

# Reformulating Queries: Theory and Practice (Full version)

Anonymous

## Abstract

We consider a setting where a user wants to pose a query against a dataset where some background knowledge, expressed as logical sentences, is available, but only a subset of the information can be used to answer the query. We thus want to *reformulate* the user query against the subvocabulary, arriving at a query that is equivalent to the user’s query assuming the background theory, but using only the restricted vocabulary. We consider two variations of the problem, one where we want any such reformulation and another where we restrict the size. We first present a classification of the complexity of the problem, showing that distinct techniques are needed for the two variations above. We then present algorithms for solving the problems in practice and evaluate their performance.

## 1 Introduction

A basic problem dealt with in databases and knowledge representation is *query reformulation with background knowledge*. The input to the problem includes a logical formula  $Q$ , the *input query*, and a set of sentences of first-order logic  $\Sigma$ , the *integrity constraints* or *background theory*. We also have a specification of an “interface restriction”: what data can we actually access? For example, we can specify the allowed data as a sub-vocabulary  $\mathcal{V}$  of the vocabulary of  $\Sigma$  and  $Q$ . We want to know if the answer to  $Q$  can be obtained using the restricted interface. It is well-known that this is the case if and only if there is another formula  $Q'$  that is logical equivalent  $Q$  over all instances satisfying  $\Sigma$  (for short,  $\Sigma$ -equivalent to  $Q$ ), and which uses only relations from  $\mathcal{V}$ . If  $Q'$  is in first-order logic, we call this an  $\text{FO}, \mathcal{V}, \Sigma$ -*reformulation of  $Q$* . We can likewise restrict  $Q'$  to be in a subset of first-order logic, e.g., by restricting the connectives or quantifiers.

Variants of the reformulation problem have been considered in several communities:

- Within databases the emphasis has been on theories given as “dependencies” — Horn sentences in predicate (first-order) logic. The bulk of the work has been on “view-based reformulation”: the vocabulary is divided up into the “base relations”, and the “view relations” which are derived from

the base ones. Each view relation  $V_i \in \mathcal{V}$  is associated with a defining formula  $\varphi_i$ , and the sentences  $\Sigma$  state that  $V_i$  corresponds exactly to its definition:  $V_i(\vec{x}) \leftrightarrow \varphi_i(\vec{x})$ . The interface relations are then the  $V_i \in \mathcal{V}$ . The problem has been extensively studied in the case where the “view definitions”  $\varphi_i$  are built up from  $\wedge, \exists$  (that is, they are *conjunctive queries*, or *CQs*). One desires a target  $Q'$  as another CQ, perhaps adding a requirement that  $Q'$  contains no non-redundant atoms. The view-based reformulation problem has been extended to reformulation with dependencies, using the *chase and backchase (C&B) method* [Deutsch *et al.*, 2006; Popa, 2000]. The method proceeds by first generating consequences under the dependences (the *chasing phase*) and then searching for a minimal subset of the consequences that are equivalent to the original formula  $Q$  (the *backchasing phase*). The C&B exploits properties specific to dependencies: the full set of consequences can often be generated compactly; further the consequences can be instrumented in such a way as to make the backchasing phase efficient [Meier, 2014; Ileana *et al.*, 2014]. Reformulation over both subsignatures and functional interfaces is considered in [Toman and Weddell, 2011; Benedikt *et al.*, 2016]. [Toman and Weddell, 2011] describes implementations of reformulation via interpolation and via the C&B method, but the implementation discussion focuses (as with [Meier, 2014; Ileana *et al.*, 2014]) on the case without disjunction. In contrast, [Benedikt *et al.*, 2016] deals strictly with theoretical issues, providing complexity bounds on the existence of a reformulation for richer languages such as the guarded fragment, but not for intermediate languages, and not for reformulation with a size bound.

- In the knowledge representation community both the background knowledge and the problem statement have been different. For background theories, the focus has been on *description logics*. The signature is usually restricted to have only unary and binary predicates; but the sentences may allow more flexibility than the Horn fragments considered in the DB community, including both disjunction and negation. The emphasis has not been on reformulation as defined above, where we desire a target formula that gives exactly the same results as the input on an instance satisfying the background theory. Instead the emphasis is on a more general problem, getting the best upper approximation of a formula that uses a given vocabulary. Variants of this problem have been studied as a *forgetting* or *uniform interpolation* problem [Lang *et al.*,

2003; Lutz and Wolter, 2011; Koopmann and Nikitina, 2017; Konev *et al.*, 2009; Lutz *et al.*, 2012; Nikitina and Rudolph, 2014]. These works are interested in approximating theories by other theories (e.g., in a restricted signature), while staying within a particular decidable fragment. This is motivated by issues of modularity in ontologies [Kontchakov *et al.*, 2009; Jiménez-Ruiz *et al.*, 2008].

In this paper our input will be a theory representing background knowledge and a simple formula (e.g., CQ or Boolean combination of CQs). We look for a formula that is *exactly* equivalent relative to the background theory, as in the DB work above. Thus the tools will be quite different from works considering the “best approximation” of a theory as in [Lutz and Wolter, 2011; Koopmann and Nikitina, 2017; Lutz and Wolter, 2011; Konev *et al.*, 2009; Lutz *et al.*, 2012; Nikitina and Rudolph, 2014]. However, unlike work in the DB community, we consider theories that allow disjunction and negation. We distinguish the problem of determining if *some* reformulation exists from the problem of finding one that is within a certain *size bound*; the latter problem has not been studied even in the case of dependencies. Both problems can be attacked by a reduction to entailment. For finding some reformulation there is a reduction using interpolation in theorem proving, dating back to the birth of interpolation [Craig, 1957b]. For size-bounded reformulation one can use a “guess-and-check” method that is analogous to the C&B, but accommodating disjunction and negation in the theories. Surprisingly, we give lower bounds showing that both these approaches provide optimal complexity for their respective problems, showing this over a range of logics and targets for reformulation. After completing an examination of the complexity of reformulation in Section 3, we turn to practice in Section 4. We discuss several ways of implementing the approaches of Section 3, focusing on an approach using a new variant of Huang’s interpolation algorithm for resolution [Huang, 1995]. We perform an extensive experimental comparison of the “Optimized Huang” approach on top of an existing theorem prover with several alternatives, including both an extension of C&B to disjunction (based on [Deutsch *et al.*, 2008]) and an approach using an interpolation algorithm that is integrated with a theorem-prover. *To the best of our knowledge our work is the first look at query reformulation over a subsignature for logics that include disjunction in arbitrary arity, either from the point of view of complexity or in practice.* **Organization:** Section 3 provides a theoretical study of the complexity of reformulation, considering both propositional and predicate logic theories. Section 4 turns to reformulation in practice, outlining our implementation strategy and experimental testbed, and overviewing our experimental results. Conclusions are given in Section 5.

**Organization:** Section 3 provides a theoretical study of the complexity of reformulation, considering both propositional and predicate logic theories. Section 4 turns to reformulation in practice, outlining our implementation strategy and experimental testbed, and overviewing our experimental results. Conclusions are given in Section 5.

## 2 Preliminaries

The most basic problem we consider is the *reformulation existence problem*, denoted  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \mathcal{L})$ , where:

- $Q$  is a first-order sentence (the input query),
- $\Sigma$  a finite set of first-order sentences (the background theory),
- $\mathcal{V}$  a subset of the relations mentioned in  $\Sigma$ ,
- $\mathcal{L}$  is a fragment of first-order logic (e.g., conjunctive queries).

The problem holds if and only if there is a query  $Q'$  in  $\mathcal{L}$  that is a  $\Sigma$ -reformulation of  $Q$  over  $\mathcal{V}$ . That is, we wish to determine if any reformulation in  $\mathcal{L}$  using vocabulary  $\mathcal{V}$  exists.

The *size-restricted reformulation problem* is denoted  $\text{REF}^{\leq}(Q, \Sigma, \mathcal{V}, \mathcal{L}, k)$ . It has an additional integer argument  $k$ , written in unary. It returns true whenever there is a  $\Sigma$ -reformulation of  $Q$  over  $\mathcal{V}$  of size at most  $k$ .

We first consider these problems when the vocabulary is in the setting of propositional logic. There we consider input queries  $Q$  consisting of a single proposition,  $\Sigma$  consisting of propositional sentences, and  $\mathcal{V}$  being a subset of the propositions mentioned in  $\Sigma$ . We can consider getting reformulations  $Q'$  that are propositional sentences ( $\text{PROP}$ ,  $\mathcal{V}$ ,  $\Sigma$ -reformulations), ones that do not use negation ( $\text{PROP}^+$ ,  $\mathcal{V}$ ,  $\Sigma$ -reformulations), sentences in disjunctive normal form ( $\text{DNF}$ ,  $\mathcal{V}$ ,  $\Sigma$ -reformulations), DNFs that are also monotone (that is, disjunctions of conjunctions of positive literals,  $\text{PROP}^{\text{UCQ}}$ ,  $\mathcal{V}$ ,  $\Sigma$ -reformulations), and sentences that use only conjunction ( $\text{PROP}^{\text{CQ}}$ ,  $\mathcal{V}$ ,  $\Sigma$ -reformulations). For our theories  $\Sigma$  we will consider general propositional formulas, Horn formulas (clauses in which there is at most one positive literal), as well as “disjunctive Horn” formulas, that is, either implications of the form  $\bigwedge_i A_i \rightarrow \bigvee_j B_j$  with  $A_i, B_j$  propositions or disjunctions  $\bigvee_j B_j$  where  $B_j$  propositions.

We also consider the first order case, where we will always consider input queries that are CQs. Our target reformulation  $Q'$  can be either a CQ, a positive existential FO ( $\text{FO}^+$ ) formula, a union of CQs (UCQ), or a general first-order sentence: thus we talk about CQ,  $\mathcal{V}$ ,  $\Sigma$ -reformulations, UCQ,  $\mathcal{V}$ ,  $\Sigma$ -reformulations, etc. Note that for  $\exists\text{REF}$ , there is no difference between  $\text{PROP}^+$  and  $\text{PROP}^{\text{UCQ}}$  targets, since these have the same expressiveness. We thus use “monotone reformulation” for both. Similarly, in the first-order case, “monotone reformulation” refers to either  $\text{FO}^+$  or UCQ.

For the first-order case, we begin with *existential rules* or *Tuple-Generating Dependencies* (TGDs), universal quantifications of formulas of the form  $\bigwedge_i R_i(\vec{x}) \rightarrow \exists \vec{y} S(\vec{x}, \vec{y})$ . Most reasoning problems concerning general TGDs are undecidable; but they become decidable when the TGDs are *full*: that is, there are no existentials on the right side. They are also decidable when rules with existentials in the head can not fire unboundedly often; this holds for the weakly-acyclic TGDs [Fagin *et al.*, 2005].

In this paper we also investigate the reformulation problem for theories consisting of *disjunctive TGDs*, sentences of the form:  $\bigwedge_i R_i(\vec{x}) \rightarrow \exists \vec{y} \bigvee_i S_i(\vec{x}, \vec{y})$ . These sentences arise naturally in data integration, both to model bi-directional relationships between a local and global schema and to capture data type restrictions. For brevity we focus on the case of *full*

disjunctive TGDs (FullDisTGDs), where there are no existentially quantified variables in the head.

When  $\varphi_1$  and  $\varphi_2$  are logical sentences we write  $\varphi_1 \models \varphi_2$  ( $\varphi_1$  entails  $\varphi_2$ ) if any model of  $\varphi_1$  is a model of  $\varphi_2$ .

### 3 The complexity of reformulation

We now investigate our reformulation problems in theory. All of our upper bounds follow from reduction to entailment of CQs for the logical fragments in question. The lower bounds will require more effort.

**Complexity of  $\exists\text{REF}$ .** The fact that the  $\exists\text{REF}$  problem for both general reformulation and monotone reformulation can be reduced to entailments from interpolation, dates back to work of Craig [Craig, 1957b]. Craig’s reduction takes as input a set of relations  $\text{Sch}$ , creating a copy  $\text{Sch}' = \{R' \mid R \in \text{Sch}\}$ . For any formula  $\varphi$  over  $\text{Sch}$ , let  $\varphi'$  be formed by replacing any relation  $R$  by  $R'$ . Given a set of relations  $\mathcal{V}$  let  $\text{ForwAx} = \bigwedge_{R \in \mathcal{V}} \forall \vec{x} R(\vec{x}) \rightarrow R'(\vec{x})$  and  $\text{BackAx} = \bigwedge_{R \in \mathcal{V}} \forall \vec{x} R'(\vec{x}) \rightarrow R(\vec{x})$ .

From the work of Craig [Craig, 1957b] and Lyndon [Lyndon, 1959] we obtain:

**Proposition 1.**  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \text{FO})$  holds iff  $Q \wedge \Sigma \wedge \text{ForwAx} \wedge \text{BackAx} \wedge \Sigma' \wedge \neg Q'$  is not satisfiable.  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \text{UCQ})$  holds iff  $Q \wedge \Sigma \wedge \text{ForwAx} \wedge \Sigma' \wedge \neg Q'$  is not satisfiable.

They also show that reformulation can be read off from a proof of unsatisfiability in a suitable proof system, using an interpolation algorithm. We discuss this further in Section 4. For now we note that this result immediately provides upper bounds when  $\Sigma$  is restricted to range over fragments with decidable entailment. For example, applying this to propositional logic we get upper bounds on determining if there is a propositional reformulation and monotone reformulation:

**Corollary 1.**  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \text{PROP})$  and  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \text{PROP}^+)$  are in CONP as  $Q$  ranges over conjunctions of propositions,  $\Sigma$  over propositional sentences.

The CONP upper bound is inherited if  $\Sigma$  consists of “disjunctive Horn clauses”  $\vee\text{Horn}$ , that is, formulas of the form  $\bigwedge A_i \rightarrow \bigvee_j B_j$ . Alternatively, these are clauses with at least one positive literal. In the case where the background theory  $\Sigma$  is in propositional Horn logic, entailment is in P and thus  $\exists\text{REF}$  for PROP and  $\text{PROP}^+$  is in P as well.

The same approach can be applied to predicate (first-order) logic. For full TGDs, entailment is in EXPTIME, and becomes NP if we fix the arity of the schema. Thus applying the reduction above,  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \text{FO})$  and  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \text{UCQ})$  are both in EXPTIME and in NP for fixed arity.

If we consider weakly acyclic TGDs, entailment is 2EXPTIME-complete [Cali *et al.*, 2010]. Using this we see that  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \text{FO})$  and  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \text{UCQ})$  are in 2EXPTIME. For full disjunctive TGDs, we can easily show that entailment is CONEXPTIME-complete, and becomes  $\Pi_2^p$ -complete when the arity of relations is fixed; note that although the syntax of disjunctive TGDs matches Disjunctive Datalog, the semantics is classical entailment, different from [Eiter *et al.*, 1997]. Thus the “Craig reduction”

above shows  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \text{FO})$  and  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \text{UCQ})$  are in CONEXPTIME and in  $\Pi_2^p$  for fixed arity.

The Craig reduction does not provide any information about reformulating when the theories are propositional but the target does not allow disjunction. But there is a simpler reduction to entailment for  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \text{PROP}^{\text{CQ}})$ : letting  $S$  be the propositions in  $\mathcal{V}$  that are entailed by  $Q \wedge \Sigma$ , this holds exactly when  $S \wedge \Sigma$  entails  $Q$ . We can guess counterexample models for implication of each proposition not in  $S$  and also a counterexample to  $S \wedge \Sigma$  entailing  $Q$  to get a CONP bound.

In the case of theories  $\Sigma$  that are Horn (in particular, those given by TGDs) it is known that if a CQ  $Q$  has a UCQ reformulation over  $\mathcal{V}$  with respect to  $\Sigma$  then it has a CQ reformulation as well. Indeed, one of the disjuncts in the UCQ reformulation must already give a CQ reformulation. Thus we get upper bounds for CQ-reformulation for classes of TGDs. For disjunctive TGDs, disjunction may be essential in the reformulation; thus for CQ-reformulation of disjunctive TGDs, we can not make use of either the Craig reduction or a reduction that considers each atom at a time (as for  $\text{PROP}^{\text{CQ}}$ ). We can show that  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \text{CQ})$  is undecidable when  $\Sigma$  varies over full disjunctive TGDs this problem is undecidable. The proof is by reduction from Datalog containment. Intuitively, solving  $\exists\text{REF}$  requires us to be able to decide whether a UCQ  $Q$  is equivalent to a CQ  $Q'$  with respect to a set of full TGDs  $\Sigma$ , and this is the same as checking that the Datalog programs  $\Sigma \cup \{Q \rightarrow \text{Goal}\}$  and  $\Sigma \cup \{Q' \rightarrow \text{Goal}\}$  are equivalent.

Table 1 summarizes our bounds for  $\exists\text{REF}$ . The second column shows the bounds one gets for reformulating as a general formula (propositional or first-order, depending on the theories  $\Sigma$ ), as a monotone formula, or as a union of conjunctions: in each case the bounds are the same. The last column shows the complexity of reformulating as conjunction of atoms. In the case of FullTGD and FullDisTGD, we give the complexity both for the general case and when the arity of the schema is fixed. One can easily show that all of the complexity bounds are tight, by a reduction from entailment to reformulation. Full proofs are in the appendix.

**Complexity of  $\text{REF}^{\leq}$ .** The problem  $\text{REF}^{\leq}$  does not have a similar direct reduction to entailment. But one can reduce to “guessing plus entailment”. The machine guesses  $Q'$  and then checks entailment using the oracle. For example, from this “guessing and verifying” technique we get an upper bound on  $\text{REF}^{\leq}$  for propositional logic.

**Proposition 2.** The problems  $\text{REF}^{\leq}(Q, \Sigma, \mathcal{V}, \text{PROP}, k)$ ,  $\text{REF}^{\leq}(Q, \Sigma, \mathcal{V}, \text{PROP}^+, k)$ , and  $\text{REF}^{\leq}(Q, \Sigma, \mathcal{V}, \text{PROP}^{\text{CQ}}, k)$  are in  $\Sigma_2^p$ , as  $\Sigma$  varies over formulas in propositional logic. For  $\Sigma$  varying over  $\text{PROP}^{\text{HORN}}$ , these are all in NP.

For full TGDs, the approach above gives a bound of EXPTIME for all of the variants of  $\text{REF}^{\leq}$ , since the verification stage is Datalog evaluation, done in EXPTIME, and this subsumes the time taken for the guessing phase. The problem is in NP when the arity is fixed. For weakly-acyclic TGDs the approach gives 2EXPTIME even for fixed arity. For FullDisTGD it gives CONEXPTIME (since the guessing phase is again negligible on top of entailment) and  $\Sigma_3^p$  for fixed arity, since the guessing stage adds an alternation on top of

Table 1: Summary of results for  $\exists\text{REF}$ 

Constraints	mon. or gen.	CQ
PROP	CoNP	CoNP
$\vee\text{Horn}$	CoNP	CoNP
$\text{PROP}^{\text{HORN}}$	P	P
FullITGD	EXPTIME, NP	EXPTIME, NP
WATGD	2EXPTIME	2EXPTIME
FullDisTGD	CoNEXPTIME, $\Pi_2^P$	Undecidable

the  $\Pi_2^P$  entailment problem.

The complexity bounds for  $\text{REF}^{\leq}$  are shown in Table 2. In the second column we list the complexity for finding a small reformulation as a monotone formula or disjunction of conjunctions; the complexity is the same in all cases. The last column shows that finding a small reformulation as a conjunction of atoms always has the same complexity. As with  $\exists\text{REF}$ , for FullITGD and FullDisTGD we consider both the general case and fixed arity.

All of the lower bounds provided by the guess-and-check method turn out to be tight, but unlike for  $\exists\text{REF}$ , this is not obvious. For  $\text{REF}^{\leq}(Q, \Sigma, \mathcal{V}, \text{PROP}, k)$  where  $\Sigma$  ranges over arbitrary propositional formulas, the  $\Sigma_2^P$  lower bound follows from results on minimization of Boolean formulas in [Umans, 2001]. Indeed, using related results [Buchfuhrer and Umans, 2011] one can show hardness for the variant of the problem where we fix the number of alternations of  $\vee$  and  $\wedge$ . The technique of [Buchfuhrer and Umans, 2011] requires formulas with hard satisfiability problems, while  $\vee\text{Horn}$  formulas are always satisfiable. Thus for  $\vee\text{Horn}$  we need a more-involved “native” reduction for  $\Sigma_2^P$ -hardness:

**Theorem 1.**  $\text{REF}^{\leq}(Q, \Sigma, \mathcal{V}, \mathcal{L}, k)$  for  $\mathcal{L} \in \{\text{PROP}^+, \text{PROP}^{\text{UCQ}}, \text{PROP}^{\text{CQ}}\}$  is  $\Sigma_2^P$ -hard when  $\Sigma$  ranges over  $\vee\text{Horn}$  formulas.

*Sketch.* The proof is by reduction of  $\exists\forall 3\text{SAT}$  problem. Suppose  $\varphi = \exists \bar{x} \forall \bar{y} \neg \psi$  is formula with  $\bar{x}$  and  $\bar{y}$  tuples of propositional variables, and  $\psi$  a conjunction of clauses  $\ell^1 \vee \ell^2 \vee \ell^3$  with each  $\ell^j$  a variable from  $\bar{x} \cup \bar{y}$  or its negation.

We construct a set of propositional disjunctive Horn formulas  $\Sigma$  and a set of propositional variables  $\mathcal{V}$  such that  $\varphi$  is valid if and only if there exists a positive formula  $Q'$  over  $\mathcal{V}$  of size  $k$  such that  $Q'$  is equivalent to a propositional variable  $Q$  over  $\Sigma$ . The vocabulary  $\mathcal{V}$  consists of variables  $\text{TrueX}_i$  and  $\text{FalseX}_i$  for each  $x_i$  in  $\bar{x}$ , while  $\Sigma$  guarantees the following:

1.  $Q$  implies all the variables in  $\mathcal{V}$ , so the equivalence of  $Q$  and  $Q'$  over  $\Sigma$  boils down to the containment of  $Q'$  in  $Q$ ;
2. containment of  $Q'$  in  $Q$  over  $\Sigma$  is only possible if  $Q'$  implies either  $\text{TrueX}_i$  or  $\text{FalseX}_i$  for each  $x_i$ ; since the size of  $Q'$  is bounded by  $k$ , it means that  $Q'$  can only be of the form  $\text{AsX}_1 \wedge \dots \wedge \text{AsX}_k$ , for  $\text{AsX}_i \in \{\text{TrueX}_i, \text{FalseX}_i\}$ —that is,  $Q'$  corresponds to an assignment of variables  $x_i$ ;
3. if a dataset and  $\Sigma$  imply  $\text{TrueX}_i$  or  $\text{FalseX}_i$  for each  $x_i$ , then they imply either  $\text{TrueY}_i$  or  $\text{FalseY}_i$  for each univer-

Table 2: Summary of results for  $\text{REF}^{\leq}$ 

Constraints	monotone	CQ
PROP	$\Sigma_2^P$	$\Sigma_2^P$
$\vee\text{Horn}$	$\Sigma_2^P$	$\Sigma_2^P$
$\text{PROP}^{\text{HORN}}$	NP	NP
FullITGD	EXPTIME, NP	EXPTIME, NP
WATGD	2EXPTIME	2EXPTIME
FullDisTGD	CoNEXPTIME, $\Sigma_3^P$	CoNEXPTIME, $\Sigma_3^P$

sally quantified  $y_i$  in  $\bar{y}$ , which corresponds to an assignment of the universally quantified variables in  $\varphi$ ;

4. if one of the minimal extensions of a dataset satisfying the rules (“disjunctive chase models”) falsifies a clause in  $\psi$  under this encoding, which corresponds to falsifying the whole of  $\psi$ , then this extension implies  $Q$ .

These guarantees imply the statement of the theorem.  $\square$

For arbitrary FullITGD, WATGD, and FullDisTGD sentences, our upper bounds match those for entailment, and there is an easy reduction from entailment that gives a matching lower bound. But when the arity is fixed, entailment for FullDisTGD is  $\Pi_2^P$ , while our reduction to entailment gives only a  $\Sigma_3^P$  upper bound for  $\text{REF}^{\leq}$ . We show  $\Sigma_3^P$ -hardness of  $\text{REF}^{\leq}$  for bounded arity FullDisTGD via a direct reduction from the  $\exists\forall 3\text{SAT}$  problem. The reduction uses a similar idea to the  $\Sigma_2^P$ -reduction in the propositional case, but it is more involved: inner existential quantification is encoded by a homomorphism check from a candidate reformulation  $Q'$  to  $Q$ . Full details are in the appendix.

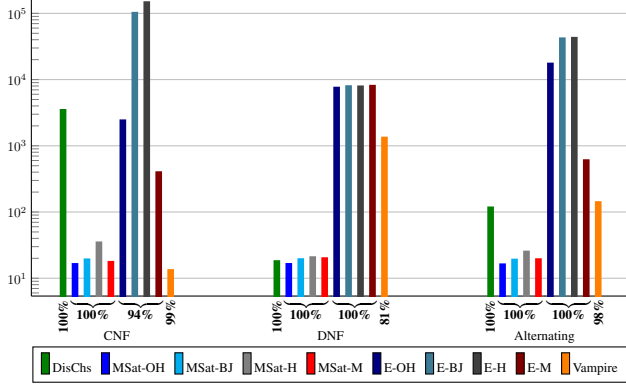
**Theorem 2.**  $\text{REF}^{\leq}(Q, \Sigma, \mathcal{V}, \mathcal{L}, k)$  for any  $\mathcal{L} \in \{\text{FO}^+, \text{UCQ}, \text{CQ}\}$  is  $\Sigma_3^P$ -hard when  $\Sigma$  ranges over sets of FullDisTGD formulas and the arity of relations is bounded by 2.

## 4 Implementing reformulation

We consider how to implement methods for finding a reformulation. The results of Section 3 show that for  $\exists\text{REF}$  there is a polynomial reduction to entailment, giving the possibility of performing reformulation on top of a theorem-prover. In contrast, our lower bounds for  $\text{REF}^{\leq}$  imply that in the presence of disjunction such a reduction is unlikely. We thus return to the method of Craig outlined in Section 3. A *Craig interpolant* for an entailment  $\varphi_1 \models \varphi_2$  is a formula  $\psi$  such that  $\varphi_1 \models \psi \models \varphi_2$  and  $\psi$  contains only relations common to  $\varphi_1$  and  $\varphi_2$ . A *Lyndon interpolant* further has the property that a relation appears positively in  $\psi$  only if it appears positively in both  $\varphi_1$  and  $\varphi_2$ , and similarly for appearing negatively. [Craig, 1957a] showed that first-order Craig interpolants exist for every first-order entailment, while [Lyndon, 1959] showed the same for Lyndon interpolants. The argument of [Craig, 1957b] shows that any Craig interpolant for the entailment

$$Q \wedge \Sigma \wedge \text{BackAx} \models (\text{ForwAx} \wedge \Sigma' \rightarrow Q')$$

Figure 1: Evaluation of reformulation approaches.



where  $\text{ForwAx}$ ,  $\Sigma'$ ,  $Q'$  are as in Section 3, can be converted to a  $\Sigma, \mathcal{V}$ -reformulation of  $Q$ ; conversely the entailment is also a necessary condition for such a reformulation to exist. A variation of the argument shows that any Lyndon interpolant for

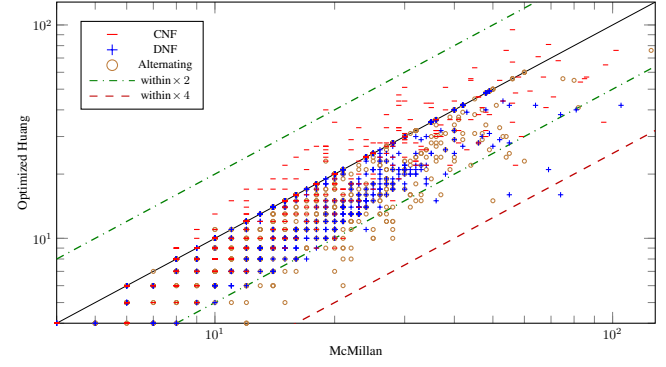
$$Q \wedge \Sigma \models (\text{ForwAx} \wedge \Sigma') \rightarrow Q'$$

is itself a monotone  $\Sigma, \mathcal{V}$ -reformulation of  $Q$ , and the entailment is likewise necessary. Note that for reformulating a conjunction of literals over propositional theory, these two notions co-incide. Thus implementing  $\exists\text{REF}$  using Craig's technique requires finding proofs witnessing the appropriate entailment, and then applying an interpolation procedure. We implemented monotone  $\Sigma, \mathcal{V}$ -reformulation for target languages with disjunction following this high-level approach, proving an entailment and applying an interpolation procedure to it. Since the goal of this work is not to build a theorem-prover from scratch, we considered both using a theorem-prover that already supports interpolation, and building our own interpolation procedure on top of existing theorem-provers. The first-order theorem prover which has the most robust native support for interpolation is Vampire [Hoder *et al.*, 2010; 2012]. Thus we test Vampire's interpolation as a "black-box" to get reformulations, using the reductions to interpolation used for all of our  $\exists\text{REF}$  problem. We complemented this with our own implementation of interpolation for resolution theorem-proving, discussed next.

A number of interpolation procedures exist for resolution, each of them creating a "partial interpolant" proceeding top-down on a resolution proof. The final interpolant is the formula associated with a leaf node representing the final contradiction in the resolution proof. Huang [Huang, 1995] and Bonacina and Johansson [Bonacina and Johansson, 2011] provide similar algorithms; the base cases assign either True or False to input clauses, while the inductive step for resolution may introduce atoms. McMillan's algorithm [McMillan, 2003], in contrast, has a more complex base case, while the propagation through resolution steps is purely compositional. For problems where we want a monotone reformulation, we need an implementation of "Lyndon interpolation" – one that does not introduce relations negated unless they appear negatively on both sides of the entailment. Huang's algorithm does not have this property, but we developed its variant that achieves this while also performing certain optimizations that we found crucial in producing good quality reformulations.

We describe the algorithm now for universal theories (the

Figure 2: Comparison between Optimized Huang and McMillan interpolation algorithms over MathSat.



extension to handle general theories is via skolemization, as described in [Huang, 1995]). Consider an entailment  $\Sigma_1 \models \Sigma_2$ , where  $\Sigma_1$  and  $\Sigma_2$  are clauses. In a resolution proof, we search for a contradiction from  $\Sigma_1 \cup \neg\Sigma_2$  where we must convert  $\neg\Sigma_2$  to a conjunction of clauses. A resolution proof is a DAG, with each node representing a derived formula within the proof. Root nodes are either elements of  $\Sigma_1$  or  $\neg\Sigma_2$ , while non-root node with formula  $\varphi$  has two parents representing formulas that resolve to  $\varphi$ . In a complete proof there is at least one leaf node which is associated with a contradiction. For purposes of interpolation, clauses in each node are first annotated with a *provenance*. Root nodes are annotated LEFT if they are in  $\Sigma_1$  and RIGHT if they are in  $\neg\Sigma_2$ . The result of a resolution step is a clause in which each literal results applying a unifier to a literal in one of the parents, and inherits the provenance of that literal.

Following [Huang, 1995], we assign partial interpolants for root nodes as: False (resp. True), if the root corresponds to a clause in  $\Sigma_1$  (resp.  $\Sigma_2$ ). For a non-root node  $c$  with parents  $p_1$  and  $p_2$ , we form the interpolant by induction. Let  $\varphi_1$  (resp.  $\varphi_2$ ) be the interpolant for  $p_1$  (resp.  $p_2$ ). Let the atoms involved in resolution be  $\alpha(\rho_1 \dots \rho_n)$  in  $p_1$  and  $\neg\alpha(\tau_1 \dots \tau_n)$  in  $p_2$ , and let  $\sigma$  be the unifier. We do case analysis on the provenance of  $\alpha(\rho_1 \dots \rho_n)$  in  $p_1$  and  $\neg\alpha(\tau_1 \dots \tau_n)$  in  $p_2$ .

- If the provenance is LEFT (resp. RIGHT) in both, we produce  $\sigma(\varphi_1) \vee \sigma(\varphi_2)$  (resp.  $\sigma(\varphi_1) \wedge \sigma(\varphi_2)$ ).
- If the provenance for  $\alpha(\rho_1 \dots \rho_n)$  in  $p_1$  is LEFT and the provenance for  $\neg\alpha(\tau_1 \dots \tau_n)$  in  $p_2$  is RIGHT, we produce  $\sigma(\varphi_1 \vee (\alpha(\rho_1, \dots, \rho_n) \wedge \varphi_2))$ .
- If the provenance for  $\alpha(\rho_1 \dots \rho_n)$  in  $p_1$  is RIGHT and the provenance for  $\neg\alpha(\tau_1 \dots \tau_n)$  in  $p_2$  is LEFT, we produce  $\sigma(\varphi_2 \vee (\neg\alpha \wedge \varphi_1))$ .

The distinction from Huang's original algorithm is that [Huang, 1995] collapses the last two cases into one, producing a unifier of  $(\neg\alpha \wedge \varphi_1) \vee (\alpha \wedge \varphi_2)$  in both cases. Indeed, it is for this reason that Huang lacks the Lyndon property. Our modification is inspired by a similar case distinction in [Bonacina and Johansson, 2011]. In the above algorithm, we assume that at each inductive step the propositional factoring and absorption rules are applied, so that  $\text{False} \vee \varphi$  is simplified to  $\varphi$ ,  $\text{False} \wedge \varphi$  to False, and so forth. Further, in the case where  $\alpha$  is propositional, we simplify the second-to-last case by producing  $\sigma(\varphi_1 \vee (\alpha \wedge \varphi_2'))$ , where  $\varphi_2'$  substitutes True for occurrences of  $\alpha$  and then simplifies as above; similarly for the last case, substituting for False in  $\varphi_1$ .

In the above algorithm, we assume that at each inductive step the propositional factoring and absorption rules are applied, so that  $\text{False} \vee \varphi$  is simplified to  $\varphi$ ,  $\text{False} \wedge \varphi$  to  $\text{False}$ , and so forth. Further, in the case where  $\alpha$  is propositional, we simplify the second-to-last case by producing  $\sigma(\varphi_1 \vee (\alpha \wedge \varphi'_2))$ , where  $\varphi'_2$  substitutes True for occurrences of  $\alpha$  and then simplifies as above; similarly for the last case, substituting for False in  $\varphi_1$ . We implement a number of other simplifications, described in [Anonymous, 2017]. We refer to this as the *Optimized Huang Algorithm* below.

**Proposition 3.** *The “Optimized Huang Algorithm” produces an interpolant  $\varphi$  for  $\Sigma_1 \models \Sigma_2$ , and any negated literal in  $\varphi$  occurs negated in  $\Sigma_1$  and  $\Sigma_2$ .*

All of these interpolation procedures require resolution proofs as inputs. For propositional logic there are several options, but we focused on MathSat [Cimatti *et al.*, 2013] to generate proofs; for first-order logic the only theorem-prover that both produces completely detailed resolution proofs and supports quantification was E [Schulz, 2013]. We also look at a variant of the chase-based techniques used for TGDs (e.g., [Deutsch *et al.*, 2006]): we take a “disjunctive chase” of the input query  $Q$  under the background theory – firing any deterministic rules and then forking into two models on each disjunction. This produces a set  $S$  of models such that another query  $Q'$  is a consequence of  $Q$  if it holds in each member of  $S$ . We can thus search for a reformulation by looking for CQs of  $\mathcal{V}$  formulas and testing whether they hold in each  $S$ . We implemented this technique using the disjunctive logic programming system DLV [Leone *et al.*, 2006] to create the set  $S$ , focusing on the propositional case.

**Testbed.** We built a generator for reformulation problems, based on two components. The first is a generator of random *target reformulations*: formulas in the desired class (e.g., PROP, PROP<sup>+</sup>). The second is a generator that takes as input a target reformulation  $\varphi$  and outputs sentences  $\Sigma$ , input query  $Q$ , and signature  $V$  such that  $\varphi$  is a  $\Sigma$ -reformulation of  $Q$  over  $V$ . The tool generates sentences that follow the structure of  $\varphi$  to witness that  $Q$  implies  $\varphi$ ; e.g., if  $\varphi$  is  $B \wedge (A_1 \vee A_2)$ ,  $\Sigma$  will contain clauses equivalent to  $Q \rightarrow C$ ,  $C \rightarrow B$ ,  $C \rightarrow D$ ,  $D \rightarrow A_1 \vee A_2$ , and additionally some “noise clauses” (that can be controlled by an input parameter).  $\Sigma$  will likewise contain sentences that witness  $\varphi$  implies  $Q$ , again with a tunable noise parameter.

We tested the following algorithms. 1. Native end-to-end interpolation using Vampire 4.1 for CASC J8. 2. Our implementations of the interpolation algorithms of [Huang, 1995], [Bonacina and Johansson, 2011], and [McMillan, 2003], and our implementation of Optimized Huang, on top of both the MathSat and E theorem provers. 3. As an alternative to interpolation for propositional theories, we implemented a variant of the C&B for logics with disjunction, based on the “disjunctive chase” of [Deutsch *et al.*, 2008], making use of the Disjunctive Datalog engine DLV. Specifically, we generate all minimal models of  $Q \wedge \Sigma$  using DLV. We then let  $\varphi$  be the disjunction over all models  $M$  of the conjunction of the literals in  $\mathcal{V}$  of  $M$ , and check if  $\varphi$  implies  $Q$  under  $\Sigma$ . It is straightforward to see that for propositional logic, if there is any reformulation over  $\mathcal{V}$ , this approach will find one. We

report experiments only for the case of propositional formulas. Since the interpolation algorithms themselves are linear, the runtime performance of interpolation-based approaches is dominated by theorem-proving time. MathSat was the fastest, requiring only a few milliseconds. The disjunctive chase approach was next, always terminating within 86secs. We timed out any other computation within  $5\times$  the time of the disjunctive chase, since the latter gives an idea of the complexity of brute-force search.

We summarize results on the size of reformulations in Figure 1, with further details on the experiments available in the appendix. On the left we see the average output size of all tools when the “ideal target reformulation” is a CNF, a DNF, or an alternation of disjunction and conjunction, with the number of alternations varying from 2 to 5. For each category the results are averages over more than 600 reformulation problems randomly generated. The bars show the averages only over problems where the tools could prove the existence of a reformulation at all, with the percentages where the theorem prover could find any reformulation prior to timeout shown below each bar. The disjunctive chase approach always generates reformulations as DNF, so it is not surprising that it performs badly in cases where the ideal reformulation is a CNF or alternates connectives. Even in the case of DNF, it is inferior to that of interpolation-based procedures; intuitively the disjunctive chase fails to identify succinct representations achieved by leveraging repetition of atoms across different minimal models. Reformulation using Vampire end-to-end failed to find any reformulation within the timeout in a significant percentage of the cases, and further the size of reformulation was, in average, quite poor. Due to the integrated nature of this approach we can only speculate on the reason; but our conjecture is that the interpolation procedures of Vampire are very geared towards examples coming from verification (and in particular the use of interpreted theories, as in [Hoder *et al.*, 2010; 2012]). Our experiments can be seen as indicating a strong difference between interpolation-for-reformulation and interpolation-for-verification.

The remaining entries in the left figure use our own interpolation implementation over distinct theorem provers. They show the superiority of MathSat over E for all interpolation algorithms, and the superiority of our optimized Huang algorithm and McMillan’s algorithm over the other implementations. Thus the right of Figure 1 contains a finer comparison of these two algorithms on top of MathSat. We see that Optimized Huang is superior on most examples, although not all. We also note that our experiments show McMillan’s algorithm to be much more robust to inefficiencies in the underlying theorem prover – that is, even with long proofs it can produce relatively smaller reformulations. This may be explained by the compositionality of McMillan’s algorithm, with the cases for resolution being simple combination of components. On the contrary, optimized Huang introduces a literal in the interpolant at each LEFT-RIGHT and RIGHT-LEFT resolution step. This may consequently induce some redundancy in the obtained formula that is not solvable by our current simplification techniques.

## 5 Conclusion

Our theoretical results shows that the most straightforward algorithms for  $\exists$ REF and  $\text{REF}^{\leq}$  provide the optimal complexity. While  $\exists$ REF can be polynomially reduced to entailment over the corresponding logic, our lower bounds show that  $\text{REF}^{\leq}$  requires an additional layer of non-determinism. There is no obvious implementation other than brute-force for this. In our experiments we focused on the interpolation-based approach, showing that for propositional theories good performance can be achieved with an optimized resolution-based interpolation algorithm on top of a high-performance theorem-prover. Limitations in the proofs exposed by first-order theorem-provers currently represent a bottleneck to application of this approach to richer logics. An interesting direction is to allow the search heuristics of a theorem-prover to be guided by the size of partial-interpolants, thus giving “interpolant-driven theorem-proving”. We are examining the implementation of this over open-source provers such as E.

## References

- [Anonymous, 2017] Anonymous. Query reformulation: Theory and practice, 2017. Available at [www.github.com/qreform/ijcail7qreform](http://www.github.com/qreform/ijcail7qreform).
- [Benedikt et al., 2016] Michael Benedikt, Balden Ten Cate, Julien Leblay, and Efthymia Tsamoura. *Generating plans from proofs: the interpolation-based approach to query reformulation*. Morgan Claypool, 2016.
- [Bonacina and Johansson, 2011] Maria Paola Bonacina and Moa Johansson. On interpolation in decision procedures. In *Automated Reasoning with Analytic Tableaux and Related Methods*. 2011.
- [Buchfuhrer and Umans, 2011] David Buchfuhrer and Christopher Umans. The complexity of boolean formula minimization. *J. Comput. Syst. Sci.*, 77(1):142–153, 2011.
- [Calì et al., 2010] Andrea Calì, Georg Gottlob, and Andreas Pieris. Query Answering under Non-guarded Rules in Datalog+/- . In *Web Reasoning and Rule Systems*, 2010.
- [Cimatti et al., 2013] Alessandro Cimatti, Alberto Griggio, Bastiaan Joost Schaafsma, and Roberto Sebastiani. The MathSAT5 SMT solver. In *TACAS*, 2013.
- [Craig, 1957a] William Craig. Linear reasoning. a new form of the Herbrand-Gentzen theorem. *The Journal of Symbolic Logic*, 22(03):250–268, 1957.
- [Craig, 1957b] William Craig. Three uses of the Herbrand-Gentzen theorem in relating model theory and proof theory. *Journal of Symbolic Logic*, 22(3):269–285, 1957.
- [Deutsch et al., 2006] Alin Deutsch, Lucian Popa, and Val Tannen. Query reformulation with constraints. *SIGMOD Record*, 35(1):65–73, 2006.
- [Deutsch et al., 2008] Alin Deutsch, Alan Nash, and Jeff Remmel. The chase revisited. In *PODS*, 2008.
- [Eiter et al., 1997] Thomas Eiter, Georg Gottlob, and Heikki Mannila. Disjunctive Datalog. *ACM Trans. Database Syst.*, 22(3):364–418, 1997.
- [Fagin et al., 2005] Ronald Fagin, Phokion G. Kolaitis, Renee J. Miller, and Lucian Popa. Data exchange: Semantics and query answering. *Theoretical Computer Science*, 336(1):89–124, 2005.
- [Hoder et al., 2010] Kryštof Hoder, Laura Kovács, and Andrei Voronkov. *Interpolation and Symbol Elimination in Vampire*. 2010.
- [Hoder et al., 2012] Krystof Hoder, Andreas Holzer, Laura Kovács, and Andrei Voronkov. Vinter: A Vampire-Based tool for interpolation. In *APLAS*, 2012.
- [Huang, 1995] Guoxiang Huang. Constructing Craig interpolation formulas. In *Computing and Combinatorics*. 1995.
- [Ileana et al., 2014] Ioana Ileana, Bogdan Cautis, Alin Deutsch, and Yannis Katsis. Complete yet practical search for minimal query reformulations under constraints. In *SIGMOD*, 2014.
- [Jiménez-Ruiz et al., 2008] Ernesto Jiménez-Ruiz, Bernardo Cuenca Grau, Ulrike Sattler, Thomas Schneider, and Rafael Berlanga. Safe and economic re-use of ontologies: A logic-based methodology and tool support. In *ESWC*, 2008.
- [Konev et al., 2009] Boris Konev, Dirk Walther, and Frank Wolter. Forgetting and uniform interpolation in large-scale description logic terminologies. In *IJCAI*, 2009.
- [Kontchakov et al., 2009] Roman Kontchakov, Luca Pulina, Ulrike Sattler, Thomas Schneider, Petra Selmer, Frank Wolter, and Michael Zakharyashev. Minimal module extraction from DL-Lite ontologies using QBF solvers. In *IJCAI*, 2009.
- [Koopmann and Nikitina, 2017] Patrick Koopmann and Nadeschda Nikitina. Small is beautiful: Computing minimal equivalent EL concepts. In *AAAI*, 2017.
- [Lang et al., 2003] Jérôme Lang, Paolo Liberatore, and Pierre Marquis. Propositional independence: Formula-variable independence and forgetting. *J. Artif. Intell. Res. (JAIR)*, 18:391–443, 2003.
- [Leone et al., 2006] N. Leone, G. Pfeifer, W. Faber, T. Eiter, G. Gottlob, S. Perri, and F. Scarcello. The DLV system for knowledge representation and reasoning. *TOCL*, 7(3):499–562, 2006.
- [Lutz and Wolter, 2011] Carsten Lutz and Frank Wolter. Foundations for uniform interpolation and forgetting in expressive description logics. In *IJCAI*, 2011.
- [Lutz et al., 2012] Carsten Lutz, Inanç Seylan, and Frank Wolter. An automata-theoretic approach to uniform interpolation and approximation in the description logic EL. In *KR*, 2012.
- [Lyndon, 1959] Roger C. Lyndon. An interpolation theorem in the predicate calculus. *Pacific Journal of Mathematics*, 9:129–142, 1959.
- [McMillan, 2003] K.L. McMillan. Interpolation and sat-based model checking. In *CAV*. 2003.

- [Meier, 2014] Michael Meier. The backchase revisited. *VLDB Journal*, 23(3):495–516, 2014.
- [Nikitina and Rudolph, 2014] Nadeschda Nikitina and Sebastian Rudolph. (non-)succinctness of uniform interpolants of general terminologies in the description logic EL. *Artif. Intell.*, 215:120–140, 2014.
- [Popa, 2000] Lucian Popa. *Object/Relational Query Optimization with Chase and Backchase*. PhD thesis, U. Penn., 2000.
- [Sagiv and Yannakakis, 1980] Yehoshua Sagiv and Mihalis Yannakakis. “Equivalences Among Relational Expressions with the Union and Difference Operators”. *Journal of the ACM*, 27(4):633–655, 1980.
- [Schulz, 2013] Stephan Schulz. System Description: E 1.8. In *LPAR*, 2013.
- [Shmueli, 1993] Oded Shmueli. Equivalence of datalog queries is undecidable. *The Journal of Logic Programming*, 15(3):231 – 241, 1993.
- [Toman and Weddell, 2011] D. Toman and G. Weddell. *Fundamentals of Physical Design and Query Compilation*. Morgan Claypool, 2011.
- [Umans, 2001] Christopher Umans. The minimum equivalent DNF problem and shortest implicants. *J. Comput. Syst. Sci.*, 63(4):597–611, 2001.



## Proofs for reformulation results for $\exists\text{REF}$

We first show both upper and lower bounds for  $\exists\text{REF}$  for general propositional logic formulas:

**Proposition 4.** *The problems  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \text{PROP})$ ,  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \text{PROP}^+)$ , and  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \text{PROP}^{\text{CQ}})$  as  $Q$  ranges over conjunctions of propositions and  $\Sigma$  over propositional sentences, are all CONP-complete.*

*Proof.* The upper bound for the first two comes from Corollary 1. For the third, to disprove the existence of such a reformulation we guess a set  $S$  of propositions in  $\mathcal{V}$ , and confirm that each proposition in the list is not implied by  $Q \wedge \Sigma$ , while confirming that the remaining propositions along with  $\Sigma$  do not imply  $Q$ . Both confirmations can be done in NP.

The lower bound for each is a reduction from propositional validity. Given  $\varphi$  for which we want to determine validity, we let  $\Sigma$  be  $Q \leftrightarrow (\varphi \vee A)$  where  $A$  is a proposition not in  $\varphi$ , and we let  $\mathcal{V} = \emptyset$ . If  $\varphi$  is valid then  $Q$  is equivalent to True for models satisfying  $\Sigma$ , and thus we can reformulate  $Q$  as True. Conversely if  $Q$  can be reformulated over  $\mathcal{V}$ , then the reformulation must be True or False, and  $\Sigma$  ensures that it cannot be False. Thus  $\varphi$  must be valid according to  $\Sigma$  and this easily implies that  $\varphi$  is valid unconditionally.  $\square$

We get the same bounds for  $\forall\text{Horn}$  formulas.

**Proposition 5.** *The problems  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \text{PROP})$ ,  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \text{PROP}^+)$  and  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \text{PROP}^{\text{CQ}})$  as  $Q$  ranges over conjunctions of propositions and  $\Sigma$  over  $\forall\text{Horn}$ , are all CONP-complete.*

*Proof.* First, we show that entailment of propositions via  $\forall\text{Horn}$  theories is CONP hard. Given an arbitrary CNF formula  $\varphi$  we can convert each clause that contains a positive literal into a  $\forall\text{Horn}$  formula. Each “all-negative” clause of the form  $\bigvee \neg P_i$  we convert into a formula  $\bigwedge_i P_i \rightarrow Q'$ . We claim that  $\varphi$  is unsatisfiable if and only if the formulas above imply  $Q'$ . In one direction, if there is a model of the formulas above and  $\neg Q'$ , then clearly the all-negative clauses must hold, and thus  $\varphi$  is satisfiable. Conversely, if there is a model of  $\varphi$ , we extend it to make  $Q'$  false to get a counterexample to entailment.

Next we note that entailment of an atomic proposition  $Q$  from  $\forall\text{Horn}$  formulas can be reduced to reformulation with  $\forall\text{Horn}$  formulas. Given  $\Sigma$  and  $Q$  we simply add  $Q \rightarrow Q'$  for  $Q'$  a fresh symbol, and ask whether  $Q'$  can be reformulated over the vocabulary  $\{Q\}$ . If  $Q'$  can be reformulated, then clearly the reformulation must be  $Q$ , and thus  $\Sigma$  entails  $Q$ . Conversely, if  $\Sigma$  entails  $Q$  then clearly  $Q'$  is equivalent to  $Q$  modulo  $\Sigma$ .  $\square$

When we turn to propositional Horn formulas, entailment becomes P. Using Corollary 1 this gives us P bounds for reformulation over PROP and  $\text{PROP}^+$ . Looking at one proposition in  $\mathcal{V}$  at a time, we get a P bound for reformulating over  $\text{PROP}^{\text{CQ}}$ . In fact, it is easily seen (a variation of [Sagiv and Yannakakis, 1980]) that if a conjunction  $Q$  can be reformulated in  $\text{PROP}^+$  over a theory given by a set of  $\text{PROP}^{\text{HORN}}$  formulas, then it can be reformulated as a conjunction. In summary, we see:

**Proposition 6.** *The problems  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \text{PROP})$ ,  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \text{PROP}^+)$  and  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \text{PROP}^{\text{CQ}})$  as  $Q$  ranges over conjunctions of propositions and  $\Sigma$  over  $\text{PROP}^{\text{HORN}}$ , are all in P.*

**$\exists\text{REF}$  for Full TGDs.** Turning to our predicate logic theories, we begin with the case of TGDs.

We first note that for any class of TGDs  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \text{UCQ})$  and  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \text{CQ})$  are the same. If a UCQ  $Q'$  is equivalent to  $Q$  modulo TGDs  $\Sigma$ ,  $Q'$  must hold in the chase of  $\Sigma$ , and thus one of its disjuncts must hold in the chase; therefore one of the disjuncts must be equivalent to  $Q$  modulo  $\Sigma$ . We make use of this implicitly in the arguments below.

**Proposition 7.** *The problems  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \text{FO})$ ,  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \text{UCQ})$  and  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \text{CQ})$  as  $Q$  ranges over CQs and  $\Sigma$  over FullTGD are all EXPTIME-complete. When the arity is fixed, the problems become NP complete.*

*Proof.* The reduction to entailment gives the EXPTIME and NP upper bounds, since chasing with full TGDs is simply a fixpoint computation as in Datalog evaluation.

Before turning to the lower bounds, we note that the first two problems do not differ in their answers: if a CQ has an FO reformulation with respect to Full TGDs, then it has a monotone one. This is shown by considering the entailments produced by the “Craig reduction” and “Lyndon reduction” mentioned before. They differ only by addition of the axioms BackAx. But in performing the chase process, there will never be a need to fire the BackAx constraints, since no new facts will be produced. Thus in our lower bound, we will only deal with the monotone and CQ cases.

Conversely, we can reduce from Datalog evaluation to get EXPTIME-completeness and NP-completeness. Consider a Datalog program  $P$  with goal  $G()$  which we want to evaluate on an input instance  $D$  over constants  $\{d_1 \dots d_k\}$ . Let  $\text{DOM}$  be a new relation of arity  $k$ . We let  $\Sigma$  be formed by:

- turning each rule of  $P$  into a full TGD,
- adding the rules  $\text{DOM}(x_1 \dots x_k) \rightarrow A(x_{j_1} \dots x_{j_k})$  for each fact  $A_i(d_{j_1} \dots d_{j_k}) \in D$ ,
- adding the rule  $\text{DOM}(x_1 \dots x_k) \wedge G() \rightarrow G'(x_1 \dots x_k)$ ,
- adding the rule  $G'(x_1 \dots x_k) \rightarrow \text{DOM}(x_1 \dots x_k)$ .

Let  $\mathcal{V} = \{G'\}$  and consider the query

$$Q = \exists v_1 \dots v_k \text{ DOM}(v_1 \dots v_k)$$

We claim that the goal predicate of  $P$  is true on  $D$  if and only if  $Q$  has a  $\Sigma$ -reformulation over  $\mathcal{V}$ .

In one direction, suppose the goal of  $P$  evaluates to true on  $D$ . We claim that  $Q' = \exists v_1 \dots v_k G'(v_1 \dots v_k)$  is a  $\Sigma$ -reformulation of  $Q$  over  $\mathcal{V}$ . If  $Q'$  holds on an instance  $I$  satisfying  $\Sigma$  then the last rule implies that  $Q$  holds. Conversely if  $Q$  holds then the facts corresponding to  $D$  are populated in  $I$  by the first set of rules; and thus  $G()$  holds using the second set of rules. Finally  $Q'$  holds using the last rule.

Conversely suppose  $Q$  has a monotone reformulation  $Q'$  over  $\mathcal{V}$ . Let  $I_0$  be the canonical database of  $Q$ , consisting of a single fact  $\text{DOM}(c_1 \dots c_k)$ . It is clear that  $Q'$  must have some atom containing  $G'$ , since  $Q$  is neither a tautology or a contradiction relative to  $\Sigma$ . Thus in particular  $Q'$  must be derived from  $I_0$  using the rules of  $\Sigma$ . Clearly a derivation must use the third rule to generate the  $G'$  atom. But this can only happen if  $G()$  is derived. From a derivation of  $G()$  from the fact  $\text{DOM}(c_1 \dots c_k)$  we can read off a derivation of the goal of  $P$  on  $D$  as required.  $\square$

**Weakly-acyclic TGDs.** We now consider weakly acyclic TGDs. Entailment of CQs under these sentences is 2EXPTIME-complete [Cali *et al.*, 2010], even when the arity of the schema is fixed. Thus using the Craig and Lyndon reductions, we get that  $\exists \text{REF}(Q, \Sigma, \mathcal{V}, \text{FO})$  and  $\exists \text{REF}(Q, \Sigma, \mathcal{V}, \text{UCQ})$  are in 2EXPTIME. A matching lower bound follows by a reduction from entailment, similar to those above.

**Full disjunctive TGDs.** For arbitrary FullDisTGD fact entailment is CONEXPTIME-complete. Although we imagine this has been known before, we provide an argument here. Note that the combined complexity of the “possibility problem” for Disjunctive Datalog has been examined in detail in [Eiter *et al.*, 1997]; but this does not co-incide either with entailment or its negation.

**Proposition 8.** *The problems  $\exists \text{REF}(Q, \Sigma, \mathcal{V}, \text{FO})$  and  $\exists \text{REF}(Q, \Sigma, \mathcal{V}, \text{UCQ})$  as  $Q$  ranges over CQs and  $\Sigma$  over FullDisTGD are all CONEXPTIME-complete. When the arity is fixed, the problems become  $\Pi_2^p$  complete.*

*Proof.* The CONEXPTIME upper bound is simply by universally choosing disjuncts in the head to fire and then proceeding until a fixpoint is reached. The lower bound is by direct reduction from the acceptance problem for a NEXPTIME machine. We give a sketch of the argument here. Consider a non-deterministic machine  $M$  running in time  $2^{p(m)}$  for polynomial  $p$  and an input  $w$  of length  $m$  to the complement. Letting  $n = p(m)$ , we code the value on a tape of such a machine with relations:

- $\text{IsOne}(t_1, \dots, t_n, s_1, \dots, s_n)$  representing that at time  $t_1 \dots t_n$ , the cell at index  $s_1 \dots s_n$  has value 1, and similarly  $\text{IsZero}(\vec{t}, \vec{s})$  saying that the cell has value 0,
- $\text{State}_q(t_1, \dots, t_n)$  asserting that the state at time  $t_1 \dots t_n$  is  $q$ ,
- $\text{Head}(t_1, \dots, t_n, s_1, \dots, s_n)$  stating that at time  $t_1 \dots t_n$  the head is on  $s$ ,
- $\text{IsOne}(x)$  and  $\text{IsZero}(x)$  stating that bit  $x$  represents 1 and 0, respectively.

We will have relations  $\text{AllEq}_i(x_1, \dots, x_i, y_1, \dots, y_i)$  representing that  $x_1, \dots, x_i$  and  $y_1, \dots, y_i$  agree on all  $i$ . The following rules will generate the needed equalities, assuming that we begin with two distinct elements inhabiting  $\text{IsOne}$  and  $\text{IsZero}$ :

$$\begin{aligned} & \text{IsOne}(x) \wedge \text{IsOne}(y) \vee \text{IsZero}(x) \vee \text{IsZero}(y) \rightarrow \text{AllEq}_1(x, y), \\ & \text{AllEq}_i(x_1, \dots, x_i, y_1, \dots, y_i) \wedge [\text{IsOne}(x_{i+1}) \wedge \text{IsOne}(y_{i+1}) \vee \text{IsZero}(x_{i+1}) \wedge \text{IsZero}(y_{i+1})] \rightarrow \\ & \quad \text{AllEq}_{i+1}(x_1, \dots, x_{i+1}, y_1, \dots, y_{i+1}). \end{aligned}$$

Using this we have successor relation  $\succ(x_1, \dots, x_n, y_1, \dots, y_n)$  and a set of rules, for all  $i$ , that code its semantics:

$$\begin{aligned} & \bigwedge_{i \leq j} \text{IsOne}(x_i) \wedge \text{IsZero}(x_{j+1}) \wedge \bigwedge_{i \leq j} \text{IsZero}(y_1, \dots, y_i) \\ & \wedge \text{IsOne}(y_{j+1}) \wedge \text{AllEq}_{n-(j+2)}(x_{j+2} \dots x_n, y_{j+2}, \dots, y_n) \\ & \rightarrow \succ(x_1, \dots, x_n, y_1, \dots, y_n). \end{aligned}$$

Now given an exponentially long successor relation we can easily write disjunctive rules stating that the configurations satisfy the transition relations.

When the arity is fixed entailment becomes  $\Pi_2^p$ -complete by an analogous argument.

Thus the Craig reduction gives CONEXPTIME and  $\Pi_2^p$  upper bounds for  $\exists \text{REF}(Q, \Sigma, \mathcal{V}, \text{FO})$  and  $\exists \text{REF}(Q, \Sigma, \mathcal{V}, \text{UCQ})$ , with the lower bounds coming via reduction from entailment.  $\square$

We now consider  $\exists \text{REF}(Q, \Sigma, \mathcal{V}, \text{CQ})$ , which we claim is undecidable.

**Proposition 9.** *The problem  $\exists \text{REF}(Q, \Sigma, \mathcal{V}, \text{CQ})$  as  $Q$  ranges over CQs and  $\Sigma$  over FullDisTGD is undecidable.*

*Proof.* We show undecidability of a simpler problem: determining if there exists a CQ  $Q'$  that is a  $\Sigma$ -reformulation over  $\mathcal{V}$  of  $Q$ , where  $\Sigma$  consists of full TGDs and  $Q$  is a UCQ. This problem can easily be reduced to  $\exists\text{REF}(Q, \Sigma, \mathcal{V}, \text{CQ})$ , by creating a single disjunctive rule relating  $Q$  to its disjuncts. We show that the problem above is undecidable via a reduction from Datalog containment, shown undecidable in [Shmueli, 1993].

Given  $P_1$  and  $P_2$  in Datalog with (nullary) goal predicates  $G_1()$  and  $G_2()$ , respectively, we consider Boolean  $Q = (G() \wedge G_1()) \vee G_2()$  where  $G$  is a new nullary relation, and our full TGDs  $\Sigma$  correspond to the two Datalog programs. We claim that  $Q$  can be reformulated modulo  $\Sigma$  to a CQ (over any vocabulary) iff  $P_1$  is contained in  $P_2$ . Clearly if  $P_1$  is contained in  $P_2$ , then  $Q$  is equivalent to the CQ  $G_2()$  modulo  $\Sigma$ . In the other direction, if  $Q$  is equivalent to a CQ  $Q'$  then  $Q'$  must be equivalent to one of the disjuncts of  $Q$  modulo  $\Sigma$ : as mentioned earlier, this is because  $Q$  must hold in the chase of the canonical database of  $Q'$ , and hence one of the disjuncts of  $Q$  must hold in the chase. If  $Q'$  is equivalent to  $G() \wedge G_1()$ , then we know that  $Q$  is equivalent to  $G() \wedge G_1()$  modulo  $\Sigma$ . But this would mean that  $G_2$  implies  $G() \wedge G_1()$  modulo  $\Sigma$ , which is impossible since the relation  $G$  is never mentioned in a rule of  $\Sigma$ . Hence we must have  $Q'$  is equivalent to  $G_2()$  modulo  $\Sigma$ . This means that  $G() \wedge G_1()$  must imply  $G_2()$  modulo  $\Sigma$ . Since  $G$  is fresh this means that  $G_1()$  must imply  $G_2()$  modulo  $\Sigma$ . From this we see that the Datalog program given  $P_1$  is contained in  $P_2$ .  $\square$

## Proofs for $\text{REF}^\leq$

The following proposition justifies the first part of Proposition 2 for each case but  $\text{PROP}^{\text{CQ}}$ , as well as the corresponding lower bounds.

**Proposition 10.** *The problems  $\text{REF}^\leq(Q, \Sigma, \mathcal{V}, \text{PROP}, k)$ ,  $\text{REF}^\leq(Q, \Sigma, \mathcal{V}, \text{PROP}^+, k)$  and  $\text{REF}^\leq(Q, \Sigma, \mathcal{V}, \text{PROP}^{\text{UCQ}}, k)$  as  $Q$  ranges over conjunctions of propositions and  $\Sigma$  over propositional sentences, are  $\Sigma_2^P$ -complete.*

*Proof.* We first consider  $\text{REF}^\leq(Q, \Sigma, \mathcal{V}, \text{PROP}, k)$ . [Buchfuhrer and Umans, 2011] shows that minimization of a propositional formula in DNF within the set of all formulas is  $\Sigma_2^P$ -hard. We can reduce this problem to  $\text{REF}^\leq(Q, \Sigma, \mathcal{V}, \text{PROP}, k)$  by producing, for a DNF  $\varphi = \bigvee_{i \leq n} \varphi_i$ , the following:

- query  $Q$ , where  $Q$  is a new symbol,
- $\Sigma$  consisting of  $Q \leftrightarrow R_i$  and  $R_i \leftrightarrow \varphi_i$ , for all  $i$ ,
- target alphabet consisting of the literals of  $\varphi$ .

It is clear that reformulations of  $Q$  correspond to formulas equivalent to  $\varphi$ .

For the second, we perform the same reduction, but we add also copies  $P'_i$  of each literal  $P_i$  in  $\varphi$ , adding to  $\Sigma$  additional formulas  $\neg P_i \leftrightarrow P'_i$ , and take our target vocabulary to use both copies of the literals. Thus formulas equivalent to  $\varphi$  correspond to monotone reformulations of  $Q$  of the same size.

For the third we make use of the fact that finding a minimal equivalent DNF is  $\Sigma_2^P$ -hard [Umans, 2001], and the “monotonization” reduction above.  $\square$

The upper bound for  $\text{PROP}^{\text{CQ}}$  in Proposition 2 is proven using the same technique as above. The lower bound will follow from the corresponding lower bound for disjunctive Horn constraints proven below.

The following proposition justifies the second part of Proposition 2 (as well as the corresponding upper bounds).

**Proposition 11.** *The problems  $\text{REF}^\leq(Q, \Sigma, \mathcal{V}, \text{PROP}, k)$ ,  $\text{REF}^\leq(Q, \Sigma, \mathcal{V}, \text{PROP}^+, k)$ ,  $\text{REF}^\leq(Q, \Sigma, \mathcal{V}, \text{PROP}^{\text{UCQ}}, k)$ ,  $\text{REF}^\leq(Q, \Sigma, \mathcal{V}, \text{PROP}^{\text{CQ}}, k)$  as  $Q$  ranges over conjunctions of propositions and  $\Sigma$  over propositional Horn sentences, are NP-complete.*

*Proof.* We reduce from the set cover problem, which has input a family  $S_1 \dots S_n$  of finite sets and a number  $k$ , and asked whether there are at most  $k$  indices  $j$  such that the union of  $S_j$  covers each element in  $\bigcup_{1 \leq i \leq n} S_i$ .

We start by giving a reduction that works for the first 3 problems above. Given an instance of set cover as above, let  $\{X_1 \dots X_d\}$  be the union of all the  $S_i$ . Our reduction produces formulas over propositions  $X_i$ , for  $1 \leq i \leq d$ ,  $S_i$ , for  $1 \leq i \leq n$ , and a symbol  $Q$ . The theory  $\Sigma$  consists of:

$$\begin{aligned} Q &\leftrightarrow \bigvee_{0 \leq i \leq d} X_i, \\ S_i &\leftrightarrow \bigvee_{X_\ell \in S_i} X_\ell. \end{aligned}$$

The signature  $\mathcal{V}$  is  $\{S_1, \dots, S_n\}$ . From a set cover with  $S_{j_1} \dots S_{j_k}$  we can use the reformulation  $\bigvee_i S_{j_i}$ , and from a reformulation with at most  $k$  sets  $S_i$ , we can read off a set cover with the same  $S_i$ .

In the above problems we used reformulations that used disjunction. For the fourth problem, we need to consider reformulations that use only conjunction. However we can “dualize” the problem above. We have extra propositions  $R$  and  $S'_i$ , for  $1 \leq i \leq n$ , and add to  $\Sigma$  formulas

$$\begin{aligned} R &\leftrightarrow \neg Q, \\ S'_i &\leftrightarrow \neg S_i. \end{aligned}$$

A conjunctive reformulation of  $R$  with the  $S'_i$  will correspond to a disjunctive reformulation of  $Q$  with the  $S_i$ .  $\square$

We now give the proof of hardness with  $\forall$ Horn formulas.

*Proof of Theorem 1.* Let us first recall the statement of the theorem:

$\text{REF}^{\leq}(Q, \Sigma, \mathcal{V}, \text{PROP}^+, k)$ ,  $\text{REF}^{\leq}(Q, \Sigma, \mathcal{V}, \text{PROP}^{\text{UCQ}}, k)$  and  $\text{REF}^{\leq}(Q, \Sigma, \mathcal{V}, \text{PROP}^{\text{CQ}}, k)$  are all  $\Sigma_2^P$ -hard when  $\Sigma$  ranges over  $\forall$ Horn formulas.

For simplicity, we assume that the size of a positive propositional formula is the number of occurrences of propositional variables in this formula.

The proof is by reduction of  $\exists\forall 3\text{SAT}$  problem. Suppose  $\varphi = \exists \bar{x} \forall \bar{y} \neg \psi$  is formula where  $\bar{x} = x_1, \dots, x_k$  and  $\bar{y} = y_1, \dots, y_m$  are propositional variables, and  $\psi$  is a formula in 3CNF—that is, a conjunction of clauses  $\ell^1 \vee \ell^2 \vee \ell^3$  with each  $\ell^j$  a variable from  $\bar{x} \cup \bar{y}$  or the negation of such a variable.

We construct a set of propositional disjunctive Horn clauses  $\Sigma$  and a set of propositional variables  $\mathcal{V}$  such that  $\varphi$  is valid if and only if there exists a positive formula  $Q'$  over  $\mathcal{V}$  of size  $k$  such that  $Q'$  is equivalent to a propositional variable  $Q$  over  $\Sigma$ .

In particular, the vocabulary  $\mathcal{V}$  consists of variables  $\text{TrueX}_i$  and  $\text{FalseX}_i$  for each  $x_i$ .

On the high level, the clauses in  $\Sigma$  guarantee that

1.  $Q$  implies all the variables in  $\mathcal{V}$ , so the equivalence of  $Q$  and  $Q'$  over  $\Sigma$  boils down to the containment of  $Q'$  in  $Q$  over  $\Sigma$ ;
2. containment of  $Q'$  in  $Q$  over  $\Sigma$  is only possible if  $Q'$  implies either  $\text{TrueX}_i$  or  $\text{FalseX}_i$  for each  $x_i$ ; since the size of  $Q'$  is bounded by  $k$ , it means that  $Q'$  can only be of the form  $\text{AsX}_1 \wedge \dots \wedge \text{AsX}_k$ , for  $\text{AsX}_i \in \{\text{TrueX}_i, \text{FalseX}_i\}$ —that is,  $Q'$  corresponds to an assignment of variables  $x_i$ ;
3. if a dataset and  $\Sigma$  imply  $\text{TrueX}_i$  or  $\text{FalseX}_i$  for each  $x_i$ , then they imply  $\text{TrueY}_i$  or  $\text{FalseY}_i$  for each  $y_i$  in  $\bar{y}$ , which corresponds to an assignment of the universally quantified variables in  $\varphi$ ;
4. if a chase of a dataset implies a *falsification* of a clause in  $\psi$  under this encoding, which corresponds to falsification of the whole  $\psi$ , then it implies  $Q$  as well.

In the following formal definition of  $\Sigma$  we use disjunction on the left of implication and conjunction on the right for brevity. The first group of  $\Sigma$ , corresponding to the guarantee 1 above, consists of the rules

$$Q \rightarrow \text{TrueX}_i \wedge \text{FalseX}_i, \quad \text{for each } i = 1, \dots, k. \quad (1)$$

The second group, that corresponds to the guarantee 2 together with the next groups, consists of the rules

$$\begin{aligned} \text{TrueX}_1 \vee \text{FalseX}_1 &\rightarrow \text{XAssignedUntil}_1, \\ \text{XAssignedUntil}_{i-1} \wedge (\text{TrueX}_i \vee \text{FalseX}_i) &\rightarrow \text{XAssignedUntil}_i, \\ &\quad \text{for each } i = 2, \dots, k, \\ \text{XAssignedUntil}_k &\rightarrow \text{XAssignedAll}. \end{aligned} \quad (2)$$

The third group, corresponding to guarantee 3, consists of the rules

$$\text{XAssignedAll} \rightarrow \text{TrueY}_j \vee \text{FalseY}_j, \quad \text{for each } j = 1, \dots, m. \quad (3)$$

Finally, the fourth group, corresponding to guarantee 4 (and a part of the guarantee 2), consists of the rules

$$\text{XAssignedAll} \wedge \text{AsZ}^1 \wedge \text{AsZ}^2 \wedge \text{AsZ}^3 \rightarrow Q, \quad (4)$$

for each clause  $\ell^1 \vee \ell^2 \vee \ell^3$  in  $\psi$ , where, for each  $j = 1, 2, 3$ ,  $\text{AsZ}^j$  is

- $\text{TrueX}_i$  if  $\ell^j$  is a negation of an existential variable  $x_i \in \bar{x}$ ,
- $\text{FalseX}_i$  if  $\ell^j$  is an existential variable  $x_i \in \bar{x}$ ,
- $\text{TrueY}_i$  if  $\ell^j$  is a negation of a universal variable  $y_i \in \bar{y}$ ,
- $\text{FalseY}_i$  if  $\ell^j$  is a universal variable  $y_i \in \bar{y}$ .

Next we give a formal proof of the correctness of the construction.

Suppose first that  $\varphi$  is valid. We need to prove that there exists  $Q'$  equivalent to  $Q$ . To this end, consider an assignment of  $\bar{x}$  witnessing the validity of  $\varphi$  and let  $Q' = \text{AsX}_1 \wedge \dots \wedge \text{AsX}_k$ , where, for each  $i = 1, \dots, k$ ,  $\text{AsX}_i = \text{TrueX}_i$  if  $x_i$  evaluates to *true* by the assignment and  $\text{AsX}_i = \text{FalseX}_i$  otherwise. Rules (1) guarantee that any dataset that contains  $Q$  and satisfies  $\Sigma$  contains also all variables in  $Q'$ . Hence we only need to prove that any dataset that implies  $Q'$  and satisfies  $\Sigma$  contains also  $Q$ . To this end, consider a dataset  $D$  that contains all variables in  $Q'$  and satisfies  $\Sigma$ . Rules (2) guarantee that  $D$  contains  $\text{XAssignedAll}$ . Hence, rules (3) guarantee that  $D$  contains also  $\text{TrueY}_i$  or  $\text{FalseY}_i$  for each  $y_i$  in  $\bar{y}$ . Hence, overall,  $D$  encodes one or more extensions of the assignment of  $\bar{x}$  to  $\bar{y}$ . Fix one of these extensions. Since  $\varphi$  is valid, there exists a clause in  $\psi$  that evaluates to *false* under the overall assignment of  $\bar{x} \cup \bar{y}$ . Therefore, rule (4), corresponding to this clause, guarantees that  $D$  contains  $Q$ , as required.

Suppose now that there exists a positive query  $Q'$  of size at most  $k$  equivalent to  $Q$ . We need to prove that  $\varphi$  is valid.

First, we prove the following claim.

**Claim 1.** *Query  $Q'$  implies  $\text{TrueX}_i \vee \text{FalseX}_i$  for all  $i \in \{1, \dots, k\}$  (the implication here does not take  $\Sigma$  into account— that is, we claim that  $Q' \rightarrow \text{TrueX}_i \vee \text{FalseX}_i$ , not  $Q', \Sigma \rightarrow \text{TrueX}_i \vee \text{FalseX}_i$ ).*

*Proof.* Assume, for the sake of contradiction, that this is not the case. Consider the smallest  $i$  such that  $Q'$  does not imply  $\text{TrueX}_i \vee \text{FalseX}_i$ , and form the dataset

$$D = \{\text{TrueX}_j, \text{FalseX}_j \mid j \neq i\} \cup \{\text{XAssignedUntil}_j \mid j < i\}.$$

On the one hand, one can see by inspection that  $D$  satisfies  $\Sigma$ . We claim that  $D$  satisfies  $Q'$  as well. Indeed,  $Q'$  can be decomposed as a disjunction where each disjunct represents which of  $\text{TrueX}_i, \text{FalseX}_i$  is true:

$$\begin{aligned} &[(\text{TrueX}_i \wedge \text{FalseX}_i) \wedge \gamma_1] \\ &\vee [\text{TrueX}_i \wedge \gamma_2] \\ &\vee [\text{FalseX}_i \wedge \gamma_3] \\ &\vee \gamma_4, \end{aligned}$$

where the  $\gamma_1, \dots, \gamma_4$  are formulas that are either *false* or monotone formulas in  $\text{TrueX}_j, \text{FalseX}_j$  for  $j \neq i$ , such that at least one of  $\gamma_1, \dots, \gamma_4$  is not *false*. If  $\gamma_4$  is *false* then  $Q'$  implies  $\text{TrueX}_i \vee \text{FalseX}_i$ , contrary to assumption. But if  $\gamma_4$  is a monotone formula in  $\text{TrueX}_j, \text{FalseX}_j$  for  $j \neq i$  then  $D$  must satisfy  $\gamma_4$  and hence must satisfy overall  $Q'$ .

On the other hand,  $D$  does not contain the variable  $Q$ —that is, does not imply  $Q$  as a query. This contradicts the fact that  $Q'$  and  $Q$  are equivalent for instances satisfying the formulas, because  $D$  is a witness for non-equivalence.  $\square$

Since by Claim 1 query  $Q'$  implies  $\text{TrueX}_i \vee \text{FalseX}_i$  for all  $i$   $Q'$  is either *True* or must mention, for each  $i$ , one of  $\text{TrueX}_i \vee \text{FalseX}_i$ . In the case  $Q'$  is *True*, the conclusion is obvious. In the other case, since the size of  $Q'$  is bounded by  $k$ , we can conclude that  $Q' = \text{AsX}_1 \wedge \dots \wedge \text{AsX}_k$ , for  $\text{AsX}_i \in \{\text{TrueX}_i, \text{FalseX}_i\}$ .

Consider the assignment of  $\bar{x}$  corresponding to  $Q'$ . Next we prove that  $\psi$  evaluates to *false* under any extension of this assignment to  $\bar{y}$ . Consider such an extension and a dataset  $D$  consisting of

- all the variables in  $Q'$ ;
- $\text{XAssignedUntil}_i$  for all  $i = 1, \dots, k$ ;
- $\text{XAssignedAll}$ ;
- for all  $i = 1, \dots, m$ ,  $\text{TrueY}_i$  if the extension assigns  $y_i$  to *true* and  $\text{FalseY}_i$  otherwise.

Let  $D'$  be the result of chasing  $D$  under all the non-disjunctive rules. Such  $D'$  satisfies all the rules of  $\Sigma$ , since the disjunctive rules was already satisfied by construction. Also,  $D'$  implies  $Q'$ , so the chase contains  $Q$  as well, because  $Q'$  and  $Q$  are equivalent over  $\Sigma$ . By rules (4), which are the only rules implying  $Q$ , this is possible only if there is a clause in  $\psi$  that evaluates to *false* under the overall assignment encoded by  $D$ . So  $\psi$  evaluates to *false* as well, as required.  $\square$

**Full TGDs.** When background knowledge is captured using sentences that are full TGDs, entailment is in EXPTIME and in NP if the arity is bounded. Thus  $\text{REF}^{\leq}(Q, \Sigma, \mathcal{V}, \text{UCQ}, k)$  all can be solved in EXPTIME and in NP using the guess-and-check method explained earlier. The matching hardness follows immediately from the reduction from entailment.

**Weakly-acyclic TGDs.** For weakly-acyclic TGDs, entailment is in 2EXPTIME even for fixed arity. Since guessing a reformulation of size  $k$  can also be simulated in a 2EXPTIME machine,  $\text{REF}^{\leq}(Q, \Sigma, \mathcal{V}, \text{UCQ}, k)$  and  $\text{REF}^{\leq}(Q, \Sigma, \mathcal{V}, \text{CQ}, k)$  are both in 2EXPTIME just by guessing the reformulation and checking it again. Since we always have a lower bound from entailment, the bounds are tight.

**Full disjunctive TGDs.** We now turn to our results for full disjunctive TGDs. For arbitrary arity relations, we have a CONEXPTIME upper bound by the “guessing-and-check” method. The complexity of “guessing” is subsumed by the

CONEXPTIME complexity of checking the entailment. Since entailment with full disjunctive TGDs is CONEXPTIME-hard, we have a matching lower bound. For bounded arity guess-and-check entailment is in  $\Pi_2^p$ , and thus guess-and-check gives us a  $\Sigma_3^p$  upper-bound. The main difficulty will be proving a matching lower bound for bounded arity.

*Proof of Theorem 2.* Let us first recall the statement of the theorem:

$\text{REF}^{\leq}(Q, \Sigma, \mathcal{V}, \mathcal{L}, k)$  for any  $\mathcal{L} \in \{\text{FO}^+, \text{UCQ}, \text{CQ}\}$  is  $\Sigma_3^p$ -hard when  $\Sigma$  ranges over sets of FullDisTGD formulas and the arity of relations is bounded by 2.

For simplicity, we assume that the size of a positive first-order formula is the number of atoms in this formula (the predicates that appear in the reduction are all at most binary).

The proof is by reduction of  $\exists\forall\exists 3\text{SAT}$  problem. Suppose  $\varphi = \exists\bar{x}\forall\bar{y}\exists\bar{z}\psi$  is a formula where  $\bar{x} = x_1, \dots, x_n$ ,  $\bar{y}$  and  $\bar{z}$  are propositional variables, and  $\psi$  is a formula in 3CNF—that is, a conjunction of clauses  $\ell^1 \vee \ell^2 \vee \ell^3$  with each  $\ell^j$  a variable from  $\bar{x} \cup \bar{y} \cup \bar{z}$  or the negation of such a variable.

For each clause  $\gamma = \ell^1 \vee \ell^2 \vee \ell^3$  in  $\psi$  let  $\alpha_\gamma^j, j = 1, \dots, 7$ , be an enumeration of the assignments to variables in  $\ell^1, \ell^2$  and  $\ell^3$  that satisfy the clause.

Let  $\sigma$  be the conjunction of

$$\forall \text{True}_s(t_s) \wedge \forall \text{False}_s(f_s)$$

for all  $s \in \bar{x} \cup \bar{y} \cup \bar{z}$ , and

$$\text{ClAssignment}_\gamma^j(v_\gamma^j) \wedge \text{Clause}_\gamma(v_\gamma^j) \wedge \text{Arg}^1(v_\gamma^j, u_\gamma^{j,1}) \wedge \text{Arg}^2(v_\gamma^j, u_\gamma^{j,2}) \wedge \text{Arg}^3(v_\gamma^j, u_\gamma^{j,3}),$$

where  $u_\gamma^{j,i}, i = 1, 2, 3$ , is  $t_s$  if the propositional variable  $s$  in  $\ell^i$  is assigned to *true* by  $\alpha_\gamma^j$  and  $f_s$  if is assigned to *false*. So, overall,  $\sigma$  contains 7 variables  $v_\gamma^j$  for each clause  $\gamma$  and 2 variables  $t_s$  and  $f_s$  for each propositional variable  $s$ . We also denote  $\bar{v}_\sigma$  the tuple of all variables of  $\sigma$ , that is,  $\{v_\gamma^1, \dots, v_\gamma^7, \mid \gamma \text{ in } \psi\} \cup \{t_s, f_s \mid s \in \bar{x} \cup \bar{y} \cup \bar{z}\}$ .

So above  $v_\gamma^i$  is a distinct variable corresponding to each assignment to clause  $\gamma$ . The rules of  $\sigma$  intuitively encode the semantics of the clause.

Let  $\xi$  be the conjunction of

$$\text{XAssigned}(t_x) \wedge \text{XAssigned}(f_x)$$

for all  $x \in \bar{x}$ .

Let  $\tau$  be the conjunction of

$$\text{Clause}_\gamma(w_\gamma) \wedge \text{Arg}^1(w_\gamma, w_{s^1}) \wedge \text{Arg}^2(w_\gamma, w_{s^2}) \wedge \text{Arg}^3(w_\gamma, w_{s^3}),$$

for each clause  $\gamma = \ell^1 \vee \ell^2 \vee \ell^3$  in  $\psi$ , where  $s^i \in \bar{x} \cup \bar{y} \cup \bar{z}, i = 1, 2, 3$ , is the variable in  $\ell^i$ .

Finally, let  $\pi$  be the conjunction of

$$\begin{aligned} \text{XAssigned}(w_x), & \quad \text{for all } x \in \bar{x}, \\ \text{YAssigned}(w_y), & \quad \text{for all } y \in \bar{y}, \text{ and} \\ \text{ZAssigned}(w_z), & \quad \text{for all } z \in \bar{z}. \end{aligned}$$

Observe that the variables of  $\tau$  and  $\pi$  are disjoint from those of  $\sigma$ .

Then, let  $\Sigma$  consist of the implication:

$$\sigma \wedge \tau \wedge \pi \rightarrow \xi, \tag{5}$$

the implications

$$\begin{aligned} \sigma \wedge (\text{XAssigned}(t_{x_1}) \vee \text{XAssigned}(f_{x_1})) & \rightarrow \text{XAssignedUntil}_1(), \\ \sigma \wedge \text{XAssignedUntil}_{i-1}() \wedge \\ (\text{XAssigned}(t_{x_i}) \vee \text{XAssigned}(f_{x_i})) & \rightarrow \text{XAssignedUntil}_i(), \\ & \quad \text{for each } i = 2, \dots, n, \\ \text{XAssignedUntil}_n() & \rightarrow \text{XAssignedAll}(), \end{aligned} \tag{6}$$

and the implications

$$\sigma \wedge \text{XAssignedAll}() \rightarrow \text{YAssigned}(t_y) \vee \text{YAssigned}(f_y), \quad \text{for all } y \in \bar{y}, \tag{7}$$

$$\sigma \wedge \text{XAssignedAll}() \rightarrow \text{ZAssigned}(t_z) \wedge \text{ZAssigned}(f_z), \quad \text{for all } z \in \bar{z}. \tag{8}$$

Finally, let  $Q$  be the Boolean CQ that is the existential closure of  $\sigma \wedge \tau \wedge \pi$ . Let

$$\begin{aligned} \mathcal{V} = \{ & \text{ClAssignment}_\gamma^1, \dots, \text{ClAssignment}_\gamma^7, \text{Clause}_\gamma \mid \gamma \text{ in } \psi\} \cup \\ & \{\forall \text{True}_s, \forall \text{False}_s \mid s \in \bar{x} \cup \bar{y} \cup \bar{z}\} \cup \{\text{Arg}^1, \text{Arg}^2, \text{Arg}^3, \text{XAssigned}\}, \end{aligned}$$

Let  $k = |\sigma| + n$ , where  $|\sigma|$  is the number of atoms in  $\sigma$ ; without loss of generality, we assume that  $k$  is the restriction on the number of atoms in  $Q'$ .

Having completed the construction, we now prove that  $\varphi$  is valid if and only if there exists a positive existential first order Boolean query  $Q'$  of size at most  $k$  that is equivalent to  $Q$  under  $\Sigma$ .

For the forward direction, suppose that  $\varphi$  is valid. That is, there exists an assignment of variables  $\bar{x}$  such that for all assignments of  $\bar{y}$  there is an assignment of  $\bar{z}$  with  $\psi$  evaluating to *true* by the overall assignment. Consider the query  $Q'$  defined as follows:

$$\exists \bar{v}_\sigma (\sigma \wedge \text{XAssigned}(u_{x_1}) \wedge \dots \wedge \text{XAssigned}(u_{x_n})),$$

where each  $u_{x_i}$ ,  $i = 1, \dots, n$ , is  $t_{x_i}$ , if  $x_i$  is assigned to *true* by the assignment, and  $f_{x_i}$ , otherwise. First, note that  $Q'$  has precisely  $k$  atoms. Therefore, it is enough to prove that  $Q'$  is equivalent to  $Q$  under  $\Sigma$ . Suppose for the sake of contradiction that this is not the case. We argue that  $Q$  and  $\Sigma$  implies  $Q'$ . Observe that the body of  $Q$ , that is  $\sigma \wedge \tau \wedge \pi$ , is precisely the body of rule (5), so this rule fires. The head of the rule (5),  $\xi$ , is the “assignment” of all variables to both *true* and *false* is  $\sigma$ , while  $Q'$  conjoins  $\sigma$  with atoms in which only one of *true* and *false* assigned for each  $x$ . Thus  $\xi$  implies  $Q'$ .

Therefore, there exists a database  $D$  that satisfies both  $Q'$  and  $\Sigma$ , but does not satisfy  $Q$ . Since  $D$  satisfies  $Q'$ , there exists a homomorphism  $h$  from  $Q'$  to  $D$ . Since  $D$  satisfies  $\Sigma$ , by rules (6) it contains  $\text{XAssignedUntil}_i()$  for all  $i$  and  $\text{XAssignedAll}()$ . Therefore, by rules (8),  $D$  contains either  $\text{YAssigned}(h(t_y))$  or  $\text{YAssigned}(h(f_y))$  for all  $y \in \bar{y}$  and both  $\text{ZAssigned}(h(t_z))$  and  $\text{ZAssigned}(h(f_z))$  for all  $z \in \bar{z}$ . Consider any extension of the assignment of  $\bar{x}$  to  $\bar{y}$  that agrees with  $D$ , in the sense that  $\text{YAssigned}(h(t_y))$  holds in  $D$  if  $y$  is assigned to *true* and  $\text{YAssigned}(h(f_y))$  holds otherwise. Since there exists an extension of the assignment to  $\bar{z}$  such that  $\psi$  evaluates to *true*, it is straightforward to construct a homomorphism from  $\tau \cup \pi$  to  $D$ . However, there is also a homomorphism from  $\sigma$  to  $D$  by construction. Since  $\tau \cup \pi$  and  $\sigma$  do not have any variables in common, there is a homomorphism from the whole of  $Q$  to  $D$ . This contradicts the fact  $Q$  does not hold in  $D$ . Therefore our assumption was wrong, and  $Q$  is indeed equivalent to  $Q'$  under  $\Sigma$ , as required.

For the backward direction of the proof, suppose that there exists a positive query  $Q'$  of size at most  $k$  equivalent to  $Q$  relative to  $\Sigma$ . We need to prove that  $\varphi$  is valid.

We make use of the following claims.

**Claim 2.** Consider any  $Q'$  that is equivalent to  $Q$  relative to the constraints  $\Sigma$  (regardless of size). For each  $i = 1, \dots, n$ ,  $Q' \subseteq Q_i$ , where

$$Q_i = \exists \bar{v}_\sigma (\sigma \wedge (\text{XAssigned}(t_{x_i}) \vee \text{XAssigned}(f_{x_i}))).$$

*Proof.* Assume, for the sake of contradiction, that this is not the case, and there is a number  $i$  with  $Q' \not\subseteq Q_i$ . Convert  $Q'$  to a UCQ and consider a CQ  $Q''$  in this UCQ such that  $Q'' \not\subseteq Q_i$ . Let  $D$  be the database that is the result of replacing each variable  $v$  in  $Q''$  by a fresh constant  $a_v$ . Let  $D'$  be one of the outputs of the disjunctive chase of  $D$  with  $\Sigma$ . On the one hand,  $D'$  satisfies both  $Q'$  and  $\Sigma$  by construction. On the other hand, we will argue that  $D'$  does not satisfy  $Q$ . Assume, for the sake of contradiction, that  $D'$  satisfies  $Q$ . Then,  $D'$  contains  $\text{YAssigned}$  and  $\text{ZAssigned}$  atoms, which require  $\text{XAssignedAll}()$  and, in turn,  $\text{XAssignedUntil}_i()$  in  $D'$  (because none of these are in  $D$ , so should be generated by  $\Sigma$ ). But this is only possible, because  $\text{XAssignedUntil}_i()$  is not in  $D$  and can be obtained only by the corresponding rule in (6). But this rule cannot fire, because it requires  $Q_i$ , which is not present in  $D$  and is not obtainable by any rules in  $\Sigma$ . So, we conclude that  $D'$  does not satisfy  $Q$ .

Since  $D'$  does not satisfy  $Q$ , this means that  $Q$  is not equivalent to  $Q$  under  $\Sigma$ . So, our assumption was wrong and  $Q' \subseteq Q_i$  indeed holds for all  $i$ .  $\square$

**Claim 3.** Suppose  $Q'$  is a positive query equivalent to  $Q$  relative to  $\Sigma$ , and also that  $Q'$  has size at most  $k$ . Then  $Q'$  has the following form (up to renaming of variables):

$$\exists \bar{v}_\sigma (\sigma \wedge \text{XAssigned}(u_{x_1}) \wedge \dots \wedge \text{XAssigned}(u_{x_n})),$$

where each  $u_{x_i}$ ,  $i = 1, \dots, n$ , is either  $t_{x_i}$  or  $f_{x_i}$ .

*Proof.* Consider again the UCQ that is obtained by flattening  $Q'$ , and any CQ  $Q''$  in this UCQ. Consider also an output  $\bar{Q}''$  from disjunctive chasing of  $Q''$  with  $\Sigma$ . Since  $Q$  is equivalent to  $Q'$  under  $\Sigma$ , there exists a homomorphism  $h$  from  $Q$  to  $\bar{Q}''$ . Since  $Q$  contains  $\sigma$ , we have  $h(\sigma) \subseteq \bar{Q}''$ . However,  $\Sigma$  does not mention any predicates of  $\sigma$  in the heads, so  $h(\sigma) \subseteq Q''$ . Since  $Q$  is equivalent to  $Q'$  under  $\Sigma$ , we also know that there exists a homomorphism  $h'$  from  $Q'$  to any instance  $\bar{Q}$  in the disjunctive chase of  $Q$ . In particular,  $h'(h(\sigma)) \subseteq \bar{Q}$ . Again, since  $\Sigma$  does not create atoms over predicates of  $\sigma$ , we have that  $h'(h(\sigma)) \subseteq Q$ . The conjunction  $\sigma$  has a separate unary predicate for each of its variables, which are not mentioned in  $\tau$  and  $\pi$ . That is, they are not mentioned in the other parts of  $Q$ . So, the composition of  $h'$  and  $h$  is the identity on the variables of  $\sigma$ , and, in particular,  $h$  is injective on these variables.

Overall, we have that each CQ  $Q''$  in the UCQ representing  $Q'$  has a subformula (possibly re-arranging the order of existential quantifiers) that is isomorphic to  $\sigma$ . By Claim 2, each such CQ has a homomorphism from at least one of

$$\sigma \wedge \text{XAssigned}(t_{x_i}) \quad \text{and} \quad \sigma \wedge \text{XAssigned}(f_{x_i})$$

for each  $i = 1, \dots, n$ .

Let  $Q'_{\min}$  be a minimal (in the number of atoms) positive existential first order sentence equivalent to  $Q'$ . Since  $\sigma$  has a separate unary predicate for each of its variables, it does not have a non-trivial automorphism. Therefore, the body of  $Q'_{\min}$

has, up to renaming of variables, the form  $\sigma \wedge \eta$ , such that, for each  $i$ ,  $\eta$  implies at least one of  $\text{XAssigned}(t_{x_i})$ ,  $\text{XAssigned}(f_{x_i})$ ,  $\text{VTrue}_{x_i}(t'_{x_i}) \wedge \text{XAssigned}(t'_{x_i})$ , and  $\text{VFalse}_{x_i}(f'_{x_i}) \wedge \text{XAssigned}(f'_{x_i})$ , where  $t_{x_i}$  and  $f_{x_i}$  are the variables of  $\sigma$ , while  $t'_{x_i}$  and  $f'_{x_i}$  are some other variables.

However, we know that the number of atoms in  $Q'$  and, therefore,  $Q'_{\min}$  is bounded by the number of atoms in  $\sigma$  plus  $n$ , the size of  $\bar{x}$ . So, the number of atoms in  $\eta$  is bounded by  $n$ , which is only possible if  $\eta = \text{XAssigned}(u_{x_1}) \wedge \dots \wedge \text{XAssigned}(u_{x_n})$  where  $u_{x_i}$  is either  $t_{x_i}$  or  $f_{x_i}$ . Therefore,  $Q'$  is equal, up to renaming of variables, to  $Q'_{\min}$ , and has the form as required in the claim.  $\square$

By Claim 3, there exists a uniquely defined assignment of variables  $\bar{x}$  corresponding to  $Q'$ . Next we prove that for any extension of this assignment to variables  $\bar{y}$  there exists an extension to  $\bar{z}$  such that  $\psi$  holds under the overall assignment. Suppose, for the sake of contradiction, that this is not the case. Thus there is an extension to  $\bar{y}$  such that  $\psi$  evaluates to *false* for all extensions to  $\bar{z}$ . Let  $D$  be the database obtained from  $Q'$  by instantiating each variable  $v$  by a fresh constant  $a_v$ . Let  $D'$  be the extension of  $D$  with the atoms

$$\begin{array}{ll} \text{XAssignedUntil}_i(), & \text{for all } i = 1, \dots, n, \\ \text{XAssignedAll}(), & \\ \text{YAssigned}(a_{t_y}), & \text{for all } y \in \bar{y} \text{ assigned to } \textit{true} \text{ by the extension,} \\ \text{YAssigned}(a_{f_y}), & \text{for all } y \in \bar{y} \text{ assigned to } \textit{false} \text{ by the extension,} \\ \text{ZAssigned}(a_{t_z}), \text{ZAssigned}(a_{f_z}), & \text{for each } z \in \bar{z}. \end{array}$$

By construction,  $D'$  satisfies  $Q'$  and all the rules in  $\Sigma$ , except, possibly, (5). Next we prove that  $D'$  does not satisfy  $Q$ . Note that if  $D'$  does not satisfy  $Q$  then (a)  $D'$  satisfies all  $\Sigma$ , including (5) (b) since  $D'$  satisfies  $Q'$ , queries  $Q$  and  $Q'$  are not equivalent, which is a contradiction.

Suppose, for the sake of contradiction, that  $D'$  satisfies  $Q$ . That is, there is a homomorphism  $h$  from the body of  $Q$  to  $D'$ . By construction,

$$\begin{array}{ll} h(w_\gamma) = a_{v_\gamma^j}, & \text{for each clause } \gamma \text{ in } \psi \text{ and some } j = 1, \dots, 7, \\ h(w_s) = a_{t_s}, & \text{for each } s \in \bar{x} \cup \bar{y} \text{ assigned to } \textit{true}, \\ h(w_s) = a_{f_s}, & \text{for each } s \in \bar{x} \cup \bar{y} \text{ assigned to } \textit{false}, \end{array}$$

and, for each  $z \in \bar{z}$ ,

$$\text{either } h(w_z) = a_{t_z} \quad \text{or} \quad h(w_z) = a_{f_z}.$$

Consider the extension of the assignment of  $\bar{x} \cup \bar{y}$  to  $\bar{z}$  according to  $h$ . By construction, it satisfies  $\psi$ , which contradicts our assumption that  $\psi$  evaluates to *false* for all extensions to  $\bar{z}$ . Therefore, our second assumption was wrong and  $D'$  does not satisfy  $Q$  (and, therefore, satisfies  $\Sigma$ ). As we already said, this contradicts the fact that  $Q$  and  $Q'$  are equivalent under  $\Sigma$ . Therefore, our first assumption was also wrong, and indeed for any extension of the assignment of  $\bar{x}$  defined by  $Q'$  to  $\bar{y}$  there is an extension to  $\bar{z}$  such that  $\psi$  holds. This means that  $\varphi$  is valid, as required.  $\square$

## Correctness of Optimized Huang

We prove Proposition 3:

The Optimized Huang Algorithm produces an interpolant  $\varphi$  for  $\Sigma_1 \models \Sigma_2$ , and any negated literal in  $\varphi$  occurs negated in both  $\Sigma_1$  and  $\Sigma_2$ .

*Proof.* One can show by induction that if a literal occurs in a proof node  $n$  with provenance LEFT, then it is the result of unifying a literal in  $\Sigma_1$ , while if it occurs with provenance RIGHT it is the result of unifying a literal in  $\neg\Sigma_2$ . From this we can easily see that every relational symbol occurring in  $\varphi$  is present in both  $\Sigma_1$  and  $\Sigma_2$ , and that if it occurs negated it must be negated in both  $\Sigma_1$  and  $\Sigma_2$ .

For a clause  $\varphi$  in which each negated or non-negated literal is annotated with LEFT or RIGHT, we let  $\varphi_{\text{LEFT}}$  be the disjunction of all LEFT-labelled literals in  $\varphi$ .  $\varphi_{\text{LEFT}}$  is defined to be False if there are no such literals. Let  $\varphi_{\text{RIGHT}}$  be defined similarly.

We will prove the following invariant that is derived from Bonacina and Johansson [Bonacina and Johansson, 2011]:

For a proof node with formula  $\alpha$  and interpolant  $\varphi$ , we have

- $\Sigma_1 \models \varphi \vee \alpha_{\text{LEFT}}$
- $\Sigma_2 \wedge \varphi \models \alpha_{\text{RIGHT}}$

In particular, if  $\alpha = \perp$  then  $\Sigma_1 \models \varphi$  and  $\varphi \models \neg\Sigma_2$  as required.

The base cases are clear for both items. We prove the inductive cases, fixing a non-root node  $c$  with parents  $p_1 = r_1 \vee \alpha(\rho_1 \dots \rho_n)$  and  $p_2 = r_2 \vee \neg\alpha(\tau_1 \dots \tau_n)$  having interpolants  $\varphi_1$  and  $\varphi_2$  respectively. Thus  $c = \sigma(r_1) \vee \sigma(r_2)$  for some unifier  $\sigma$ .



- We consider the case where provenance is LEFT for  $\alpha(\rho_1 \dots \rho_n)$  and  $\neg\alpha(\tau_1 \dots \tau_n)$ . The interpolant is then  $\sigma(\varphi_1) \vee \sigma(\varphi_2)$ . For the first item, we must show that  $\Sigma_1 \models \sigma(\varphi_1) \vee \sigma(\varphi_2) \vee (\sigma(r_1) \vee \sigma(r_2))_{\text{LEFT}}$ . So assume we have an instance  $I$  that satisfies  $\Sigma_1$ ,  $\neg\sigma(\varphi_1)$ , and  $\neg\sigma(\varphi_2)$ . We know  $I \models \sigma(r_1 \vee \alpha(\rho_1 \dots \rho_n))_{\text{LEFT}} = \sigma(r_1)_{\text{LEFT}} \vee \sigma(\alpha(\rho_1 \dots \rho_n))$  by induction (the equality being since  $\alpha$  is LEFT). Similarly we know  $I \models \sigma(r_2)_{\text{LEFT}} \vee \neg\sigma(\alpha(\tau_1 \dots \tau_n))$  by induction. Since one of  $\sigma(\alpha(\rho_1 \dots \rho_n))$ ,  $\neg\sigma(\alpha(\tau_1 \dots \tau_n))$  holds we get  $I \models \sigma(r_1 \vee r_2)_{\text{LEFT}}$ , as required. For the second item, we must show  $\Sigma_2 \wedge (\varphi_1 \vee \varphi_2) \models \sigma(r_1 \vee r_2)_{\text{RIGHT}}$ . So assume  $I \models \Sigma_2$  and, without loss of generality  $\varphi_1$  (the case where  $I \models \varphi_2$  is symmetric). By induction  $I \models \sigma(r_1 \vee \alpha(\rho_1 \dots \rho_n))_{\text{RIGHT}} = \sigma(r_1)_{\text{RIGHT}}$ . Thus  $I \models \sigma(r_1 \vee r_2)_{\text{RIGHT}}$  as required.
- Now consider the case where the provenance is RIGHT for both  $\alpha(\rho_1 \dots \rho_n)$  and  $\neg\alpha(\tau_1 \dots \tau_n)$ . The interpolant is  $\sigma(\varphi_1) \wedge \sigma(\varphi_2)$ . For the first item, we must show that  $\Sigma_1 \models \sigma(\varphi_1 \wedge \varphi_2) \vee \sigma(r_1 \vee r_2)_{\text{LEFT}}$ . If  $I$  does not satisfy  $\sigma(\varphi_1 \wedge \varphi_2)$ , it must fail one of  $\sigma(\varphi_1)$ ,  $\sigma(\varphi_2)$ . Without loss of generality, assume  $I$  does not satisfy  $\sigma(\varphi_1)$ . By induction  $\Sigma_1 \models \varphi_1 \vee (r_1 \vee \alpha(\rho_1 \dots \rho_n))_{\text{LEFT}}$ , and so  $I \models (r_1)_{\text{LEFT}}$ . Thus  $I \models \sigma(r_1 \vee r_2)_{\text{LEFT}}$ , as required. For the second item, we must show that  $\Sigma_2 \wedge \sigma(\varphi_1 \wedge \varphi_2) \models \sigma(r_1 \vee r_2)_{\text{RIGHT}}$ , and thus assume  $I \models \Sigma_2 \wedge \sigma(\varphi_1 \wedge \varphi_2)$ . By induction  $I \models \sigma(r_1 \vee \alpha(\rho_1 \dots \rho_n))_{\text{RIGHT}} = \sigma(r_1)_{\text{RIGHT}} \vee \sigma(\alpha(\rho_1 \dots \rho_n))$ . Also by induction  $I \models \sigma(r_2 \vee \neg\alpha(\tau_1 \dots \tau_n))_{\text{RIGHT}} = \sigma(r_2)_{\text{RIGHT}} \vee \neg\sigma(\alpha(\tau_1 \dots \tau_n))$ . Putting these together we can conclude that  $I \models \sigma(r_1 \vee r_2)_{\text{RIGHT}}$ , hence  $I \models \sigma(r_1 \vee r_2)_{\text{RIGHT}}$  as required.
- We turn to the case where the provenance is LEFT for  $\alpha(\rho_1 \dots \rho_n)$  and RIGHT for  $\neg\alpha(\tau_1 \dots \tau_n)$ . Then the interpolant is  $\sigma(\varphi_1 \vee (\alpha(\rho_1 \dots \rho_n) \wedge \varphi_2))$ . For the first item we must show that  $\Sigma_1 \models \sigma(\varphi_1) \vee \sigma(\alpha(\rho_1 \dots \rho_n) \wedge \varphi_2) \vee \sigma(r_1 \vee r_2)_{\text{LEFT}}$ . Suppose  $I$  satisfies  $\Sigma_1 \wedge \neg\sigma(\varphi_1) \wedge \neg\sigma(\alpha(\rho_1 \dots \rho_n) \wedge \varphi_2)$ . Then using the induction hypothesis,  $I \models \sigma(r_1 \vee \alpha(\rho_1 \dots \rho_n))_{\text{LEFT}} = \sigma(r_1)_{\text{LEFT}} \vee \sigma(\alpha(\rho_1 \dots \rho_n))$ .  $I$  must satisfy either  $\neg\sigma(\alpha(\rho_1 \dots \rho_n))$  or  $\neg\sigma(\varphi_2)$ , by assumption. If  $I \models \sigma(\varphi_2)$  then by induction  $I \models \sigma(r_2 \vee \neg\alpha(\tau_1 \dots \tau_n))_{\text{LEFT}} = \sigma(r_2)_{\text{LEFT}}$ , and thus  $I \models \sigma(r_1 \vee r_2)_{\text{LEFT}}$  and we are done. On the contrary, if  $I \models \neg\sigma(\alpha(\rho_1 \dots \rho_n))$ , then since  $I \models \sigma(r_1)_{\text{LEFT}} \vee \sigma(\alpha(\rho_1 \dots \rho_n))$  we deduce that  $I \models \sigma(r_1)_{\text{LEFT}}$ ; again  $I \models \sigma(r_1 \vee r_2)_{\text{LEFT}}$  and we are done. For the second item we need to show that  $\Sigma_2 \wedge (\sigma(\varphi_1) \vee \sigma(\alpha(\rho_1 \dots \rho_n) \wedge \varphi_2)) \models \sigma(r_1 \vee r_2)_{\text{RIGHT}}$ . Thus assume  $I \models \Sigma_2 \wedge \sigma(\varphi_1) \vee \sigma(\alpha(\rho_1 \dots \rho_n) \wedge \varphi_2)$ . In the first case, suppose  $I \models \sigma(\varphi_1)$ . Then by induction  $I \models \sigma(r_1 \vee \alpha(\rho_1 \dots \rho_n))_{\text{RIGHT}} = \sigma(r_1)_{\text{RIGHT}}$  and we are done. In the second case,  $I \models \sigma(\alpha(\rho_1 \dots \rho_n) \wedge \varphi_2)$ . Then by induction  $I \models \sigma(r_2 \wedge \neg\alpha(\tau_1 \dots \tau_n))_{\text{RIGHT}}$ . Since  $I \models \sigma(\alpha(\rho_1 \dots \rho_n))$  this implies that  $I \models \sigma(r_2)_{\text{RIGHT}}$ , and from there we easily conclude.
- Finally, we consider the case where the provenance is RIGHT for  $\alpha(\rho_1 \dots \rho_n)$  and LEFT for  $\neg\alpha(\tau_1 \dots \tau_n)$ . The interpolant is  $\sigma(\varphi_2) \vee (\neg\sigma(\alpha(\tau_1 \dots \tau_n)) \wedge \sigma(\varphi_1))$ . For the first item we must show that  $\Sigma_1 \models \sigma(\varphi_2) \vee (\neg\sigma(\alpha(\tau_1 \dots \tau_n)) \wedge \sigma(\varphi_1)) \vee \sigma(r_1 \vee r_2)_{\text{LEFT}}$ . Towards this goal we assume  $I \models \Sigma_1 \wedge \neg\sigma(\varphi_2) \wedge \neg(\neg\sigma(\alpha(\tau_1 \dots \tau_n)) \wedge \sigma(\varphi_1))$ . By induction  $I \models \sigma(r_2 \vee \neg\alpha(\tau_1 \dots \tau_n))_{\text{LEFT}} = \sigma(r_2)_{\text{LEFT}} \vee \neg\sigma(\alpha(\tau_1 \dots \tau_n))$ . If  $I \models \sigma(r_2)_{\text{LEFT}}$  we have that  $I \models \sigma(r_1 \vee r_2)_{\text{LEFT}}$  and we are done. So we can assume  $I \models \neg\sigma(\alpha(\tau_1 \dots \tau_n))$ . Since  $I \models \neg(\neg\sigma(\alpha(\tau_1 \dots \tau_n)) \wedge \sigma(\varphi_1))$  we deduce that  $I \models \neg\sigma(\varphi_1)$ . Thus by induction  $I \models \sigma(r_2 \vee \alpha(\tau_1 \dots \tau_n))_{\text{LEFT}} = \sigma(r_2)_{\text{LEFT}}$  and we are done. For the second item we must show that  $\Sigma_2 \wedge \sigma(\varphi_2 \vee (\neg\alpha(\tau_1 \dots \tau_n) \wedge \varphi_1)) \models \sigma(r_1 \vee r_2)_{\text{RIGHT}}$ . Again we assume  $I$  satisfies  $\Sigma_2 \wedge \sigma(\varphi_2 \vee (\neg\alpha(\tau_1 \dots \tau_n) \wedge \varphi_1))$ . If  $I \models \sigma(\varphi_2)$ , then inductively  $I \models \sigma(r_2 \vee \neg\alpha(\tau_1 \dots \tau_n))_{\text{RIGHT}} = \sigma(r_2)_{\text{RIGHT}}$  and we are done. If  $I \models \neg\sigma(\alpha(\tau_1 \dots \tau_n)) \wedge \sigma(\varphi_1)$ , then by induction  $I \models \sigma(r_1 \vee \alpha(\tau_1 \dots \tau_n))_{\text{RIGHT}} = \sigma(r_1)_{\text{RIGHT}} \vee \sigma(\alpha(\tau_1 \dots \tau_n))$ . Using  $I \models \neg\sigma(\alpha(\tau_1 \dots \tau_n))$  we can conclude.

□

## Propositional Optimizations

We now report about all propositional optimizations that were employed during the inductive computation of the interpolant. To do this, let  $\alpha$  be the proposition involved in the resolution of the clause  $p_1$ , having interpolant  $\varphi_1$ , with the clause  $p_2$ , having interpolant  $\varphi_2$ . In particular, the occurrence of  $\alpha$  is positive in  $p_1$ , i.e.,  $\alpha \in p_1$ , and negative in  $p_2$ , i.e.,  $\neg\alpha \in p_2$ . The obtained clause is  $p = (p_1 \setminus \{\alpha\}) \cup (p_2 \setminus \{\neg\alpha\})$ .

- If the provenance of  $\alpha$  is LEFT in  $p_1$  and that of  $\neg\alpha$  is RIGHT in  $p_2$ , we have that the interpolant  $\varphi$  associated with  $p$  is computed as follows.
  - $\varphi = \varphi_1$ , if  $\varphi_1$  simplifies to True when all occurrences of  $\alpha$  in it are replaced by True.
  - $\varphi = \alpha \wedge (\varphi'_1 \vee \varphi'_2)$ , if the previous condition does not hold and  $\varphi_1$  simplifies to False when all occurrences of  $\alpha$  in it are replaced by False, where  $\varphi'_1$  and  $\varphi'_2$  are obtained as the simplifications of  $\varphi_1$  and  $\varphi_2$ , respectively, after the substitution of all their occurrences  $\alpha$  with True.

- $\varphi = \varphi_1 \vee (\alpha \wedge \varphi'_2)$ , if the previous two conditions do not hold, where  $\varphi'_2$  is the obtained as the simplification of  $\varphi_2$  after the substitution of all its occurrences  $\alpha$  with True.
- If the provenance of  $\alpha$  is RIGHT in  $p_1$  and that of  $\neg\alpha$  is LEFT in  $p_2$ , we have that the interpolant  $\varphi$  associated with  $p$  is computed as follows.
  - $\varphi = \varphi_2$ , if  $\varphi_2$  simplifies to True when all occurrences of  $\alpha$  in it are replaced by False.
  - $\varphi = \neg\alpha \wedge (\varphi'_1 \vee \varphi'_2)$ , if the previous condition does not hold and  $\varphi_2$  simplifies to False when all occurrences of  $\alpha$  in it are replaced by True, where  $\varphi'_1$  and  $\varphi'_2$  are obtained as the simplifications of  $\varphi_1$  and  $\varphi_2$ , respectively, after the substitution of all their occurrences  $\alpha$  with False.
  - $\varphi = \varphi_2 \vee (\neg\alpha \wedge \varphi'_1)$ , if the previous two condition do not hold, where  $\varphi'_2$  is the obtained as the simplification of  $\varphi_2$  after the substitution of all its occurrences  $\alpha$  with False.

## Details on Experiments

We are interested in generating test problems for reformulation consisting of an input query  $Q$ , a theory  $\Sigma$  in some particular format (e.g. FullDisTGD) and a set of target relations  $\mathcal{V}$ . To the best of our knowledge, no such testbed exists. The challenge is that we want to generate  $Q, \Sigma, \mathcal{V}$  so that: we know  $Q$  has a  $\Sigma$ -reformulation over  $\mathcal{V}$ , but it is non-trivial for a reasoning system to find such a reformulation. For example, if we randomly generated a theory  $\Sigma$ , then surely it would rarely be the case that a reformulation existed. Conversely there are many simple recipes for “re-arranging”  $Q$ : we could simply rename relations in  $Q$ , added the renamed relations to  $\mathcal{V}$ , and added equivalences between the original relations and the renamed ones in  $\Sigma$ , and we would guarantee that a reformulation exists. But the simple nature of our encoding might make this atypically easy for a reasoning system.

Our testbed is based on two components. First, we have a *target reformulation generator* TargetReformGen, that generates a random sentence  $Q'$  in some desired target language (e.g. CQ, UCQ) over some vocabulary  $\mathcal{V}$ .

Second, we have a *theory generator* TheoryGen that inputs a target  $Q'$  and generates  $Q, \Sigma$  such that  $Q$  is equivalent to  $Q'$  over  $\Sigma$ .

Our testbed is then formed by first randomly generating target reformulations using TargetReformGen and then feeding each into TheoryGen to get the final reformulation problem.

The generator TargetReformGen produces three types of sentences: CNF, DNF and sentences with a given alternation depth AltDepth. The generator works for both first-order logic and propositional logic. We focus on the parameters that are relevant for the propositional case below. In the propositional CNF case, the generator is parametrised by

1. the total number of conjuncts,
2. the number of literals in each conjunct,
3. the probability of reusing the same relation name in different conjuncts from  $Q'$ ,
4. the probability of having negative literals and, finally,
5. a seed for producing random numbers.

Given these parameters, TargetReformGen first produces a vocabulary  $\mathcal{V}$  and then a sentence  $Q'$  using relations from  $\mathcal{V}$ . TargetReformGen works similarly for the DNF case. In the alternating case, TargetReformGen produces  $\text{AltDepth} + 1$  CNF and DNF subsentences and “glues” them using conjunction or disjunction such that the alternation depth of the resulting sentence  $Q'$  is AltDepth. TargetReformGen is parametrised by (i) the alternation depth AltDepth, (ii) the number of literals in each subsentence, (iii) the type of the first subsentence of  $Q'$  (CNF or DNF) and, finally, by the last three arguments also used in the CNF and DNF cases.

Our generator TheoryGen will generate new relations in TheoryDepth layers, and then rules going in each direction between layer  $i$  and layer  $i + 1$ . In the most basic configuration, the rules will match the structure of the input reformulation  $\varphi$ , with disjunctive rules matching disjunction in  $\varphi$  and conjunctive rules matching conjunction in  $\varphi$ .

For example, if our target reformulation is  $\varphi = (T_1 \wedge T_2) \vee T_3$ , we might add:

$$\begin{aligned}
 Q &\rightarrow R_1 \vee R_2 \\
 R_1 &\rightarrow T_1 \wedge T_2 \\
 R_2 &\rightarrow T_3 \\
 T_3 &\rightarrow R'_2 \\
 T_1 \wedge T_2 &\rightarrow R'_1 \\
 R'_1 &\rightarrow Q \\
 R'_2 &\rightarrow Q
 \end{aligned}$$

Clearly, the rules ensure that  $\varphi$  is a reformulation of  $Q$ . Additional parameters allow the generated rules to be more complex:

- For each “layer” of alternation, we can arrange that its generation occurs not in one steps (e.g.  $R_2 \rightarrow T_2$ ) but as a sequence of  $n$  steps of TGDs, where  $n$  is a parameter.
- We can generate additional “noise” rules, with heads and bodies randomly chosen. Since these are added in addition to the “core rules” above, they will not impact the fact that  $\varphi$  is a reformulation of  $Q$ , but will make this harder to detect. The number of such rules is determined by an additional configuration parameter.

In the experiments regarding CNF reformulations, we had both the numbers of conjuncts and literals ranging from 2 to 6. The probability of reusing the same relation was set either to 0 or to  $3/4$ . Similarly, the probability of using negated literals was either 0 or  $5/8$ . The same parameters were used for the DNF and alternating reformulations. In particular, for the latter, the nesting of conjunctions and disjunctions varied from 2 to 5.