_i) = A_i$ ,  $1 \leq i \leq n$ . For a bound query, the result  $\text{EVAL}(Q, D)$  is just a single non-negative integer. For bag-set semantics, we remove all duplicates of facts in the instance  $D$  so that each relation is a set. Let  $\tilde{D}$  denote the database obtained this way. The *result of evaluating  $Q$  under bag-set semantics* is the  $n$ -ary relation  $\text{EVAL}_{BS}(Q, D) = \text{EVAL}_B(Q, \tilde{D})$ . When the subscript is omitted,  $\text{EVAL}(Q, D)$  simply refers to query evaluation under bag semantics.

**DEFINITION 2 (QUERY CONTAINMENT).** Let  $Q$  and  $Q'$  be two  $n$ -ary conjunctive queries with inequalities over some schema  $\mathcal{S}$ . We say that  $Q$  is *contained in  $Q'$  under bag semantics*, denoted by  $Q \subseteq_B Q'$  if for every database  $D$  over  $\mathcal{S}$  and for every  $n$ -tuple  $(A_1, \dots, A_n)$  of elements from the domain, we have that

$$|\text{EVAL}_B(Q, D)(A_1, \dots, A_n)| \leq |\text{EVAL}_B(Q', D)(A_1, \dots, A_n)|.$$

In other words, the multiplicities of facts of tuples corresponding to  $Q$  are no more than those corresponding to  $Q'$ . The concept  $Q$  is *contained in  $Q'$  under bag-set semantics*, denoted by  $Q \subseteq_{BS} Q'$ , is defined in an analogous way using  $\text{EVAL}_{BS}(Q, D)$  and  $\text{EVAL}_{BS}(Q', D)$ .

The main problem that we consider in this paper is query containment under bag semantics and under bag-set semantics. In fact, we parameterize this problem into a family of sub-problems using the characteristics of the database schema and the number of inequalities in the queries as parameters.

**DEFINITION 3 (QUERY CONTAINMENT PROBLEM).** Let  $k$  be a non-negative integer and  $m, d$  be two positive integers. We write  $\text{ConQC}_B(k, m, d)$  to denote the following decision problem: given a schema  $\mathcal{S}$  having at most  $m$  relations, each of arity at most  $d$ , and given two conjunctive queries  $Q$  and  $Q'$  each of which has at most  $k$  inequalities, is  $Q \subseteq_B Q'$ ?

We write  $\text{ConQC}_{BS}(k, m, d)$  to denote the same problem under bag-set semantics.

We also allow (some of) these parameters to be unbounded; we denote this by setting the appropriate parameter(s) to  $\infty$ . For example,

$\text{ConQC}_B(\infty, \infty, \infty)$  denotes the conjunctive query problem under bag semantics where there are no restrictions on any one of the three parameters.  $\square$

To simplify our presentation, we will make systematic use of the concept of a *view*. Informally, a view  $V$  is defined by a rule in which other previously defined views can be part of the body of  $V$ . We do not allow views to be defined in terms of themselves, that is, recursive views are prohibited. In effect, views are queries defined by a non-recursive Datalog program with inequalities. A formal self-contained inductive definition of views over a schema  $\mathcal{S}$  is as follows. For the base case, any relation symbol in  $\mathcal{S} \cup \{\neq\}$  with variables as arguments is a view. Inductively, suppose  $V_1, \dots, V_k$  are views that are already defined. Then, a view  $V$  can be defined by any rule of the

$$V(x_1, \dots, x_n) :- T_1(\vec{z}_1), \dots, T_p(\vec{z}_p),$$

where each  $\vec{z}_i$  is a tuple of variables from the set  $Z = \{x_1, \dots, x_n, y_1, \dots, y_m\}$ , and each  $T_i(\vec{z}_i)$  is obtained from one of the views  $V_j$ ,  $1 \leq j \leq k$ , by replacing the free variables of  $V_j$  by corresponding variables from  $\vec{z}_i$ .

It is clear that every view can be *unfolded* to a conjunctive query with inequalities. This is done by tracing the inductive definition of views and replacing bound variables in such a way that the subgoals of the views do not have common (“local” to the subgoal) bound variables. For example, consider the views

$$U(u_1, u_2) :- A(u_1, u_2, w),$$

$$W(u_1, u_2) :- A(u_1, u_2, w), A(u_1, u_2, w), B(u_1, w), u_1 \neq w$$

$$V(x_1, x_2) :- U(x_1, y), U(x_1, y), W(x_2, y)$$

on the schema  $\{A, B\}$ . Unfolding  $V$ , yields the following conjunctive query with inequalities:

$$Q(x_1, x_2) :- A(x_1, y, w_1), A(x_1, y, w_2), A(x_2, y, w_3), \\ A(x_2, y, w_3), B(x_2, w_3), x_2 \neq w_3.$$

Given view  $V$  and database  $D$ , we can define  $\text{EVAL}_B(V, D)$ , the result of evaluating  $V$  on  $D$  under bag semantics, by first unfolding  $V$  to a query  $Q$  and then evaluating  $Q$  on  $D$  under bag semantics. In a similar way, we can define the result  $\text{EVAL}_{BS}(V, D)$  of evaluating  $V$  on  $D$  under bag-set semantics.

By induction on the definition of views, it is not hard to show that if we have a view

$$V(x_1, \dots, x_n) :- T_1(\vec{z}_1), \dots, T_p(\vec{z}_p),$$

then for every database  $D$  and every  $n$ -tuple  $(A_1, \dots, A_n)$  of elements from the domain, we have that

$$|\text{EVAL}_B(V, D)(A_1, \dots, A_n)| = \sum_{\tau} \prod_{i=1}^p |\text{EVAL}_B(T_i, D)(\tau(\vec{z}_i))|,$$

where  $\tau$  ranges over all assignments such that  $\tau(x_i) = A_i$ , for  $1 \leq i \leq n$ . Moreover,  $|\text{EVAL}_{BS}(V, D)(A_1, \dots, A_n)|$  may be evaluated similarly. In the sequel, we will use these two facts repeatedly.

We conclude this section with two more concepts. Let  $Q$  be a conjunctive query with inequalities over a schema  $\mathcal{S}$  defined by the rule  $Q(x_1, \dots, x_n) :- T_1(\vec{z}_1), \dots, T_p(\vec{z}_p)$ , where each  $T_i \in \mathcal{S} \cup \{\neq\}$ . For an assignment  $\tau$  that satisfies the inequalities in  $Q$ , we define  $\tau(Q)$  to be the database over  $\mathcal{S}$  whose relations are the bags of facts  $T_i(\tau(\vec{z}_i))$ , for every  $T_i \in \mathcal{S}$ . Finally, we say that  $\tau$  is a *homomorphism* from  $Q$  to a database  $D$ , if every fact in  $\tau(Q)$  appears in  $D$ .

### 3. Bag vs. bag-set semantics

Depending on whether we use bag semantics or bag-set semantics, evaluating queries can give different answers. So it is natural to wonder whether the complexity of the conjunctive query containment problem (with or without inequalities) under bag semantics is different than its complexity under bag-set semantics.

The following theorem shows that in fact, the two problems are polytime reducible to each other, so their complexities are essentially the same. In fact, the theorem shows something stronger: The conjunctive query containment problem (with or without inequalities), under both bag-set and bag semantics, in which we have no restrictions on the number of relations or their arity, is polytime reducible to the conjunctive query containment problem (with or without inequalities, respectively) in which queries are restricted to use a single binary relation. Previously, Chadhuri and Vardi [4] noted that  $\text{ConQC}_B(0, \infty, \infty)$  is polytime reducible to  $\text{ConQC}_{BS}(0, \infty, \infty)$ . Our claim is stronger, not just because it applies to  $\neq$ -constraints, but because it also provides a reduction that reduces the queries to have a single binary relation.

**THEOREM 1.** *For any  $k \geq 0, m \geq 1$ , and  $d \geq 1$ , we have the following:*

$$\begin{aligned} \text{ConQC}_{BS}(k, m, d) &\leq_P \text{ConQC}_B(2k, m, d) \\ \text{ConQC}_B(k, \infty, \infty) &\leq_P \text{ConQC}_{BS}(k, 1, 2) \end{aligned}$$

*The reductions hold even if we allow any combination of  $k, m, d$  to be  $\infty$ . (If  $k = \infty$ , we interpret  $2k = \infty$  as well.)*

**PROOF SKETCH.** We first give the reduction for the first item. Suppose we are given an instance of  $\text{ConQC}_{BS}(k, r, m)$  in the form of queries  $Q_1, Q_2$ , which may have up to  $k \neq$ -constraints each. (In the case that  $k = \infty$ , we allow them to have arbitrarily many  $\neq$ -constraints.) Since we are considering bag-set semantics, we may assume without loss of generality that every subgoal of  $Q_1$  and  $Q_2$  appears with multiplicity one. Let  $Q'_1$  be query  $Q_1$  modified so that each of its subgoals appears with multiplicity exactly two. Define

$$Q_1^B : - Q_1 \wedge Q_1 \quad \text{and} \quad Q_2^B : - Q_2 \wedge Q_2$$

Note that queries  $Q_1^B$  and  $Q_2^B$  both have at most  $2k \neq$ -constraints. We claim that  $\text{EVAL}_{BS}(Q_1, D) \leq \text{EVAL}_{BS}(Q_2, D)$  for all  $D$  if and only if  $\text{EVAL}_B(Q_1^B, D') \leq \text{EVAL}_B(Q_2^B, D')$  for all  $D'$ .

The key to this reduction is that if all of the facts in database  $D$  have multiplicity one, then  $\text{EVAL}_B(Q_1, D) = \text{EVAL}_B(Q_1^B, D)$ . On the other hand, if  $D$  has facts with multiplicity greater than one, then query  $Q_2^B$  will benefit more than  $Q_1^B$ . Specifically, suppose that

$$\text{EVAL}_{BS}(Q_1, D) > \text{EVAL}_{BS}(Q_2, D)$$

for some database  $D$ . Of course, this implies by definition that

$$\text{EVAL}_B(Q_1, \tilde{D}) > \text{EVAL}_B(Q_2, \tilde{D})$$

(Recall that  $\tilde{D}$  is the database obtained from  $D$  by setting the multiplicity of all facts in  $D$  to one.) Hence,

$$\begin{aligned} \text{EVAL}_B(Q_1^B, \tilde{D}) &= \text{EVAL}_B(Q_1, \tilde{D}) \cdot \text{EVAL}_B(Q_1, \tilde{D}) \\ &> \text{EVAL}_B(Q_2, \tilde{D}) \cdot \text{EVAL}_B(Q_1^B, \tilde{D}) \\ &= \text{EVAL}_B(Q_2^B, \tilde{D}) \end{aligned}$$

On the other hand, suppose

$$\text{EVAL}_{BS}(Q_1, D) \leq \text{EVAL}_{BS}(Q_2, D)$$

for all  $D$ . For query  $Q$  given by  $Q : -A_1(\tilde{z}_1), \dots, A_p(\tilde{z}_p)$ , define the *weight of homomorphism*  $\phi$  mapping from  $Q$  to  $D$  to be the value

$\prod_{i \in [p]} |A_i(\phi(\tilde{z}_i))|_D$ . Then for any  $w > 0$ , let  $k_w$  be the number of homomorphisms from  $Q_1$  to  $D$  that have weight  $w$ . Clearly,

$$\begin{aligned} \text{EVAL}_{BS}(Q_1, D) &= \sum_w k_w \\ \text{EVAL}_B(Q_1, D) &= \sum_w w k_w \\ \text{EVAL}_B(Q'_1, D) &= \sum_w w^2 k_w \end{aligned}$$

Noting that  $\text{EVAL}_B(Q_2, D) \geq \text{EVAL}_{BS}(Q_2, D) \geq \text{EVAL}_{BS}(Q_1, D)$ , we see

$$\begin{aligned} \text{EVAL}_B(Q_2 \wedge Q'_1, D) &= \text{EVAL}_B(Q_2, D) \cdot \text{EVAL}_B(Q'_1, D) \\ &\geq \text{EVAL}_{BS}(Q_1, D) \cdot \text{EVAL}_B(Q'_1, D) \\ &= \left( \sum_w k_w \right) \left( \sum_w w^2 k_w \right) \\ &\geq \left( \sum_w w \cdot k_w \right)^2 \text{ by Cauchy-Schwartz} \\ &= \text{EVAL}_B(Q_1 \wedge Q_1, D) \end{aligned}$$

That is,  $\text{EVAL}_B(Q_1^B, D) \leq \text{EVAL}_B(Q_2^B, D)$ . So the reduction works as claimed. This completes the proof for the first item.

We now give the reduction for the second item. Suppose we are given  $Q_1$  and  $Q_2$ , each with at most  $k$  inequalities. Note that  $Q_1$  and  $Q_2$  may have subgoals with multiplicity greater than one. Further, suppose that  $Q_1, Q_2$  are defined over  $r$  relations,  $R_1, \dots, R_r$  with arities  $k_1, \dots, k_r$  respectively, and that the queries use variables  $\mathcal{V} = \{v_1, \dots, v_n\}$ . Without loss of generality, we assume  $k_i \geq 2$  for all  $i \in [r]$ .

Given such a query,  $Q$ , suppose we may write  $Q$  as

$$Q : - \bigwedge_{i=1}^r \bigwedge_{j=1}^{k_i} R_i(\tilde{z}_{ij})$$

where each  $\tilde{z}_{ij}$  belongs to  $\mathcal{V}$ . Again, note that the above expression may have repeated subgoals. We define a corresponding query, denoted  $Q^{BS}$ , which uses a single binary relation  $R$ , in terms of views  $\text{View}_1, \dots, \text{View}_r$ :

$$Q^{BS} : - \bigwedge_{i=1}^r \bigwedge_{j=1}^{k_i} \text{View}_i(\tilde{z}_{ij})$$

where for each  $i \in [r]$ , we define  $\text{View}_i$  using binary relation  $R$  as follows.

$$\begin{aligned} \text{Path}_\ell(s, t) &: - R(s, v_1), R(v_1, v_2), \dots, R(v_{\ell-1}, t) \\ \text{View}_i(u_1, \dots, u_{k_i}) &: - \text{Path}_{r+2}(s_1, u_1), \dots, \text{Path}_{r+2}(s_{k_i}, u_{k_i}) \\ &\quad \text{Path}_{i+1}(s_1, t_1), \dots, \text{Path}_{i+1}(s_{k_i}, t_{k_i}) \\ &\quad R(s_2, t_1), \dots, R(s_{k_i}, t_{k_i-1}) \end{aligned}$$

Notice that  $Q^{BS}$  has no repeated subgoals. We claim that

$$\text{EVAL}_B(Q_1, D) \leq \text{EVAL}_B(Q_2, D)$$

for all  $D$  if and only if

$$\text{EVAL}_{BS}(Q_1^{BS}, D') \leq \text{EVAL}_{BS}(Q_2^{BS}, D')$$

for all  $D'$ . The details of the proof of this claim appear in the full version of this paper. Here, we only give a rough intuition.

First of all, observe that each subgoal of the original queries, e.g.  $R_i(\bar{z}_{ij})$ , is replaced by an analogous view, e.g.  $\text{View}_i(\bar{z}_{ij})$ . Given a database  $D$ , we can perform a similar operation: For each fact in  $D$ , say  $R_i(Z_1, \dots, Z_{k_i})$ , we replace it with a set of facts corresponding to the *canonical database* associated with  $\text{View}_i(Z_1, \dots, Z_{k_i})$ . (Specifically, for each subgoal  $R(x, y)$  in the unfolding of  $\text{View}_i$ , add fact  $R(X, Y)$  to the database, where  $X$  and  $Y$  are constants. Furthermore, we require that each bound variable be mapped to its own unique constant, and that each free variable  $u_\ell$  is mapped to the constant  $Z_\ell$ .) Call this new database  $D^{\text{BS}}$ . It is not hard to see that every homomorphism from  $\mathbb{Q}$  to  $D$  has a corresponding homomorphism from  $\mathbb{Q}^{\text{BS}}$  to  $D^{\text{BS}}$ . With some more work, it is possible to show that the converse is true as well: for every homomorphism from  $\mathbb{Q}^{\text{BS}}$  to  $D^{\text{BS}}$ , there is a corresponding homomorphism from  $\mathbb{Q}$  to  $D$ . (The exact structure of  $\text{View}_i$  was chosen carefully to ensure this.) That is,

$$\text{EVAL}_B(\mathbb{Q}, D) = \text{EVAL}_{\text{BS}}(\mathbb{Q}^{\text{BS}}, D^{\text{BS}}) \quad (2)$$

On the other hand, suppose we are given a database  $D$ , and we wish to construct a database, which we denote  $D^{\text{B}}$ , with the property that

$$\text{EVAL}_{\text{BS}}(\mathbb{Q}^{\text{BS}}, D) = \text{EVAL}_B(\mathbb{Q}, D^{\text{B}}) \quad (3)$$

In this case, for each tuple  $X_1, \dots, X_{k_i}$  consisting of constants from the domain of  $D$ , we add fact  $R_i(X_1, \dots, X_{k_i})$  to database  $D^{\text{B}}$  with multiplicity  $\text{EVAL}_{\text{BS}}(\text{View}_i(X_1, \dots, X_{k_i}), D)$ . (Strictly speaking, we do not add the fact if the evaluation has value 0.) It is not difficult to see that produces a database  $D^{\text{B}}$  satisfying equation (3).

So, if there is a  $D$  such that  $\text{EVAL}_B(\mathbb{Q}_1, D) > \text{EVAL}_B(\mathbb{Q}_2, D)$ , then  $\text{EVAL}_{\text{BS}}(\mathbb{Q}_1^{\text{BS}}, D^{\text{BS}}) > \text{EVAL}_{\text{BS}}(\mathbb{Q}_2^{\text{BS}}, D^{\text{BS}})$ , by equation (2). Conversely, if  $\text{EVAL}_{\text{BS}}(\mathbb{Q}_1^{\text{BS}}, D) > \text{EVAL}_{\text{BS}}(\mathbb{Q}_2^{\text{BS}}, D)$  for some  $D$ , then  $\text{EVAL}_B(\mathbb{Q}_1, D^{\text{B}}) > \text{EVAL}_B(\mathbb{Q}_2, D^{\text{B}})$ , by equation (3). The theorem thus follows.

□

COROLLARY 2. For all  $k \geq 0$ ,

$$\begin{aligned} \text{ConQC}_B(k, \infty, \infty) &\leq_P \text{ConQC}_{\text{BS}}(k, 1, 2) \\ &\leq_P \text{ConQC}_B(2k, 1, 2) \\ \text{ConQC}_{\text{BS}}(k, \infty, \infty) &\leq_P \text{ConQC}_{\text{BS}}(2k, 1, 2) \\ &\leq_P \text{ConQC}_B(4k, 1, 2) \end{aligned}$$

In the next section, we will construct queries using many relations, in order to show that  $\text{ConQC}_{\text{BS}}(k, \infty, \infty)$  is undecidable for some bounded  $k$ . Corollary 2 implies that both  $\text{ConQC}_{\text{BS}}(2k, 1, 2)$  and  $\text{ConQC}_{\text{BS}}(4k, 1, 2)$  are undecidable as well.

## 4. Undecidability of conjunctive query containment

In this section, we show that the conjunctive query containment problem with inequalities is undecidable under bag-set semantics by exhibiting a reduction from Hilbert's Tenth Problem. Our reduction will use homogeneous polynomials of degree  $d$  (i.e., each term is a product of  $d$  variables, not necessarily distinct) with non-negative coefficients. We will use the following version of Hilbert's Tenth Problem in our reduction, which can be obtained by combining the results of Matiyasevich [9] with some additional arguments (details will appear in the full version of this paper).

THEOREM 3. Let  $P_1(x_1, \dots, x_n)$ , and  $P_2(x_1, \dots, x_n)$  be homogeneous polynomials of degree  $d$  each having the same set of terms with positive integer coefficients. Further, assume that  $x_1$  divides both  $P_1(\bar{x})$  and  $P_2(\bar{x})$ . Then it is undecidable to determine whether there are non-negative integers  $x_1, \dots, x_n$  such that

$$P_1(x_1, \dots, x_n) > x_1^d \cdot P_2(x_1, \dots, x_n)$$

This problem is undecidable even if we restrict ourselves to homogeneous polynomials of degree  $d = 5$  and with  $n = 59$  variables. □

Throughout the remainder of the paper, both  $P_1(x_1, \dots, x_n)$  and  $P_2(x_1, \dots, x_n)$  will refer to degree  $d$  homogeneous polynomials each having the same set of  $m$  terms with positive integer coefficients. For  $j \in [m]$ , we associate the  $j$ -th term with an ordered  $d$ -tuple  $\mathcal{T}_j$  in  $[n]^d$ ; since  $x_1$  divides both  $P_1(\bar{x})$  and  $P_2(\bar{x})$ , we further assume that the first entry of  $\mathcal{T}_j$  is 1. Abusing notation slightly, we will also think of  $\mathcal{T}_j$  as a multiset. For example, if  $\mathcal{T}_j = (1, 1, 3, 7)$ , then  $\prod_{i \in \mathcal{T}_j} x_i = x_1^2 x_3 x_7$ . We write

$$P_1(\bar{x}) = \sum_{j=1}^m \alpha_j \prod_{i \in \mathcal{T}_j} x_i \quad \text{and} \quad P_2(\bar{x}) = \sum_{j=1}^m \beta_j \prod_{i \in \mathcal{T}_j} x_i,$$

where  $\alpha_j, \beta_j \geq 0$  for all  $j$ . For the rest of the paper, we will fix the  $m$ -dimensional vectors  $\alpha$  and  $\beta$  corresponding to  $P_1$  and  $P_2$ , respectively.

### 4.1 Proof Overview

Our goal is to show that the undecidable problem on polynomials described in Theorem 3 can be reduced to the problem of query containment under bag semantics. Since the proof of the reduction is fairly involved, we will break it into several steps. We begin in Section 4.2 by describing two queries  $\text{Poly}_1$  and  $\text{Poly}_2$  that do not have any inequalities. Further, we will restrict the evaluations of these queries over a very specific class of database instances, called *polynomial encoders*. For each  $\xi \in \mathbb{N}_0^n$ , we will construct a polynomial encoder, denoted  $D_\xi$ , and show in Theorem 5 that  $\text{EVAL}(\text{Poly}_1, D_\xi) = P_1(\xi)$  and  $\text{EVAL}(\text{Poly}_2, D_\xi) = \xi_1^d P_2(\xi)$ . By Theorem 3, there is no algorithm for the query containment problem if the evaluation is restricted to the class of polynomial encoders.

For the full problem where all instances are considered, we have to work considerably harder, and this is where we use the power of inequalities. In Section 4.3, we will extend the queries  $\text{Poly}_1$  and  $\text{Poly}_2$  to produce queries  $\mathbb{Q}_1$  and  $\mathbb{Q}_2$ , respectively, with inequalities. For technical reasons, we will work with *augmented polynomial encoders*, which are database instances that consist of polynomial encoders augmented by a small set of facts. Now, when we consider databases that are augmented polynomial encoders, the arguments we laid out in Section 4.2 go through as before. But what if the database does not have the desired structure? Then the left-hand query,  $\mathbb{Q}_1$ , could potentially *cheat* by mapping to the database in ways that we did not anticipate. Although we cannot stop this, we can guarantee that for every map of  $\mathbb{Q}_1$  to the database that cheats, there is a corresponding map from  $\mathbb{Q}_2$  to the database. We do this by defining a view `CounterCheating` that will be part of  $\mathbb{Q}_2$ . This view contains many copies of  $\text{Poly}_1$  along with appropriate inequality constraints. For every map from  $\mathbb{Q}_1$  to the database that cheats, there is a corresponding map from one of the copies of  $\text{Poly}_1$  in `CounterCheating` that maps in an identical fashion. This ensures that cheating helps  $\mathbb{Q}_2$  as much as it helps  $\mathbb{Q}_1$ . The definition of `CounterCheating` and the notion of cheating, together with the proofs are described in Section 4.3.

## 4.2 Conjunctive Queries over Polynomial Encoders

The polynomial encoder for  $\xi \in \mathbb{N}_0^n$ , denoted  $D_\xi$ , will use domain elements  $X_1, \dots, X_n$ , corresponding to variables of the polynomial, domain elements  $T_1, \dots, T_m$ , corresponding to terms of the polynomial, as well as auxiliary elements  $T_0$  and  $U_{k,1}, \dots, U_{k,\xi_k}$  for each  $k \in [m]$ . We will describe its structure later in this section.

Recall that  $\alpha$  and  $\beta$  denote the vector of coefficients for  $P_1(\bar{x})$  and  $P_2(\bar{x})$ , respectively. Both  $\text{Poly}_1$  and  $\text{Poly}_2$  will be described in terms of the following views:  $\text{Term}(u_1, \dots, u_d, z_0)$ ,  $\text{Value}(v)$ ,  $\text{Coeff}_\alpha(z_0)$ , and  $\text{Coeff}_\beta(z_0)$ . The following lemma states the several key properties of these views when evaluated on a polynomial encoder. We defer its proof to the next subsection.

LEMMA 4. *Let  $D_\xi$  be a polynomial encoder. Then*

1.  $\text{EVAL}(\text{Term}, D_\xi) = \{(X_{i_1}, \dots, X_{i_d}, T_j) \mid (i_1, \dots, i_d) = T_j\}$ , with each element having multiplicity one.
2.  $\text{EVAL}(\text{Value}, D_\xi) = \{X_1, \dots, X_n\}$ , with element  $X_k$  having multiplicity  $\xi_k$ , for all  $k \in [n]$ .
3.  $\text{EVAL}(\text{Coeff}_\alpha, D_\xi) = \{T_0, \dots, T_j\}$ , where  $T_j$  has multiplicity  $\alpha_j$  for all  $j \in [m]$ , and  $T_0$  has multiplicity one.
4. Likewise,  $\text{EVAL}(\text{Coeff}_\beta, D_\xi) = \{T_0, \dots, T_j\}$ , where  $T_j$  has multiplicity  $\beta_j$  for all  $j \in [m]$ , and  $T_0$  has multiplicity one.  $\square$

We now define  $\text{Poly}_1$  and  $\text{Poly}_2$  in terms of the above views.

$$\text{Poly}_1 := \text{Term}(u_1, \dots, u_d, z_0), \text{Coeff}_\alpha(z_0), \\ \text{Value}(u_1), \dots, \text{Value}(u_d)$$

and

$$\text{Poly}_2 := \text{Term}(u_1, \dots, u_d, z_0), \text{Coeff}_\beta(z_0), \\ \underbrace{\text{Value}(u_1), \dots, \text{Value}(u_1)}_{d \text{ times}}, \\ \text{Value}(u_1), \dots, \text{Value}(u_d)$$

Given Lemma 4, we can prove the key theorem of this section that the evaluation of  $\text{Poly}_1$  and  $\text{Poly}_2$  over polynomial encoders indeed results in the evaluations of polynomials  $P_1(\bar{x})$  and  $P_2(\bar{x})$ , respectively.

THEOREM 5. *Let  $\xi \in \mathbb{N}_0^n$ , and let  $D_\xi$  be a polynomial encoder. Then  $\text{EVAL}(\text{Poly}_1, D_\xi) = P_1(\xi)$  and  $\text{EVAL}(\text{Poly}_2, D_\xi) = \xi_1^d P_2(\xi)$ .*

PROOF. Fix a  $\xi \in \mathbb{N}_0^n$ , and let  $D = D_\xi$ . First, consider the value of  $\text{EVAL}(\text{Poly}_1, D)$ :

$$\sum_{\substack{i_1, \dots, i_d, T_j: \\ (i_1, \dots, i_d) = T_j}} |\text{Term}(X_{i_1}, \dots, X_{i_d}, T_j)| \cdot |\text{Coeff}_\alpha(T_j)| \\ \cdot |\text{Value}(X_{i_1})| \dots |\text{Value}(X_{i_d})|$$

By Lemma 4, the above expression evaluates to  $\sum_j \alpha_j \prod_{i \in T_j} \xi_i = P_1(\xi)$ , as we claimed. Applying Lemma 4, a similar proof shows that  $\text{EVAL}(\text{Poly}_2, D) = \sum_j \beta_j \xi_1^d \prod_{i \in T_j} \xi_i = \xi_1^d P_2(\xi)$ , as desired.  $\square$

The rest of this subsection is devoted to the proof of Lemma 4. We first give the definitions of the views  $\text{Term}$ ,  $\text{Value}$ ,  $\text{Coeff}_\alpha$ , and  $\text{Coeff}_\beta$  and the polynomial encoders.

Let  $R, R_1, \dots, R_d, S_0, \dots, S_m$  be binary relation names, and define the views as follows. Here, we let  $N_k \triangleq \sum_{j=1}^k \alpha_j$  for all

$k \in [m]$ .

$$\text{Term}(\bar{u}, z_0) := R_1(u_1, z_0), \dots, R_d(u_d, z_0) \\ \text{Value}(v) := R(v, v') \\ \text{Coeff}_\alpha(z_0) := S_1(z_0, z_1), \dots, S_1(z_{N_1-1}, z_{N_1}), \\ S_2(z_{N_1}, z_{N_1+1}), \dots, S_2(z_{N_2-1}, z_{N_2}), \\ \vdots \\ S_m(z_{N_{m-1}}, z_{N_{m-1}+1}), \dots \\ \dots, S_m(z_{N_m-1}, z_{N_m}), \\ S_0(z_{N_m}, z_{N_m})$$

We define  $\text{Coeff}_\beta(z_0)$  in analogy with  $\text{Coeff}_\alpha(z_0)$ , where we replace each  $N_k$  with  $M_k \triangleq \sum_{j=1}^k \beta_j$  for all  $k \in [m]$ .

We are now ready to describe the polynomial encoder  $D_\xi$  as the union of two databases  $D^*$  and  $\overline{D}_\xi$  with disjoint relations.  $D^*$  is independent of  $\xi$  and is defined below:

$$R_k = \{(X_i, T_j) \mid i \text{ is the } k\text{-th entry of } T_j\}, \quad \forall k \in [d] \\ S_j = \{(T_j, T_j), (T_{j+1}, T_{j+1}), \dots, (T_m, T_m)\} \\ \cup \{(T_j, T_0)\} \cup \{(T_0, T_0)\}, \quad \forall j \in [m] \\ S_0 = \{(T_0, T_0)\}$$

For each  $\xi \in \mathbb{N}_0^n$ , we define  $\overline{D}_\xi$  as follows:

$$R = \{(X_1, U_{1,1}), \dots, (X_1, U_{1,\xi_1}), \\ \vdots \\ (X_n, U_{n,1}), \dots, (X_n, U_{n,\xi_n})\}$$

Observe that each fact in the polynomial encoder  $D_\xi$  has multiplicity 1. Moreover,  $D_\xi = D^*$ , when  $\xi = 0$ .

PROOF OF LEMMA 4. For part (1), notice that  $R_k(X_i, T_j)$  is a fact of  $D_\xi$  if and only if  $i$  is the  $k$ th entry of  $T_j$ , for each  $i \in [n]$ ,  $j \in [m]$ ,  $k \in [d]$ . Hence,  $\text{EVAL}(\text{Term}, D_\xi)$  returns the tuple  $(X_{i_1}, \dots, X_{i_d}, T_j)$  if and only if  $(i_1, \dots, i_d) = T_j$ . Since all variables in this view are free, each tuple has multiplicity one, as we claimed.

For part (2), since  $\text{Value}(v) := R(v, v')$ ,  $v$  must map to  $X_k$  for some  $k \in [n]$ . Then,  $v'$  can map to exactly one of the  $\xi_k$  elements  $U_{k,1}, \dots, U_{k,\xi_k}$ . Therefore, the query returns the set  $\{X_1, \dots, X_n\}$  with element  $X_k$  having multiplicity precisely  $\xi_k$ .

The proof for parts (3) and (4), are similar so we prove only part (3). For every  $j \in [m]$ , the query  $\text{Coeff}_\alpha$  has a subgoal  $S_j(z_{i-1}, z_i)$  for every  $i$  such that  $N_{j-1} < i \leq N_j$ . Since only facts in  $D_\xi$  where  $T_0$  is the first component are  $S_j(T_0, T_0)$ , for all  $j \in [m]$ , we have by induction that (\*) if  $z_i$  is mapped to  $T_0$ , then  $z_k$  must be mapped to  $T_0$  for all  $k > i$ . Similarly, the only facts in  $D_\xi$  where  $T_j$  is the second component are  $S_{j'}(T_j, T_j)$ , for all  $j' \leq j$ . Therefore, (\*\*) if  $z_i$  is mapped to  $T_j$  for some  $j$ , then  $i \leq N_j$ , and  $z_k$  is mapped to  $T_j$  for all  $k < i$ .

In  $D_\xi$ , the only choices for  $z_0$  are  $T_0, T_1, \dots, T_m$ . Consider the case when  $z_0$  maps to  $T_0$ . By (\*), every  $z_k$  maps to  $T_0$ , and this satisfies all the subgoals so the multiplicity of  $T_0$  is 1.

Next, consider the case when  $z_0$  is mapped to  $T_j$  for some  $j \in [m]$ . To satisfy the subgoal  $S_0$ ,  $z_N$  must map to  $T_0$ . In the database  $S_j(T_j, T_0)$  is the only fact in which  $T_j$  is paired with any other element. Together with (\*) and (\*\*), it follows that the only maps that satisfy the query are such that for some  $i$  where  $N_{j-1} < i \leq N_j$ ,  $z_k$  is mapped to  $T_j$  for all  $k < i$ , and  $z_k$  is mapped to  $T_j$  for all  $k \geq i$ . Each of these  $N_j - N_{j-1} = \alpha_j$  choices of  $i$  result in  $z_0$  being mapped to  $T_j$ , so the multiplicity of  $T_j$  is  $\alpha_j$ , as desired.  $\square$

### 4.3 The Full Construction

We start by formally defining an *augmented polynomial encoder*. Let  $D^{\text{sink}}$  be the database on universe  $\{C_0, C_1\}$  containing

$$\begin{aligned} R &= \{(C_0, C_0)\} \\ R_i &= \{(C_0, C_0), (C_1, C_0)\}, \quad \forall i \in [d] \\ S_j &= \{(C_0, C_0), (C_0, C_1)\}, \quad \forall j \in [m] \\ S_0 &= \{(C_0, C_0)\} \end{aligned}$$

The class of augmented polynomial encoders are simply  $D_\xi^{\text{aug}} = D_\xi \cup D^{\text{sink}}$ , for each  $\xi \in \mathbb{N}_0^n$ .

LEMMA 6. *Let  $\text{Poly}_1, \text{Poly}_2$  be defined as in the last subsection. For all  $\xi \in \mathbb{N}_0^n$ , we have*

$$\begin{aligned} \text{EVAL}(\text{Poly}_1, D_\xi^{\text{aug}}) &= 1 + P_1(\xi) \\ \text{EVAL}(\text{Poly}_2, D_\xi^{\text{aug}}) &= 1 + \xi_1^d P_2(\xi) \end{aligned}$$

PROOF. First of all, consider  $\text{EVAL}(\text{Poly}_1, D^{\text{sink}})$ . For the subgoal  $\text{Value}(u_i)$ , the only fact in  $D^{\text{sink}}$  involving  $R$  is  $R(C_0, C_0)$ , therefore each  $u_i$  must map to  $C_0$ . Now, examining the subgoal  $\text{Term}(\vec{u}, z_0)$ , it is clear that with each  $u_i$  mapping to  $C_0$ ,  $z_0$  must map to  $C_0$  as well. Further, the subgoal  $S_0(z_{N_m}, z_{N_m})$  can only be satisfied by mapping  $z_{N_m}$  to  $C_0$ . So  $z_\ell$  must map to  $C_0$  as well, for all  $\ell \in [N_m]$ . Hence,  $\text{EVAL}(\text{Poly}_1, D^{\text{sink}})$  returns the tuple  $(C_0, \dots, C_0)$  with multiplicity 1.

To finish the proof, observe that the constants appearing in  $D_\xi$  and  $D^{\text{sink}}$  are disjoint, and the constraints are such that no single map of the variables can use the constants in both databases at the same time. Since  $D_\xi^{\text{aug}} = D_\xi \cup D^{\text{sink}}$ , we have

$$\begin{aligned} \text{EVAL}(\text{Poly}_1, D_\xi^{\text{aug}}) &= \text{EVAL}(\text{Poly}_1, D^{\text{sink}}) + \text{EVAL}(\text{Poly}_1, D_\xi) = 1 + P_1(\xi) \end{aligned}$$

The proof showing  $\text{EVAL}(\text{Poly}_2, D_\xi^{\text{aug}}) = 1 + \xi_1^d P_2(\xi)$  follows similarly.  $\square$

We would like to guarantee that the databases we consider contain a copy of the augmented polynomial encoder  $D_0^{\text{aug}} = D^* \cup D^{\text{sink}}$ , where  $D^*$  is defined in Section 4.2. Recall that the *canonical query*  $Q_D$  corresponding to a database  $D$  is one where there is a distinct variable  $x$  corresponding to a domain element  $X$  in the database; the query is a conjunction of subgoals  $R(x_1, \dots, x_k)$  such that  $R(X_1, \dots, X_k)$  is a fact in the database. All variables are free in  $Q_D$ . The crucial property is that given any database  $D'$ , any map for evaluating  $Q_D$  on  $D'$  induces a copy of  $D$  in  $D'$ . We define  $\text{AugDB}$  as the conjunction of two sub-queries (1) the canonical query associated with  $D_0^{\text{aug}}$ , in which we identify variable  $x_i$  with constant  $X_i$  for each  $i \in [n]$ , variable  $t_j$  with constant  $T_j$  for each  $j \in [m]$ , and variables  $c_0, c_1$  with constants  $C_0, C_1$  respectively, and (2) inequalities constraints on every pair of distinct variables. All the variables are free in  $\text{AugDB}$ .

Recall the definitions of a homomorphism  $\tau$  and the database  $\tau(Q)$  for a query  $Q$ , which were given at the end of Section 2. Notice that the inequalities of (2) guarantee that if  $\psi$  is a homomorphism from  $\text{AugDB}$  to  $D$ , then  $\psi(\text{AugDB})$  is isomorphic to  $D_0^{\text{aug}}$ , since distinct variables are mapped to distinct constants. In particular, this means that if  $\text{EVAL}(\text{AugDB}, D) > 0$ , then  $D$  must contain an isomorphic copy of  $D_0^{\text{aug}}$ . We restate this in the following lemma

LEMMA 7. *Let  $D$  be a database, and let  $\psi$  be a homomorphism from  $\text{AugDB}$  to  $D$ . Then  $\psi(\text{AugDB})$  is isomorphic to  $D_0^{\text{aug}}$ . Therefore, if  $\text{EVAL}(\text{AugDB}, D) > 0$ , then  $D$  must contain an isomorphic copy of  $D_0^{\text{aug}}$ .*

Next, we define our queries

$$\begin{aligned} Q_1 &:- \text{Poly}_1, \text{AugDB}(\vec{x}, \vec{t}, \vec{c}) \\ Q_2 &:- \text{Poly}_2, \text{CounterCheating}(\vec{x}, \vec{t}, \vec{c}) \end{aligned}$$

We are now ready to define cheating more formally.

Let  $D$  be a database, and suppose  $D' \subseteq D$  (i.e. all facts in  $D'$  belong to  $D$ ) is a database isomorphic to  $D_0^{\text{aug}}$ . Let  $\phi$  be a homomorphism from  $\text{Poly}_1$  to  $D$ . Roughly speaking, we think of  $\phi$  as cheating with respect to  $D'$  if any of the subgoals—ignoring those subgoals involving  $\text{Value}$ —are mapped to facts not in  $D'$ . More formally, define  $Q := \text{Term}(\vec{u}, z_0), \text{Coeff}_\alpha(z_0)$ . Notice that  $Q$  is essentially the query  $\text{Poly}_1$  with all subgoals involving  $\text{Value}$  removed. Further, let  $\phi'$  be the restriction of  $\phi$  to the variables in  $\text{vars}(Q) = \{u_1, \dots, u_d, z_0, \dots, z_{N_m}\}$ . Notice that  $\phi'$  is a homomorphism from  $Q$  to  $D$ . We say that  $\phi$  *cheats with respect to  $D'$*  if  $\phi'$  is not a homomorphism from  $Q$  to  $D'$ .

Now, for each homomorphism,  $\psi$ , from  $\text{AugDB}$  to  $D$ , the database  $\psi(\text{AugDB})$  is isomorphic to  $D_0^{\text{aug}}$ , by Lemma 7. The following key technical lemma guarantees that for each of these  $\psi(\text{AugDB})$ , the value of  $\text{CounterCheating}$  is large enough to counter the number of homomorphisms from  $\text{Poly}_1$  to  $D$  that cheat with respect to  $\psi(\text{AugDB})$ . We defer the proof of this, as well as the definition of  $\text{CounterCheating}$ , to the next subsection.

LEMMA 8. *Let  $D$  be a database, and let  $\text{Poly}_1$  be defined as in the last subsection.*

- *If  $D$  is an augmented polynomial encoder, then*

$$\text{EVAL}(\text{AugDB}, D) = \text{EVAL}(\text{CounterCheating}, D)$$

*Both return only the tuple  $(\vec{X}, \vec{T}, \vec{C})$ , with multiplicity one.*

- *For general  $D$ , let  $\psi$  be a homomorphism from  $\text{AugDB}$  to  $D$ , if one exists. If there are  $\gamma$  homomorphisms from  $\text{Poly}_1$  to  $D$  that cheat with respect to  $\psi(\text{AugDB})$ , then*

$$|\text{CounterCheating}(\psi(\vec{x}, \vec{t}, \vec{c}))|_D \geq 1 + \gamma$$

With Lemma 6 and Lemma 8 in hand, we are ready for the main result of this paper. Once  $\text{CounterCheating}$  is defined, it will be easy to see that the queries  $Q_1$  and  $Q_2$  use at most  $n^{2d}$  inequality constraints. Remember here that we only need  $n = 59$  and  $d = 5$ .

THEOREM 9. *For some  $k$  bounded by  $n^{2d}$ , where  $n \leq 59$  and  $d \leq 5$ , the problem  $\text{ConQC}_{\text{BS}}(k, \infty, \infty)$  is undecidable.*

PROOF. We will show that  $Q_1 \subseteq_{\text{BS}} Q_2$  if and only if  $P_1(\xi) \leq (\xi_1)^d P_2(\xi)$  for all  $\xi \in \mathbb{N}_0^n$ . It will thus follow by Theorem 3 that  $\text{ConQC}_{\text{BS}}(k, \infty, \infty)$  is undecidable.

First, suppose that there is a  $\xi \in \mathbb{N}_0^n$  such that  $P_1(\xi) > \xi_1^d P_2(\xi)$ . Lemma 6 shows that  $\text{EVAL}(\text{Poly}_1, D_\xi^{\text{aug}}) = 1 + P_1(\xi) > 1 + \xi_1^d P_2(\xi) = \text{EVAL}(\text{Poly}_2, D_\xi^{\text{aug}})$ . Further, Lemma 8 guarantees that  $\text{EVAL}(\text{AugDB}, D_\xi^{\text{aug}}) = \text{EVAL}(\text{CounterCheating}, D_\xi^{\text{aug}})$ . Hence,  $\text{EVAL}(Q_1, D_\xi^{\text{aug}}) > \text{EVAL}(Q_2, D_\xi^{\text{aug}})$ , as we wanted.

Now, suppose that  $P_1(\xi) \leq \xi_1^d P_2(\xi)$  for all  $\xi \in \mathbb{N}_0^n$ . Let  $D$  be a database, and let  $\Psi$  be the set of homomorphisms from  $\text{AugDB}$  to  $D$ . For each  $\psi \in \Psi$ , let  $N_\psi$  be the number of homomorphisms from  $\text{Poly}_1$  to  $D$  that cheat with respect to  $\psi$ . Further, we need to count the number of homomorphisms from  $\text{Poly}_1$  to  $D$  that do not cheat with respect to  $\psi$ .

Recall that in our definition of cheating, we intentionally ignored the view  $\text{Value}$ , which has the effect of ignoring relation  $R$ . But homomorphisms from  $\text{Poly}_1$  to  $D$  must respect  $R$ . So we need a way to “add  $R$  back.”

To this end, define  $D^\psi$  iteratively as follows:

1. Start with  $D^\psi$  equal to  $\psi(\text{AugDB})$ .
2. For each fact of  $D$ , check if it is of the form  $R(\psi(x_i), C)$  for some  $i$  and some constant  $C$ . If it is, then add that fact to  $D^\psi$  (with multiplicity one)

It is not hard to see that  $D^\psi$  is isomorphic to  $D_{\xi_\psi}^{\text{aug}}$  for some  $\xi_\psi$ . (In particular, the value of its  $i$ th entry,  $(\xi_\psi)_i = |\text{Value}(\psi(x_i))|_D$ .) It is also not hard to see that if  $\phi$  does not cheat with respect to  $\psi(\text{AugDB})$ , then  $\phi$  is a homomorphism from  $\text{Poly}_1$  to  $D^\psi$ . Thus, the number of homomorphisms from  $\text{Poly}_1$  to  $D$  that do not cheat is

$$\text{EVAL}(\text{Poly}_1, D^\psi) = 1 + P_1(\xi_\psi)$$

by Lemma 6. Hence, we see

$$\text{EVAL}(\text{Q}_1, D) = \sum_{\psi \in \Psi} (1 + P_1(\xi_\psi) + N_\psi)$$

Denoting the first coordinate of  $\xi_\psi$  by  $\xi_{\psi,1}$ , we also see that

$$\begin{aligned} \text{EVAL}(\text{Poly}_2, D^\psi) &\geq 1 + (\xi_{\psi,1})^d P_2(\xi_\psi) \\ &\geq 1 + P_1(\xi_\psi) \end{aligned}$$

Putting this together with Lemma 8 completes the proof:

$$\begin{aligned} \text{EVAL}(\text{Q}_2, D) &\geq \sum_{\psi \in \Psi} (1 + P_1(\xi_\psi)) \cdot (1 + N_\psi) \\ &\geq \sum_{\psi \in \Psi} (1 + P_1(\xi_\psi) + N_\psi) \\ &= \text{EVAL}(\text{Q}_1, D) \end{aligned}$$

□

Thus, applying Corollary 2 (and counting the number of inequalities in  $\text{Q}_1$  and  $\text{Q}_2$  a little more carefully), we have the following corollary.

**COROLLARY 10.** *For some  $k$  bounded by  $n^{2d}$ , where  $n \leq 59$  and  $d \leq 5$ , we have that both*

$$\begin{aligned} \text{ConQC}_B(k, 1, 2) &\text{ is undecidable, and} \\ \text{ConQC}_{BS}(k, 1, 2) &\text{ is undecidable.} \end{aligned}$$

#### 4.4 Proof Sketch of Lemma 8

We now sketch the proof for the main technical lemma. We begin by defining CounterCheating. We break this into 8 main components. One of the components is simply AugDB. Each of the other 7 is to counter a different type of cheating that could occur. We outline these 7 ways of cheating below.

Recall that  $\text{Poly}_1$  was a bound query. In defining our anti-cheating views, it will be necessary to use inequalities involving the variables of  $\text{Poly}_1$ . Hence, we define  $\text{Poly}'_1$  analogously to  $\text{Poly}_1$ , but in such a way that many of its variables become free. Specifically, let  $\text{Coeff}'_\alpha(\vec{z})$  be defined as the query with the same body as  $\text{Coeff}_\alpha$ , but with free variables  $z_0, \dots, z_{N_m}$ . Then define

$$\begin{aligned} \text{Poly}'_1(\vec{u}, \vec{z}) \\ :- \quad &\text{Term}(\vec{u}, z_0), \text{Coeff}'_\alpha(\vec{z}), \text{Value}(u_1), \dots, \text{Value}(u_d) \end{aligned}$$

For all  $k \in [d]$  and all  $j \in [m]$ , we define the following views.

$$\begin{aligned} \text{Counter}_1(t_1, \dots, t_m) &:- \\ &\text{Poly}'_1(\vec{u}, \vec{z}), \wedge_{\ell \in [m]} z_0 \neq t_\ell \\ \text{Counter}_2(t_0) &:- \\ &\text{Poly}'_1(\vec{u}, \vec{z}), z_{N_m} \neq t_0 \\ \text{Counter}_3^k(\vec{x}) &:- \\ &\text{Poly}'_1(\vec{u}, \vec{z}), \wedge_{i \in [n]} u_k \neq x_i \\ \text{Counter}_4^k &:- \\ &\text{Poly}'_1(\vec{u}, \vec{z}), R_k(w, z_0), w \neq u_k \\ \text{Counter}_5^j(t_0) &:- \\ &\text{Poly}'_1(\vec{u}, \vec{z}), S_j(z_0, w), w \neq z_0, w \neq t_0 \\ \text{Counter}_6^j(t_j) &:- \\ &\text{Poly}'_1(\vec{u}, \vec{z}), S_j(z_0, w), w \neq z_0, z_0 \neq t_j \\ \text{Counter}_7^j(t_0) &:- \\ &\text{Poly}'_1(\vec{u}, \vec{z}), S_j(z_{N_m}, w), w \neq t_0 \end{aligned}$$

Notice that each of the above views contains  $\text{Poly}'_1$  as a subgoal. So if it were not for the  $\neq$ -constraints, each view would return a multiplicity at least as large as  $\text{Poly}_1$  since  $\text{Poly}_1$  is just  $\text{Poly}'_1$  with all variables bound. The  $\neq$ -constraints guarantee that there is only one homomorphism, unless there are maps from  $\text{Poly}_1$  to the database that cheat.

Define  $\text{CounterCheating}(\vec{x}, \vec{t}, \vec{c})$  to be the conjunction of the above views and  $\text{AugDB}(\vec{x}, \vec{t}, \vec{c})$ . Note that the free variables  $(\vec{x}, \vec{t}, \vec{c})$  are the same for each view, but the variables  $(\vec{u}, \vec{z})$  are local to each view.

We first sketch the proof that if  $D$  is an augmented polynomial encoder, then

$$\text{EVAL}(\text{AugDB}, D) = \text{EVAL}(\text{CounterCheating}, D)$$

with both expressions returning only the tuple  $(\vec{X}, \vec{T}, \vec{C})$ , with multiplicity one. It is not difficult to check that  $\text{EVAL}(\text{AugDB}, D)$  returns only the tuple  $(\vec{X}, \vec{T}, \vec{C})$  with multiplicity one. Further, since CounterCheating contains AugDB as a subgoal, we see that evaluating CounterCheating on  $D$  returns only the tuple  $(\vec{X}, \vec{T}, \vec{C})$ . However, we still need to verify that its multiplicity is one.

We do this by showing that there is exactly one homomorphism from each view  $\text{Counter}_i$  to  $D$  that maps  $(\vec{x}, \vec{t}, \vec{c})$  to  $(\vec{X}, \vec{T}, \vec{C})$ . Consider the view  $\text{Counter}_6^j$  for some  $j$ . The proof for each of the other views  $\text{Counter}_i$  is similar. Since  $\text{Counter}_6^j$  contains  $\text{Poly}'_1(\vec{u}, \vec{z})$  as a subgoal, there is precisely one homomorphism that maps  $(\vec{u}, \vec{z})$  to  $(C_0, \dots, C_0)$  and  $w$  to  $C_1$ . In analogy with the proof of Lemma 4, we see that any other homomorphism  $\phi$  must map  $z_0$  to  $T_{j'}$  for some  $j'$ . Now, if  $j' \neq j$ , then  $\phi(w) = T_{j'}$ , contradicting the  $\neq$ -constraint  $w \neq z_0$ . If  $j' = j$ , then  $\phi(z_0) = T_j$ , contradicting  $z_0 \neq t_j$ .

We now sketch the proof of the second item. Let  $D$  be any database, and let  $\psi$  be a homomorphism from AugDB to  $D$ . Further, suppose that there are  $\gamma$  homomorphisms from  $\text{Poly}_1$  to  $D$  that cheat with respect to  $\psi(\text{AugDB})$ . We will show

$$|\text{CounterCheating}(\psi(\vec{u}, \vec{t}, \vec{c}))|_D \geq 1 + \gamma$$

For convenience and readability, relabel the constants in  $D$  so that  $\psi(x_i) = X_i$  for all  $i \in [n]$ ,  $\psi(t_j) = T_j$  for all  $j \in [m]$ , and  $\psi(c_\ell) = C_\ell$  for  $\ell = 0, 1$ . Under the relabeling of the constants,  $\psi(\text{AugDB})$  is the canonical database  $D_0^{\text{aug}}$ .

We partition cheating homomorphisms into seven classes. Notice that each class  $N_i$  corresponds to the respective  $\text{Counter}_i$ . Let  $\phi_0$  be the homomorphism mapping every variable in  $\text{Q}'_1$  to  $C_0$ , and let  $\Phi$



be the set of homomorphisms from  $Q'_1$  to  $D$  excluding  $\phi_0$ . For each  $j \in [m], k \in [d]$ , define

1.  $N_1 = \{\phi \in \Phi \mid \text{for all } j \in [m], \phi(z_0) \neq T_j\}$ .
2.  $N_2 = \{\phi \in \Phi \mid \phi(z_{N_m}) \neq T_0\}$ .
3.  $N_3^k = \{\phi \in \Phi \mid \text{for all } i \in [n], \phi(u_k) \neq X_i\}$ .
4.  $N_4^k = \{\phi \in \Phi \mid \phi(u_1, \dots, u_d, t) = (X_{i_1}, \dots, X_{i_d}, T_{j'}) \text{ for some } j' \in [m], \text{ but } i_k \text{ is not the } k\text{th entry of } T_{j'}\}$ .
5.  $N_5^j = \{\phi \in \Phi \mid \phi(z_{N_m}) = T_0, \exists \ell \text{ such that } S_j(z_\ell, z_{\ell+1}) \text{ is a subgoal of } \text{Poly}_1, \exists j' \in [m] \phi(z_0) = \phi(z_\ell) = T_{j'}, \text{ but } \phi(z_{\ell+1}) \text{ is neither } T_{j'} \text{ nor } T_0\}$ .
6.  $N_6^j = \{\phi \in \Phi \mid \phi(z_{N_m}) = T_0, \exists \ell \text{ such that } S_j(z_\ell, z_{\ell+1}) \text{ is a subgoal of } \text{Poly}_1, \exists j' \in [m] \phi(z_0) = \phi(z_\ell) = T_{j'} \text{ with } j' \neq j, \text{ but } \phi(z_{\ell+1}) = T_0\}$ .
7.  $N_7^j = \{\phi \in \Phi \mid \phi(z_{N_m}) = T_0, \exists \ell \text{ such that } S_j(z_\ell, z_{\ell+1}) \text{ is a subgoal of } \text{Poly}_1, \phi(z_\ell) = T_0, \text{ but } \phi(z_{\ell+1}) \neq T_0\}$ .

CLAIM 11. *If  $\phi$  is a homomorphism from  $\text{Poly}_1$  to  $D$  that cheats with respect to  $\psi(\text{AugDB})$ , then  $\phi$  belongs to  $N_1 \cup N_2 \cup \bigcup_{k \in [d]} (N_3^k \cup N_4^k) \cup \bigcup_{j \in [m]} (N_5^j \cup N_6^j \cup N_7^j)$ .*

PROOF. Let  $\phi$  be a homomorphism that cheats with respect to  $\psi(\text{AugDB})$ . Throughout this proof, we assume that  $\phi(z_{N_m}) = T_0$ , for otherwise  $\phi \in N_1$ . Likewise, we assume  $\phi(z_0) = T_j$  for some  $j$ , for otherwise  $\phi \in N_2$ . Finally, we assume for each  $k$ , there is an  $\ell$  such that  $\phi(u_k) = X_\ell$ , for otherwise  $\phi \in N_3^k$ .

If  $\phi$  cheats, then some subgoal of  $\text{Poly}_1$ , not involving relation  $R$ , is mapped to a fact that is not in  $\psi(\text{AugDB})$ . We say such a subgoal is *mapped incorrectly*.

First, suppose subgoal  $R_k(a, b)$  is mapped incorrectly by  $\phi$ . Note that  $a = u_k$  and  $b = z_0$ . By our assumption,  $\phi(u_k) = X_\ell$  for some  $\ell \in [n]$  and  $\phi(z_0) = T_j$  for some  $j$ . Hence,  $R_k(X_\ell, T_j)$  is not a fact for  $\psi(\text{AugDB})$ . So  $\ell$  is not the  $k$ th entry in  $T_j$ . That is,  $\phi \in N_4^k$ .

Now, suppose a subgoal involving relation  $S_{k'}$  is mapped incorrectly by  $\phi$ , for some  $k'$ . Let  $\ell$  be the smallest index such that, for some  $j \in [m]$ ,  $S_j(z_\ell, z_{\ell+1})$  is a subgoal of  $\text{Poly}_1$  that is mapped incorrectly by  $\phi$ . Note that  $\phi(z_\ell)$  is either  $T_0$  or  $T_{j'}$  for some  $j' \in [m]$ . If it is  $T_0$ , then  $\phi(z_{\ell+1}) \neq T_0$  since  $S_j(z_\ell, z_{\ell+1})$  is mapped incorrectly. Hence,  $\phi \in N_7^j$ .

If  $\phi(z_\ell) = T_{j'}$ , then  $\phi(z_0) = T_{j'}$  by the minimality of  $\ell$ . Further, either (1)  $\phi(z_{\ell+1}) = T_0$ , (2)  $\phi(z_{\ell+1}) \neq T_0, T_{j'}$ , or (3)  $\phi(z_{\ell+1}) = T_{j'}$ . In the first case, we see  $j' \neq j$  since  $S_j(z_\ell, z_{\ell+1})$  is mapped incorrectly. Hence,  $\phi \in N_6^j$ . In the second case, we see  $\phi \in N_5^j$ . In the third case, let  $\ell'$  be the smallest index such that  $\phi(z_{\ell'+1}) \neq T_{j'}$ , and let  $k'$  be such that  $S_{k'}(z_{\ell'}, z_{\ell'+1})$  is a subgoal of  $\text{Poly}_1$ . Notice that  $\phi(z_0) = \phi(z_{\ell'}) = T_{j'}$ . Further, notice that  $j > j'$  since  $S_j(T_{j'}, T_{j'})$  is not a fact for  $D_\psi$ , hence  $k' \neq j'$ . So if  $\phi(z_{\ell'+1}) = T_0$ , we see  $\phi \in N_6^{k'}$ . On the other hand, if  $\phi(z_{\ell'+1}) \neq T_0$ , then  $\phi \in N_5^{k'}$ .  $\square$

Claim 11 thus gives us that  $\gamma \leq |N_1| + |N_2| + \sum_{k \in [d]} (|N_3^k| + |N_4^k|) + \sum_{j \in [m]} (|N_5^j| + |N_6^j| + |N_7^j|)$ .

Our next claim simply says that each Counter does its job.

CLAIM 12. *For all  $k \in [d], j \in [m]$ ,*

$$\begin{aligned} |\text{Counter}_1(T_1, \dots, T_m)|_D &\geq 1 + |N_1| \\ |\text{Counter}_2(T_0)|_D &\geq 1 + |N_2| \\ |\text{Counter}_3^k(\vec{X})|_D &\geq 1 + |N_3^k| \\ |\text{Counter}_4^k|_D &\geq 1 + |N_4^k| \\ |\text{Counter}_5^j(T_0)|_D &\geq 1 + |N_5^j| \\ |\text{Counter}_6^j(T_j)|_D &\geq 1 + |N_6^j| \\ |\text{Counter}_7^j(T_0)|_D &\geq 1 + |N_7^j| \end{aligned}$$

PROOF SKETCH. Consider a homomorphism  $\phi$  mapping  $\text{Poly}_1$  to  $D$ . Then each local variable  $u_i$  in  $\text{Poly}_1$  gets mapped to  $\phi(u_i)$ , and likewise, each local variable  $z_i$  in  $Q'_1$  gets mapped to  $\phi(z_i)$ . Suppose, e.g.,  $\phi \in N_1$ . Then we may specify a corresponding homomorphism  $\hat{\phi}$  from  $\text{Counter}_1$  to  $D$  such that  $\hat{\phi}$  maps each local variable  $u_i$  of  $\text{Counter}_1$  to  $\phi(u_i)$ , and likewise,  $\hat{\phi}$  maps each local variable  $z_i$  of  $\text{Counter}_1$  to  $\phi(z_i)$ . Finally, we set  $\hat{\phi}(t_j) = T_j$  for each  $j$ . It is not hard to see that  $\hat{\phi}$  is indeed a homomorphism, and further, that it respects all the  $\neq$ -constraints of  $\text{Counter}_1$ . Hence, for every cheating homomorphism in  $N_1$ , there is a corresponding homomorphism from  $\text{Counter}_1$  to  $D$ . Furthermore, there is a homomorphism from  $\text{Counter}_1$  to  $D$  in which every variable is mapped to  $C_0$ . It thus follows that

$$|\text{Counter}_1(T_1, \dots, T_m)|_D \geq 1 + |N_1|$$

An analogous proof holds for each of the other cases. We include a proof sketch for the last four inequalities, since in each of those cases, we must also specify where the local variable  $w$  is mapped.

4. Let  $k \in [d]$ , let  $\phi \in N_4^k$ , and consider where the bound variables of  $\text{Poly}_1$  are mapped by  $\phi$ . In particular,  $\phi(u_k) = X_i$  for some  $i \in [n]$  and  $\phi(z_0) = T_j$  for some  $j \in [m]$ , but  $i$  is not the  $k$ th entry of  $T_j$ . Hence, we may map the corresponding bound variables of  $\text{Counter}_4^k$  in precisely the same way without violating any  $\neq$ -constraints. As for the variable  $w$ , we may map it to  $X_{i'}$ , where  $i'$  is the  $k$  entry of  $T_j$ . Since  $\phi(u_k) \neq \phi(w)$ , no  $\neq$ -constraint is violated.

Furthermore, there is an additional homomorphism mapping from  $\text{Counter}_4^k$  to  $D$  in which we map  $w$  to  $C_1$  and all the other bound variables to  $C_0$ , giving an extra homomorphism. Hence,

$$|\text{Counter}_4^k|_D \geq 1 + |N_4^k|$$

5. Let  $j \in [m]$ , let  $\phi \in N_5^j$ , and consider where the bound variables of  $\text{Poly}_1$  are mapped by  $\phi$ . In particular, there are variables  $z_\ell, z_{\ell+1}$  such that  $S_j(z_\ell, z_{\ell+1})$  is a subgoal of  $\text{Poly}_1$ ,  $\phi(z_\ell) = \phi(z_0) = T_{j'}$  for some  $j' \in [m]$ , but  $\phi(z_{\ell+1})$  is mapped to neither  $T_0$  nor  $T_{j'}$ . We may map the corresponding bound variables of  $\text{Counter}_5^j$  in precisely the same way without violating any  $\neq$ -constraints. As for the variable  $w$ , we may map it to  $\phi(z_{\ell+1})$ . Since  $\phi(w) \neq T_0$  and  $\phi(w) \neq \phi(z_0)$ , no  $\neq$ -constraint is violated.

Furthermore, there is an additional homomorphism mapping from  $\text{Counter}_5^j$  to  $D$  in which we map  $w$  to  $C_1$  and all the other bound variables to  $C_0$ , giving an extra homomorphism. Hence,

$$|\text{Counter}_5^j(T_0)|_D \geq 1 + |N_5^j|$$

6. Let  $j \in [m]$ , let  $\phi \in N_6^j$ , and consider where the bound variables of  $\text{Poly}_1$  are mapped by  $\phi$ . In particular, there are variables  $z_\ell, z_{\ell+1}$  such that  $S_j(z_\ell, z_{\ell+1})$  is a subgoal of  $\text{Poly}_1$ ,  $\phi(z_\ell) = \phi(z_0) = T_{j'}$  for some  $j' \in [m]$  with  $j' \neq j$ , but  $\phi(z_{\ell+1}) = T_0$ . We may map the corresponding bound variables of  $\text{Counter}_6^j$  in precisely the same way without violating any  $\neq$ -constraints. As for the variable  $w$ , we may map it to  $T_0$ , since  $S_j(T_{j'}, T_0)$  must be a fact of  $D$ . Since  $\phi(z_0) \neq T_j$  and  $\phi(w) \neq \phi(z_0)$ , no  $\neq$ -constraint is violated.

Furthermore, there is an additional homomorphism mapping from  $\text{Counter}_6^j$  to  $D$  in which we map  $w$  to  $C_1$  and all the other bound variables to  $C_0$ , giving an extra homomorphism. Hence,

$$|\text{Counter}_6^j(T_j)|_D \geq 1 + |N_6^j|$$

7. Let  $j \in [m]$ , let  $\phi \in N_7^j$ , and consider where the bound variables of  $\text{Poly}_1$  are mapped by  $\phi$ . In particular, there are variables  $z_\ell, z_{\ell+1}$  such that  $S_j(z_\ell, z_{\ell+1})$  is a subgoal of  $\text{Poly}_1$ ,  $\phi(z_\ell) = T_0$ , but  $\phi(z_{\ell+1}) \neq T_0$ . We may map the corresponding bound variables of  $\text{Counter}_7^j$  in precisely the same way without violating any  $\neq$ -constraints. As for the variable  $w$ , we may map it to  $\phi(z_{\ell+1})$ . Since  $\phi(w) \neq T_0$ , no  $\neq$ -constraint is violated.

Furthermore, there is an additional homomorphism mapping from  $\text{Counter}_7^j$  to  $D$  in which we map  $w$  to  $C_1$  and all the other bound variables to  $C_0$ , giving an extra homomorphism. Hence,

$$|\text{Counter}_7^j(T_0)|_D \geq 1 + |N_7^j|$$

□

To complete the proof:

$$\begin{aligned} & \text{EVAL}(\text{CounterCheating}(\psi(\vec{x}, \vec{b}, \vec{c})), D) \\ & \geq (1 + |N_1|)(1 + |N_2|) \prod_{k=1}^d [(1 + |N_3^k|)(1 + |N_4^k|)] \\ & \quad \cdot \prod_{j=1}^m [(1 + |N_5^j|)(1 + |N_6^j|)(1 + |N_7^j|)] \\ & \geq 1 + |N_1| + |N_2| + \sum_{k=1}^d (|N_3^k| + |N_4^k|) \\ & \quad + \sum_{j=1}^m (|N_5^j| + |N_6^j| + |N_7^j|) \\ & \geq 1 + \gamma \end{aligned}$$

## 5. Concluding Remarks

We showed that, under bag semantics, the query containment problem for conjunctive queries with inequalities is undecidable. Actually, even drastic restrictions of this problem are undecidable; specifically, undecidability persists even if the following two restrictions hold at the same time: (1) the queries involve a single binary relation symbol; and (2) the total number of inequalities is bounded by a certain fixed (albeit large) value. Furthermore, the same undecidability results hold under bag-set semantics.

These strong undecidability results reveal that there is no hope of using query containment as a tool to optimize conjunctive queries with inequalities in real database systems. At the same time, these results motivate several different lines of investigation. From a database-theoretic point of view, it would be interesting to identify syntactic or structural conditions that may give rise to classes of conjunctive queries with inequalities for which the containment problem under bag semantics is decidable (or, even better, has low complexity). From a graph-theoretic point of view, note that we made heavy use of directed graphs in the reduction from Hilbert’s 10th Problem. Thus, it is natural to ask whether, under bag semantics, the containment problem for conjunctive queries with inequalities is decidable when the queries involve a single binary relation symbol that is interpreted by an undirected graph.

Our work was originally motivated from the decidability question for the conjunctive-query containment problem (without inequalities) under bag semantics. While this question remains unanswered, we hope that the combinatorial tools developed here may turn out to be of use in resolving this long-standing question.

## 6. References

- [1] N. Brisaboa and H. Hernández. Testing bag-containment of conjunctive queries. *Acta Informatica*, 34(7):557–578, 1997.
- [2] N. Brisaboa, H. Hernández, J. Paramá, and M. Penabad. Conjunctive queries with built-in predicates with variables and constants over any ordered domain. In *ABDIS Research Communications*, pages 46–57. 1998.
- [3] A. Chandra and P. Merlin. Optimal implementation of conjunctive queries in relational databases. In *Proc. 9th ACM Symp. on Theory of Computing*, pages 77–90, 1977.
- [4] S. Chaudhuri and M. Vardi. Optimization of real conjunctive queries. In *Proceedings of the 12th ACM Symposium on the Principles of Database Systems*, pages 59–70, 1993.
- [5] S. Cohen, W. Nutt, and A. Serebrenik. Rewriting aggregate queries using views. In *Proc. 18th ACM Symp. on Principles of Database Systems*, pages 155–166, 1999.
- [6] Y. Ioannidis and R. Ramakrishnan. Containment of conjunctive queries: beyond relations as sets. *ACM Transactions on Database Systems*, 20(3):288–324, September 1995.
- [7] A. Klug. On conjunctive queries containing inequalities. *Journal of the Association for Computing Machinery*, 35(1):146–160, January 1988.
- [8] P. Kolaitis, D. Martin, and M. Thakur. On the complexity of the containment problem for conjunctive queries with built-in predicates. In *Proceedings of the 17th ACM Symposium on the Principles of Database Systems*, pages 197–204, 1998.
- [9] Y. Matiyasevich. *Hilbert’s 10th Problem*. MIT Press, 1993.
- [10] W. Nutt, Y. Sagiv, and S. Shurin. Deciding equivalences among aggregate queries. In *Proc. 17th ACM Symp. on Principles of Database Systems*, pages 214–223, 1998.
- [11] C. Papadimitriou. *Computational Complexity*. Addison-Wesley, 1994.
- [12] M. Penabad, N. Brisaboa, H. Hernández, and J. Paramá. A general procedure to check conjunctive query containment. *Acta Informatica*, 38(7):489–529, 2002.
- [13] Y. Sagiv and M. Yannakakis. Equivalences among relational expressions with the union and difference operators. *Journal of the Association for Computing Machinery*, 27(4):633–655, October 1980.
- [14] R. van der Meyden. The complexity of querying infinite data about linearly ordered domains. *Journal of Computer and System Sciences*, 54(1):113–135, Feb. 1997.
- [15] E. Vee. The equivalence problem for real conjunctive queries with inequalities. Submitted, 2006.
- [16] E. Vee. A note on “Testing bag-containment of conjunctive queries”. Submitted to *Acta Informatica*, 2006.