

# The Space-Efficient Core of Vadalog

Gerald Berger

Institute of Logic and Computation  
TU Wien

Andreas Pieris

School of Informatics  
University of Edinburgh

Georg Gottlob

Department of Computer Science  
University of Oxford & TU Wien

Emanuel Sallinger

Department of Computer Science  
University of Oxford

## ABSTRACT

Vadalog is a system for performing complex reasoning tasks such as those required in advanced knowledge graphs. The logical core of the underlying Vadalog language is the warded fragment of tuple-generating dependencies (TGDs). This formalism ensures tractable reasoning in data complexity, while a recent analysis focusing on a practical implementation led to the reasoning algorithm around which the Vadalog system is built. A fundamental question that has emerged in the context of Vadalog is the following: can we limit the recursion allowed by wardedness in order to obtain a formalism that provides a convenient syntax for expressing useful recursive statements, and at the same time achieves space-efficiency? After analyzing several real-life examples of warded sets of TGDs provided by our industrial partners, as well as recent benchmarks, we observed that recursion is often used in a restricted way: the body of a TGD contains at most one atom whose predicate is mutually recursive with a predicate in the head. We show that this type of recursion, known as piece-wise linear in the Datalog literature, is the answer to our main question. We further show that piece-wise linear recursion alone, without the wardedness condition, is not enough as it leads to the undecidability of reasoning. We finally study the relative expressiveness of the query languages based on (piece-wise linear) warded sets of TGDs.

### ACM Reference Format:

Gerald Berger, Georg Gottlob, Andreas Pieris, and Emanuel Sallinger. 2018. The Space-Efficient Core of Vadalog. In *Proceedings of ACM Conference (Conference'17)*. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

## 1 INTRODUCTION

In recent times, thousands of companies world-wide wish to manage their own knowledge graphs (KGs), and are looking for adequate knowledge graph management systems (KGMS). The term knowledge graph originally only referred to Google's Knowledge Graph, i.e., "a knowledge base used by Google and its services to enhance its search engine's results with information gathered from

a variety of sources.<sup>1</sup>" In the meantime, several other large companies have constructed their own knowledge graphs, and many more companies would like to maintain a private corporate knowledge graph incorporating large amounts of data in form of database facts, both from corporate and public sources, as well as rule-based knowledge. Such a corporate knowledge graph is expected to contain relevant business knowledge, for example, knowledge about customers, products, prices, and competitors, rather than general knowledge from Wikipedia and similar sources. It should be managed by a KGMS, that is, a knowledge base management system, which performs complex rule-based reasoning tasks over very large amounts of data and, in addition, provides methods and tools for data analytics and machine learning [5].

### 1.1 The Vadalog System

Vadalog is a system for performing complex reasoning tasks such as those required in advanced knowledge graphs [6]. It is Oxford's contribution to the VADA research project,<sup>2</sup> a joint effort of the universities of Oxford, Manchester, and Edinburgh, as well as around 20 industrial partners such as Facebook, BP, and the NHS (UK national health system). One of the most fundamental reasoning tasks performed by Vadalog is *ontological query answering*: given a database  $D$ , an ontology  $\Sigma$  (which is essentially a set of logical assertions that allow us to derive new intensional knowledge from  $D$ ), and a query  $q(\bar{x})$  (typically a conjunctive query), the goal is to compute the certain answers to  $q$  w.r.t. the knowledge base consisting of  $D$  and  $\Sigma$ , i.e., the tuples of constants  $\bar{c}$  such that, for every relational instance  $I \supseteq D$  that satisfies  $\Sigma$ ,  $I$  satisfies the Boolean query  $q(\bar{c})$  obtained after instantiating  $\bar{x}$  with  $\bar{c}$ . Due to Vadalog's ability to perform ontological query answering, it is currently used as the core deductive database component of the overall Vadalog KGMS, as well as at various industrial partners including the finance, security, and media intelligence industries.

The logical core of the underlying Vadalog language is a rule-based formalism known as *warded Datalog*<sup>3</sup> [16], which is a member of the Datalog<sup>±</sup> family of knowledge representation languages [10]. Warded Datalog<sup>3</sup> generalizes Datalog with existential quantification in rule heads, and at the same time applies a restriction on how certain "dangerous" variables can be used; details are given in Section 3. Such a restriction is needed as basic reasoning tasks, e.g., ontological query answering, under arbitrary Datalog<sup>3</sup> rules become undecidable; see, e.g., [4, 9]. Let us clarify that Datalog<sup>3</sup> rules are essentially *tuple-generating dependencies* (TGDs) of the form

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).  
Conference'17, July 2017, Washington, DC, USA

© 2018 Association for Computing Machinery.  
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00  
<https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

<sup>1</sup>[https://en.wikipedia.org/wiki/Knowledge\\_Graph](https://en.wikipedia.org/wiki/Knowledge_Graph)

<sup>2</sup><http://vada.org.uk/>

$\forall \bar{x} \forall \bar{y} (\varphi(\bar{x}, \bar{y}) \rightarrow \exists \bar{z} \psi(\bar{x}, \bar{z}))$ , where  $\varphi$  (the *body*) and  $\psi$  (the *head*) are conjunctions of relational atoms. Therefore, knowledge representation and reasoning should be seen as a modern application of TGDs, which have been introduced decades ago as a unifying framework for database integrity constraints.

The key properties of warded Datalog<sup>3</sup>, which led to its adoption as the logical core on top of which the Vadalog language is built, can be summarized as follows:

- (1) *Recursion over KGs*. It is able to express full recursion and joins, which are needed to express complex reasoning tasks over KGs. Moreover, navigational capabilities, empowered by recursion, are vital for graph-based structures.
- (2) *Ontological Reasoning over KGs*. After adding a very mild and easy to handle negation, the language is able to express SPARQL reasoning under the OWL 2 QL entailment regime. Recall that SPARQL is the standard language for querying the Semantic Web,<sup>3</sup> while OWL 2 QL is a prominent profile of the OWL 2 Web Ontology Language, the standard formalism for modeling Semantic Web ontologies.<sup>4</sup>
- (3) *Low Complexity*. Reasoning, in particular, ontological query answering, is tractable (in fact, polynomial time) in data complexity, which is a minimal requirement for allowing scalability over large volumes of data.

Warded Datalog<sup>3</sup> turned out to be powerful enough for expressing all the tasks given by our industrial partners, while a recent analysis of it focusing on a practical implementation led to the reasoning algorithm around which the Vadalog system is built [6].

## 1.2 Research Challenges

With the aim of isolating more refined formalisms, which will lead to yet more efficient reasoning algorithms, the following fundamental question has emerged in the context of Vadalog:

*Can we limit the recursion allowed by wardedness in order to obtain a formalism that provides a convenient syntax for expressing useful statements, importantly, most of the scenarios provided by our industrial partners, and at the same time achieves space-efficiency, in particular, NLOGSPACE data complexity?*

Let us stress that NLOGSPACE data complexity is the best that we can hope for, since navigational capabilities are vital for graph-based structures, and already graph reachability is NLOGSPACE-hard. It is known that NLOGSPACE is contained in the class NC<sub>2</sub> of highly parallelizable problems. This means that reasoning in the more refined formalism that we are aiming is principally parallelizable, unlike warded Datalog<sup>3</sup>, which is PTIME-complete and intrinsically sequential. Our ultimate goal is to exploit this in the future for the parallel execution of reasoning tasks in both multi-core settings and in the map-reduce model. In fact, we are currently in the process of implementing a multi-core implementation for the refined formalism proposed by the present work.

Extensive benchmark results are available for the Vadalog system, based on a variety of scenarios, both synthetic and industrial scenarios, including: ChaseBench [7], a benchmark that targets data exchange and query answering problems; iBench, a data exchange

benchmark developed at the University of Toronto [3]; iWarded, a benchmark specifically targeted at warded sets of TGDs; a DBpedia based benchmark; and a number of other synthetic and industrial scenarios [6]. Let us stress that all the above benchmarks contain only warded sets of TGDs. In fact, a good part of them are not *warded by chance*, i.e., they contain joins among “harmful” variables, which is one of the distinctive features of wardedness [6]. After analyzing the above benchmarks, we observed that recursion is often used in a restricted way. Approximately 70% of the TGD-sets use recursion in the following way: the body of a TGD contains at most one atom whose predicate is mutually recursive with a predicate in the head. More specifically, approximately 55% of the TGD-sets directly use the above type of recursion, while 15% can be transformed into warded sets of TGDs that use recursion as explained above. This transformation relies on a standard elimination procedure of unnecessary non-linear recursion. For example,

$$\forall x \forall y (E(x, y) \rightarrow T(x, y)) \quad \forall x \forall y \forall z (T(x, y) \wedge T(y, z) \rightarrow T(x, z)),$$

which compute the transitive closure of the extensional binary relation  $E$  using non-linear recursion, can be rewritten as the set

$$\forall x \forall y (E(x, y) \rightarrow T(x, y)) \quad \forall x \forall y \forall z (E(x, y) \wedge T(y, z) \rightarrow T(x, z))$$

that uses linear recursion. Interestingly, the type of recursion discussed above has been already studied in the context of Datalog, and is known as *piece-wise linear*; see, e.g., [1]. It is a refinement of the well-known *linear* recursion [23, 24], already mentioned in the above example, which allows only one intensional predicate to appear in the body, while all the other predicates are extensional.

Based on this key observation, the following research questions have immediately emerged:

- (1) Does warded Datalog<sup>3</sup> with piece-wise linear recursion achieve space-efficiency for ontological query answering?<sup>5</sup>
- (2) Is the combination of wardedness and piece-wise linearity justified? In other words, can we achieve the same with piece-wise linear Datalog<sup>3</sup> without the wardedness condition?
- (3) What is the expressiveness of the query language based on warded Datalog<sup>3</sup> with piece-wise linear recursion relative to prominent query languages such as Datalog?

These are top-priority questions in the context of the Vadalog system since they may provide useful insights towards more efficient reasoning algorithms, in particular, towards parallel execution of reasoning tasks. The ultimate goal of this work is to analyze piece-wise linearity, and provide definite answers to the above questions.

## 1.3 Summary of Contributions

Our main results can be summarized as follows:

- (1) Ontological query answering under warded Datalog<sup>3</sup> with piece-wise linear recursion is NLOGSPACE-complete in data complexity, and PSPACE-complete in combined complexity, which provides a definite answer to our first question. This is a rather involved result that heavily relies on a novel notion of resolution-based proof tree, which is of independent interest. In particular, we show that ontological query answering under warded Datalog<sup>3</sup> with piece-wise linear recursion boils

<sup>3</sup><http://www.w3.org/TR/rdf-sparql-query>

<sup>4</sup><https://www.w3.org/TR/owl2-overview/>

<sup>5</sup>The idea of combining wardedness with piece-wise linearity has been already mentioned in the invited paper [5], while the obtained formalism is called strongly warded.

down to the problem of checking whether a proof tree that enjoys certain properties exists, which in turn can be done via a space-bounded non-deterministic algorithm. Interestingly, our machinery allows us to re-establish the complexity of ontological query answering under warded Datalog<sup>∃</sup> via an algorithm that is significantly simpler than the one employed in [16]. This algorithm is essentially the non-deterministic algorithm for piece-wise linear warded Datalog<sup>∃</sup> with the crucial difference that it employs alternation.

- (2) To our surprise, ontological query answering under piece-wise linear Datalog<sup>∃</sup>, without the wardedness condition, is undecidable. This result, which is shown via a reduction from the standard unbounded tiling problem, provides a definite answer to our second question: the combination of wardedness and piece-wise linearity is indeed justified.
- (3) We finally investigate the relative expressive power of the query language based on warded Datalog<sup>∃</sup> with piece-wise linear recursion, which consists of all the queries of the form  $Q = (\Sigma, q)$ , where  $\Sigma$  is a warded set of TGDs with piece-wise linear recursion, and  $q$  is a conjunctive query, while the evaluation of  $Q$  over a database  $D$  is precisely the certain answers to  $q$  w.r.t.  $D$  and  $\Sigma$ . By exploiting our novel notion of proof tree, we show that it is equally expressive to piece-wise linear Datalog. The same approach allows us to elucidate the relative expressiveness of the query language based on warded Datalog<sup>∃</sup> (with arbitrary recursion), showing that it is equally expressive to Datalog. We also adopt the more refined notion of program expressive power, introduced in [2], which aims at the decoupling of the set of TGDs and the actual conjunctive query, and show that the query language based on warded Datalog<sup>∃</sup> (with piece-wise linear recursion) is strictly more expressive than Datalog (with piece-wise linear recursion). This result exposes the advantage of value invention that is available in Datalog<sup>∃</sup>-based languages.

Selected proofs are deferred to a clearly marked appendix.

## 2 PRELIMINARIES

**Basics.** We consider the disjoint countably infinite sets  $\mathbf{C}$ ,  $\mathbf{N}$ , and  $\mathbf{V}$  of *constants*, (*labeled*) *nulls*, and *variables*, respectively. The elements of  $(\mathbf{C} \cup \mathbf{N} \cup \mathbf{V})$  are called *terms*. An *atom* is an expression of the form  $R(\bar{t})$ , where  $R$  is an  $n$ -ary predicate, and  $\bar{t}$  is an  $n$ -tuple of terms. We write  $\text{var}(\alpha)$  for the set of variables in an atom  $\alpha$ ; this notation extends to sets of atoms. A *fact* is an atom that contains only constants. A *substitution* from a set of terms  $T$  to a set of terms  $T'$  is a function  $h: T \rightarrow T'$ . The restriction of  $h$  to a subset  $S$  of  $T$ , denoted  $h|_S$ , is the substitution  $\{t \mapsto h(t) \mid t \in S\}$ . A *homomorphism* from a set of atoms  $A$  to a set of atoms  $B$  is a substitution  $h$  from the set of terms in  $A$  to the set of terms in  $B$  such that  $h$  is the identity on  $\mathbf{C}$ , and  $R(t_1, \dots, t_n) \in A$  implies  $h(R(t_1, \dots, t_n)) = R(h(t_1), \dots, h(t_n)) \in B$ . We write  $h(A)$  for the set of atoms  $\{h(\alpha) \mid \alpha \in A\}$ . For brevity, we may write  $[n]$  for the set  $\{1, \dots, n\}$ , where  $n \geq 0$ .

**Relational Databases.** A *schema*  $\mathbf{S}$  is a finite set of relation symbols (or predicates), each having an associated *arity*. We write  $R/n$  to denote that  $R$  has arity  $n \geq 0$ . A position  $R[i]$  in  $\mathbf{S}$ , where  $R/n \in \mathbf{S}$  and  $i \in [n]$ , identifies the  $i$ -th argument of  $R$ . An *instance* over  $\mathbf{S}$  is a (possibly infinite) set of atoms over  $\mathbf{S}$  that contain constants and nulls, while a *database* over  $\mathbf{S}$  is a finite set of facts over  $\mathbf{S}$ . The

*active domain* of an instance  $I$ , denoted  $\text{dom}(I)$ , is the set of all terms occurring in  $I$ .

**Conjunctive Queries.** A *conjunctive query* (CQ) over  $\mathbf{S}$  is a first-order formula of the form

$$q(\bar{x}) := \exists \bar{y} (R_1(\bar{z}_1) \wedge \dots \wedge R_n(\bar{z}_n)),$$

where each  $R_i(\bar{z}_i)$ , for  $i \in [n]$ , is an atom without nulls over  $\mathbf{S}$ , each variable mentioned in the  $\bar{z}_i$ 's appears either in  $\bar{x}$  or  $\bar{y}$ , and  $\bar{x}$  are the *output variables* of  $q$ . For convenience, we adopt the rule-based syntax of CQs, i.e., a CQ as the one above will be written as the rule

$$Q(\bar{x}) \leftarrow R_1(\bar{z}_1), \dots, R_n(\bar{z}_n),$$

where  $Q$  is a predicate used only in the head of CQs. We write  $\text{atoms}(q)$  for the set of atoms  $\{R_1(\bar{z}_1), \dots, R_n(\bar{z}_n)\}$ . The *evaluation* of  $q(\bar{x})$  over an instance  $I$ , denoted  $q(I)$ , is the set of all tuples  $h(\bar{x})$  of constants with  $h$  being a homomorphism from  $\text{atoms}(q)$  to  $I$ .

**Tuple-Generating Dependencies.** A *tuple-generating dependency* (TGD)  $\sigma$  is a first-order sentence

$$\forall \bar{x} \forall \bar{y} (\phi(\bar{x}, \bar{y}) \rightarrow \exists \bar{z} \psi(\bar{x}, \bar{z})),$$

where  $\bar{x}, \bar{y}, \bar{z}$  are tuples of variables of  $\mathbf{V}$ , and  $\phi, \psi$  are conjunctions of atoms without constants and nulls. For brevity, we write  $\sigma$  as  $\phi(\bar{x}, \bar{y}) \rightarrow \exists \bar{z} \psi(\bar{x}, \bar{z})$ , and use comma instead of  $\wedge$  for joining atoms. We refer to  $\phi$  and  $\psi$  as the *body* and *head* of  $\sigma$ , denoted  $\text{body}(\sigma)$  and  $\text{head}(\sigma)$ , respectively. The *frontier* of the TGD  $\sigma$ , denoted  $\text{front}(\sigma)$ , is the set of variables that appear both in the body and the head of  $\sigma$ . We also write  $\text{var}_{\exists}(\sigma)$  for the existentially quantified variables of  $\sigma$ . The schema of a set  $\Sigma$  of TGDs, denoted  $\text{sch}(\Sigma)$ , is the set of predicates occurring in  $\Sigma$ . An instance  $I$  satisfies a TGD  $\sigma$  as the one above, written  $I \models \sigma$ , if the following holds: whenever there exists a homomorphism  $h$  such that  $h(\phi(\bar{x}, \bar{y})) \subseteq I$ , then there exists  $h' \supseteq h|_{\bar{x}}$  such that  $h'(\psi(\bar{x}, \bar{z})) \subseteq I$ .<sup>6</sup> The instance  $I$  satisfies a set  $\Sigma$  of TGDs, written  $I \models \Sigma$ , if  $I \models \sigma$  for each  $\sigma \in \Sigma$ .

**Query Answering under TGDs.** The main reasoning task under TGD-based languages is *conjunctive query answering*. Given a database  $D$  and a set  $\Sigma$  of TGDs, a *model* of  $D$  and  $\Sigma$  is an instance  $I$  such that  $I \supseteq D$  and  $I \models \Sigma$ . Let  $\text{mods}(D, \Sigma)$  be the set of all models of  $D$  and  $\Sigma$ . The *certain answers* to a CQ  $q$  w.r.t.  $D$  and  $\Sigma$  is

$$\text{cert}(q, D, \Sigma) := \bigcap \{q(I) \mid I \in \text{mods}(D, \Sigma)\}.$$

Our main task is to compute the certain answers to a CQ w.r.t. a database and a set of TGDs from a certain class  $\mathbf{C}$  of TGDs; concrete classes of TGDs are discussed below. As is customary when studying the complexity of this problem, we focus on its decision version:

PROBLEM :	CQAns( $\mathbf{C}$ )
INPUT :	A database $D$ , a set $\Sigma \in \mathbf{C}$ of TGDs, a CQ $q(\bar{x})$ , and a tuple $\bar{c} \in \text{dom}(D)^{ \bar{x} }$ .
QUESTION :	Is it the case that $\bar{c} \in \text{cert}(q, D, \Sigma)$ ?

We consider the standard complexity measures: *combined complexity* and *data complexity*, where the latter measures the complexity of the problem assuming that the set of TGDs and the CQ are fixed.

A useful algorithmic tool for tackling the above problem is the well-known *chase procedure*; see, e.g., [9, 14, 18, 22]. We start by

<sup>6</sup>By abuse of notation, we sometimes treat a tuple of variables as a set of variables, and a conjunction of atoms as a set of atoms.

defining a single chase step. Let  $I$  be an instance and  $\sigma = \phi(\bar{x}, \bar{y}) \rightarrow \exists \bar{z} \psi(\bar{x}, \bar{z})$  a TGD. We say that  $\sigma$  is *applicable* w.r.t.  $I$  if there exists a homomorphism  $h$  such that  $h(\phi(\bar{x}, \bar{y})) \subseteq I$ . In this case, *the result of applying  $\sigma$  over  $I$  with  $h$*  is the instance  $J = I \cup \{h'(\psi(\bar{x}, \bar{z}))\}$ , where  $h'(z)$  is a fresh null not occurring in  $I$ , for every  $z \in \bar{z}$ . Such a single chase step is denoted  $I \langle \sigma, h \rangle J$ . Consider now an instance  $I$ , and a set  $\Sigma$  of TGDs. A *chase sequence for  $I$  under  $\Sigma$*  is a sequence  $(I_i \langle \sigma_i, h_i \rangle I_{i+1})_{i \geq 0}$  of chase steps such that: (1)  $I = I_0$ ; (2) for each  $i \geq 0$ ,  $\sigma_i \in \Sigma$ ; and (3)  $\bigcup_{i \geq 0} I_i \models \Sigma$ . We call  $\bigcup_{i \geq 0} I_i$  the *result* of this chase sequence, which always exists. Although the result of a chase sequence is not necessarily unique (up to isomorphism), each such result is equally useful for query answering purposes, since it can be homomorphically embedded into every other result. Hence, we denote by  $\text{chase}(I, \Sigma)$  the result of an arbitrary chase sequence for  $I$  under  $\Sigma$ . The following is a classical result:

**PROPOSITION 2.1.** *Given a database  $D$ , a set  $\Sigma$  of TGDs, and a CQ  $q$ ,  $\text{cert}(q, D, \Sigma) = q(\text{chase}(D, \Sigma))$ .*

### 3 THE LOGICAL CORE OF VADALOG

A crucial component of the Vadalogue system is its reasoning engine, which in turn is built around the Vadalogue language, a general-purpose formalism for knowledge representation and reasoning. The logical core of this language is the well-behaved class of *warded* sets of TGDs that has been proposed in [16].

**An Intuitive Description.** Wardedness applies a syntactic restriction on how certain “dangerous” variables of a set of TGDs are used. These are body variables that can be unified with a null during the chase, and that are also propagated to the head. For example, given

$$P(x) \rightarrow \exists z R(x, z) \quad \text{and} \quad R(x, y) \rightarrow P(y)$$

the variable  $y$  in the body of the second TGD is dangerous. Indeed, once the chase applies the first TGD, an atom of the form  $R(\_, \perp)$  is generated, where  $\perp$  is a null value, and then the second TGD is triggered with the variable  $y$  being unified with  $\perp$  that is propagated to the obtained atom  $P(\perp)$ . It has been observed that the liberal use of dangerous variables leads to a prohibitively high computational complexity of the main reasoning tasks, in particular of CQ answering [9]. The main goal of *wardedness* is to limit the use of dangerous variables with the aim of taming the way that null values are propagated during the execution of the chase procedure. This is achieved by posing the following conditions:

- (1) all the dangerous variables should appear together in a single body atom  $\alpha$ , called a *ward*, and
- (2)  $\alpha$  can share only harmless variables with the rest of the body, i.e., variables that unify only with constants.

We proceed to formalize the above description.

**The Formal Definition.** We first need some auxiliary notions. The set of positions of a schema  $S$ , denoted  $\text{pos}(S)$ , is defined as  $\{R[i] \mid R/n \in S, \text{ with } n \geq 1, \text{ and } i \in [n]\}$ . Given a set  $\Sigma$  of TGDs, we write  $\text{pos}(\Sigma)$  instead of  $\text{pos}(\text{sch}(\Sigma))$ . The set of *affected positions* of  $\text{sch}(\Sigma)$ , denoted  $\text{aff}(\Sigma)$ , is inductively defined as follows:

- if there exists  $\sigma \in \Sigma$  and a variable  $x \in \text{var}_{\exists}(\sigma)$  at position  $\pi$ , then  $\pi \in \text{aff}(\Sigma)$ , and

- if there exists  $\sigma \in \Sigma$  and a variable  $x \in \text{front}(\sigma)$  in the body of  $\sigma$  only at positions of  $\text{aff}(\Sigma)$ , and  $x$  appears in the head of  $\sigma$  at position  $\pi$ , then  $\pi \in \text{aff}(\Sigma)$ .

Let  $\text{nonaff}(\Sigma) = \text{pos}(\Sigma) \setminus \text{aff}(\Sigma)$ . We can now classify the variables in the body of a TGD into harmless, harmful, and dangerous. Fix a TGD  $\sigma \in \Sigma$  and a variable  $x$  in  $\text{body}(\sigma)$ :

- $x$  is *harmless* if at least one occurrence of it appears in  $\text{body}(\sigma)$  at a position of  $\text{nonaff}(\Sigma)$ ,
- $x$  is *harmful* if it is not harmless, and
- $x$  is *dangerous* if it is harmful and belongs to  $\text{front}(\sigma)$ .

We are now ready to formally introduce *wardedness*.

**Definition 3.1 (Wardedness).** A set  $\Sigma$  of TGDs is *warded* if, for each TGD  $\sigma \in \Sigma$ , there are no dangerous variables in  $\text{body}(\sigma)$ , or there exists an atom  $\alpha \in \text{body}(\sigma)$ , called a *ward*, such that:

- all the dangerous variables in  $\text{body}(\sigma)$  occur in  $\alpha$ , and
- each variable of  $\text{var}(\alpha) \cap \text{var}(\text{body}(\sigma) \setminus \{\alpha\})$  is harmless.

We denote by **WARD** the class of all (finite) *warded* sets of TGDs. ■

The problem of CQ answering under *warded* sets of TGDs has been recently investigated in [16]:

**PROPOSITION 3.2.** *CQAns(WARD) is EXPTIME-complete in combined complexity, and PTIME-complete in data complexity.*

Note that [16] deals only with data complexity. However, it is implicit that the same algorithm provides an EXPTIME upper bound in combined complexity, while the lower bounds are inherited from Datalog since a set of Datalog rules (seen as TGDs) is *warded*.

**A Key Application.** One of the distinctive features of *wardedness*, which is crucial for the purposes of the Vadalogue system, is the fact that it can express every SPARQL query under the OWL 2 QL direct semantics entailment regime, which is inherited from the OWL 2 direct semantics entailment regime; for details, see [2, 16, 21]. Recall that SPARQL is the standard language for querying the Semantic Web,<sup>7</sup> while OWL 2 QL is a prominent profile of OWL 2.<sup>8</sup>

### 4 LIMITING RECURSION

We now focus on our main research question: can we limit the recursion allowed by *wardedness* in order to obtain a formalism that provides a convenient syntax for expressing useful recursive statements, and at the same time achieve space-efficiency? The above question has been extensively studied in the 1980s for Datalog programs, with *linear Datalog* being a key fragment that achieves a good balance between expressivity and complexity; see, e.g., [23, 24]. A Datalog program  $\Sigma$  is *linear* if, for each rule in  $\Sigma$ , its body contains at most one intensional predicate, i.e., a predicate that appears in the head of at least one rule of  $\Sigma$ . In other words, linear Datalog allows only for linear recursion, which is able to express many real-life recursive queries. However, for our purposes, linear recursion does not provide the convenient syntax that we are aiming at. After analyzing several real-life examples of *warded* sets of TGDs, provided by our industrial partners, we observed that the employed recursion goes beyond linear recursion. On the other hand, most of the examples coming from our industrial partners, use recursion

<sup>7</sup><http://www.w3.org/TR/rdf-sparql-query>

<sup>8</sup><https://www.w3.org/TR/owl2-overview/>

in a restrictive way: each TGD has at most one body atom whose predicate is mutually recursive with a predicate occurring in the head of the TGD. Interestingly, this more liberal version of linear recursion has been already investigated in the context of Datalog, and it is known as *piece-wise linear*; see, e.g., [1]. Does this type of recursion lead to the space-efficient fragment of warded sets of TGDs that we are looking for? The rest of this section is devoted to showing this rather involved result.

Let us start by formally defining the class of piece-wise linear sets of TGDs. To this end, we need to define when two predicates are mutually recursive, which in turn relies on the well-known notion of the predicate graph. The *predicate graph* of a set  $\Sigma$  of TGDs, denoted  $\text{pg}(\Sigma)$ , is a directed graph  $(V, E)$ , where  $V = \text{sch}(\Sigma)$ , and there exists an edge from a predicate  $P$  to a predicate  $R$ , i.e.,  $(P, R) \in E$ , iff there exists a TGD  $\sigma \in \Sigma$  such that  $P$  occurs in  $\text{body}(\sigma)$  and  $R$  occurs in  $\text{head}(\sigma)$ . Two predicates  $P, R \in \text{sch}(\Sigma)$  are *mutually recursive* (w.r.t.  $\Sigma$ ) if there exists a cycle in  $\text{pg}(\Sigma)$  that contains both  $P$  and  $R$  (i.e.,  $R$  is reachable from  $P$ , and vice versa). We are now ready to define piece-wise linearity for TGDs.

**Definition 4.1 (Piece-wise Linearity).** A set  $\Sigma$  of TGDs is *piece-wise linear* if, for each TGD  $\sigma \in \Sigma$ , there exists at most one atom in  $\text{body}(\sigma)$  whose predicate is mutually recursive with a predicate in  $\text{head}(\sigma)$ . Let PWL be the class of piece-wise linear sets of TGDs. ■

The main result of this section follows:

**THEOREM 4.2.**  $\text{CQAns}(\text{WARD} \cap \text{PWL})$  is PSPACE-complete in combined complexity, and NLOGSPACE-complete in data complexity.

The lower bounds are inherited from linear Datalog. The difficult task is to establish the upper bounds. This relies on a novel notion of proof tree, which is of independent interest. As we shall see, our notion of proof tree leads to space-bounded algorithms that allow us to show the upper bounds in Theorem 4.2, and also re-establish in a transparent way the upper bounds in Proposition 3.2. Moreover, in Section 6, we are going to use proof trees for studying the relative expressive power of (piece-wise linear) warded sets of TGDs.

## 4.1 Query Answering via Proof Trees

It is known that given a CQ  $q$  and a set  $\Sigma$  of TGDs, we can unfold  $q$  using the TGDs of  $\Sigma$  into an infinite union of CQs  $q_\Sigma$  such that, for every database  $D$ ,  $\text{cert}(q, D, \Sigma) = q_\Sigma(D)$ ; see, e.g., [15, 20]. Let us clarify that in our context, an unfolding, which is essentially a resolution step, is more complex than in the context of Datalog due to the existentially quantified variables in the head of TGDs. The intention underlying our notion of proof tree is to encode in a tree the sequence of CQs, generated during the unfolding of  $q$  with  $\Sigma$ , that leads to a certain CQ  $q'$  of  $q_\Sigma$  in such a way that each intermediate CQ, as well as  $q'$ , is carefully decomposed into smaller subqueries that form the nodes of the tree, while the root corresponds to  $q$  and the leaves to  $q'$ . As we shall see, if we focus on well-behaved classes of TGDs such as (piece-wise linear) warded sets of TGDs, we can establish upper bounds on the size of these subqueries, which in turn allow us to devise space-bounded algorithms for query answering. In what follows, we define the notion of proof tree (Definition 4.7), and establish its correspondence with query answering (Theorem 4.8). To this end, we need to introduce the main building blocks of a proof tree: chunk-based resolution

(Definition 4.3), a query decomposition technique (Definition 4.5), and the notion of specialization for CQs (Definition 4.6).

**Chunk-based Resolution.** Let  $A$  and  $B$  be non-empty set of atoms that mention only constants and variables. The sets  $A$  and  $B$  *unify* if there is a substitution  $\gamma$ , which is the identity on  $C$ , called *unifier* for  $A$  and  $B$ , such that  $\gamma(A) = \gamma(B)$ . A *most general unifier* (MGU) for  $A$  and  $B$  is a unifier for  $A$  and  $B$ , denoted  $\gamma_{A,B}$ , such that, for each unifier  $\gamma$  for  $A$  and  $B$ ,  $\gamma = \gamma' \circ \gamma_{A,B}$  for some substitution  $\gamma'$ . Notice that if two sets of atoms unify, then there exists always a MGU, which is unique (modulo variable renaming).

Given a CQ  $q(\bar{x})$  and a set of atoms  $S \subseteq \text{atoms}(q)$ , we call a variable  $y \in \text{var}(S)$  *shared* if  $y \in \bar{x}$ , or  $y \in \text{var}(\text{atoms}(q) \setminus S)$ . A *chunk unifier* of  $q$  with a TGD  $\sigma$  (where  $q$  and  $\sigma$  do not share variables) is a triple  $(S_1, S_2, \gamma)$ , where  $\emptyset \subset S_1 \subseteq \text{atoms}(q)$ ,  $\emptyset \subset S_2 \subseteq \text{head}(\sigma)$ , and  $\gamma$  is a unifier for  $S_1$  and  $S_2$  such that, for each  $x \in \text{var}(S_2) \cap \text{var}(\exists(\sigma))$ ,

- (1)  $\gamma(x) \notin C$ , i.e.,  $\gamma(x)$  is not constant, and
- (2)  $\gamma(x) = \gamma(y)$  implies  $y$  occurs in  $S_1$  and is not shared.

The chunk unifier  $(S_1, S_2, \gamma)$  is *most general* (MGCU) if  $\gamma$  is an MGU for  $S_1$  and  $S_2$ . Notice that the variables of  $\text{var}(\exists(\sigma))$  occurring in  $S_2$  unify (via  $\gamma$ ) only with non-shared variables of  $S_1$ . This ensures that  $S_1$  is a “chunk” of  $q$  that can be resolved as a whole via  $\sigma$  using  $\gamma$ .<sup>9</sup> Without the additional conditions on the substitution  $\gamma$ , we may get unsound resolution steps. Consider, e.g., the CQ and TGD

$$Q(x) \leftarrow R(x, y), S(y) \quad \text{and} \quad P(x') \rightarrow \exists y' R(x', y').$$

Resolving the atom  $R(x, y)$  in the query with the given TGD using the substitution  $\gamma = \{x \mapsto x', y \mapsto y'\}$  would be an unsound step since the shared variable  $y$  is lost. This is because  $y'$  is unified with the shared variable  $y$ . On the other hand,  $R(x, y), S(y)$  can be resolved with the TGD  $\sigma = P(x') \rightarrow \exists y' R(x', y'), S(y')$  using  $\gamma$ ; in fact, the chunk unifier is  $(\text{atoms}(q), \text{head}(\sigma), \gamma)$ .

**Definition 4.3 (Chunk-based Resolution).** Let  $q(\bar{x})$  be a CQ and  $\sigma$  a TGD. A  $\sigma$ -*resolvent* of  $q$  is a CQ  $q'(\bar{x}')$  with  $\text{body}(q') = \gamma((\text{atoms}(q) \setminus S_1) \cup \text{body}(\sigma))$  for a MGCU  $(S_1, S_2, \gamma)$  of  $q$  with  $\sigma$ . ■

**Query Decomposition.** As discussed above, the purpose of a proof tree is to encode a finite branch of the unfolding of a CQ  $q$  with a set  $\Sigma$  of TGDs, which is obtained by applying chunk-based resolution. Such a branch is a sequence  $q_0, \dots, q_n$  of CQs, where  $q = q_0$ , while, for each  $i \in [n]$ ,  $q_i$  is a  $\sigma$ -resolvent of  $q_{i-1}$  for some  $\sigma \in \Sigma$ . **Here is a simple example, which will serve as a running example in the rest of the section, that illustrates the notion of unfolding.**

**Example 4.4.** Consider the set  $\Sigma$  of TGDs consisting of

$R(x) \rightarrow \exists y T(y, x) \quad T(x, y), S(y, z) \rightarrow T(x, z) \quad T(x, y), P(y) \rightarrow G$   
and the CQ that simply asks whether  $G()$  is entailed, i.e., the CQ  $Q \leftarrow G()$ . Since the unfolding of  $q$  with  $\Sigma$  should give the correct answer for every input database, and thus for databases of the form

$$\{R(c^n), S(c^n, c^{n-1}), \dots, S(c^2, c^2), P(c^1)\} \quad \text{for some } n > 0,$$

one of its branches should be  $q = q_0, q_1, \dots, q_n$ , where

$$q_1 = Q \leftarrow T(x, y^1), P(y^1)$$

obtained by resolving  $q_0$  using the third TGD,

$$q_i = Q \leftarrow T(x, y^i), S(y^i, y^{i-1}), \dots, S(y^2, y^2), P(y^1),$$

<sup>9</sup>A similar notion known as piece unifier has been defined in [20].

for  $1 < i < n$ , obtained by resolving  $q_{i-1}$  using the second TGD,

$$q_n = Q \leftarrow R(y^n), S(y^n, y^{n-1}), \dots, S(y^2, y^2), P(y^1)$$

obtained by resolving  $q_{n-1}$  using the first TGD. ■

At this point, one may think that the proof tree that encodes the branch  $q_0, \dots, q_n$  of the unfolding of  $q$  with  $\Sigma$  is the finite labeled path  $v_0, \dots, v_n$ , where each  $v_i$  is labeled by  $q_i$ . However, another crucial goal of such a proof tree, which is not achieved via the naive path encoding, is to split each resolvent  $q_i$ , for  $i > 0$ , into smaller subqueries  $q_i^1, \dots, q_i^{n_i}$ , which are essentially the children of  $q_i$ , in such a way that they can be processed independently by resolution. The crux of this encoding is that it provides us with a mechanism for keeping the CQs that must be processed by resolution small. It should be clear from Example 4.4 that by following the naive path encoding, without splitting the resolvents into smaller subqueries, we may get CQs of unbounded size.

The key question here is how a CQ  $q$  can be decomposed into subqueries that can be processed independently. The subtlety is that, after splitting  $q$ , occurrences of the same variable may be separated into different subqueries. Thus, we need a way to ensure that a variable in  $q$ , which appears in different subqueries after the splitting, is indeed treated as the same variable, i.e., it has the same meaning. We deal with this issue by restricting the set of variables in  $q$  of which occurrences can be separated during the splitting step. In particular, we can only separate occurrences of an output variable. This relies on the convention that output variables correspond to fixed constant values of  $C$ , and thus their name is “frozen” and never renamed by subsequent resolution steps. Hence, we can separate occurrences of an output variable into different subqueries, i.e., different branches of the proof tree, without losing the connection between them.

Summing up, the idea underlying query decomposition is to split the CQ at hand into smaller subqueries that keep together all the occurrences of a non-output variable, but with the freedom of separating occurrences of an output variable.

**Definition 4.5 (Query Decomposition).** Given a CQ  $q(\bar{x})$ , a *decomposition* of  $q$  is a set of CQs  $\{q_1(\bar{y}_1), \dots, q_n(\bar{y}_n)\}$ , where  $n \geq 1$  and  $\bigcup_{i \in [n]} \text{atoms}(q_i) = \text{atoms}(q)$ , such that, for each  $i \in [n]$ :

- (1)  $\bar{y}_i$  is the restriction of  $\bar{x}$  on the variables in  $q_i$ , and
- (2) for every  $\alpha, \beta \in \text{atoms}(q)$ , if  $\alpha \in \text{atoms}(q_i)$  and  $\text{var}(\alpha) \cap \text{var}(\beta) \not\subseteq \bar{x}$ , then  $\beta \in \text{atoms}(q_i)$ . ■

**Query Specialization.** From the above discussion, one expects that a proof tree of a CQ  $q$  w.r.t. a set  $\Sigma$  of TGDs can be constructed by starting from  $q$ , which is the root, and applying two steps: resolution and decomposition. Unfortunately, this is not enough for our purposes as we may run into the following problem: some of the subqueries will mistakenly remain large since we have no way to realize that a non-output variable corresponds to a fixed constant value, which in turn allows us to “freeze” its name and separate different occurrences of it during the decomposition step. This is illustrated by Example 4.4. Observe that the size of the CQs  $\{q_i\}_{i>0}$  grows arbitrarily, while our query decomposition has no effect on them since there are Boolean queries, i.e., queries without output variables, and thus, we cannot split them into smaller subqueries.

The above issue can be solved by having an intermediate step

between resolution and decomposition, the so-called specialization step. A specialization of a CQ is obtained by converting some non-output variables of it into output variables, while keeping their name, or taking the name of an existing output variable.

**Definition 4.6 (Query Specialization).** Let  $q(\bar{x})$  be a CQ with  $\text{atoms}(q) = \{\alpha_1, \dots, \alpha_n\}$ . A *specialization* of  $q$  is a CQ

$$Q(\bar{x}, \bar{y}) \leftarrow \rho_{\bar{z}}(\alpha_1, \dots, \alpha_n)$$

where  $\bar{y}, \bar{z}$  are (possibly empty) disjoint tuples of non-output variables of  $q$ , and  $\rho_{\bar{z}}$  is a substitution from  $\bar{z}$  to  $\bar{x} \cup \bar{y}$ . ■

Consider, for example, the CQ  $q_1$  from Example 4.4

$$Q \leftarrow T(x, y^1), P(y^1)$$

obtained by resolving  $q = q_0$  using the third TGD. The query decomposition cannot split it onto smaller subqueries since the variable  $y^1$  is a non-output variable, and thus, all its occurrences should be kept together. We can consider the following specialization of  $q_1$

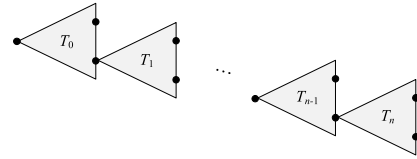
$$Q(y^1) \leftarrow T(x, y^1), P(y^1),$$

which simply converts  $y^1$  into an output variable, and now by query decomposition we can split it into the atomic queries

$$Q(y^1) \leftarrow T(x, y^1) \quad Q(y^1) \leftarrow P(y^1),$$

which represent the original query  $q_1$ .

**Proof Trees.** We are now ready to introduce our new notion of proof tree. But let us first explain the high-level idea by exploiting our running example. Consider the set  $\Sigma$  of TGDs and the CQ  $q$  from Example 4.4. The branch  $q_0, \dots, q_n$  of the unfolding of  $q$  with  $\Sigma$  given in Example 4.4 is encoded via a proof tree of the form



where each  $T_i$ , for  $i \in \{0, \dots, n\}$ , is a rooted tree with only two leaf nodes. The actual trees are depicted in Figure 1; the left one is  $T_0$ , the middle one is  $T_i$  for  $i \in [n-1]$ , while the right one is  $T_n$ . For each  $i \in \{0, \dots, n\}$ , the child of the root of  $T_i$  is obtained via resolution, then we specialize it by converting the variable  $y^i$  into an output variable, and then we decompose the specialized CQ into two subqueries. In the tree  $T_n$ , we also apply an additional resolution step in order to obtain the leaf node  $Q(y^n) \leftarrow R(y^n)$ . The underlined CQs are actually the subqueries that represent the CQ  $q_n$  of the unfolding. Indeed, the conjunction of the atoms occurring in the underlined CQs is precisely the CQ  $q_n$ .

We proceed to give the formal definition. Given a partition  $\pi = \{S_1, \dots, S_m\}$  of a set of variables, we write  $\text{eq}_\pi$  for the substitution that maps the variables of  $S_i$  to the same variable  $x_i$ , where  $x_i$  is a distinguished element of  $S_i$ . We should not forget the convention that output variables cannot be renamed, and thus, a resolution step should use a MGCU that preserves the output variables. In particular, given a CQ  $q$  and a TGD  $\sigma$ , a  $\sigma$ -resolvent of  $q$  is called *IDO* if the underlying MGCU uses a substitution that is the identity on the output variables of  $q$  (hence the name IDO). Finally, given a TGD  $\sigma$  and some arbitrary object  $o$  (e.g.,  $o$  can be the node of a



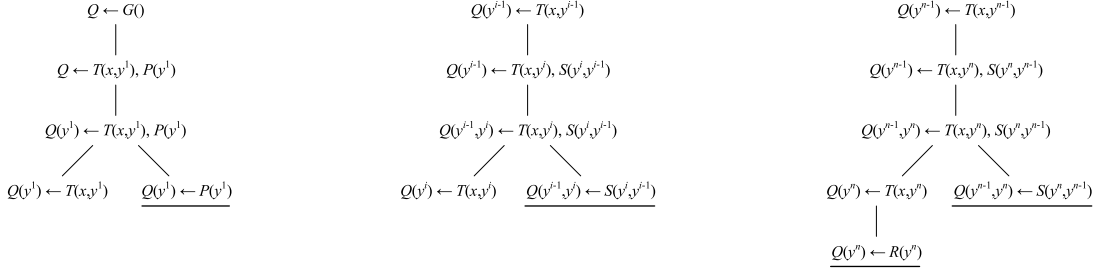


Figure 1: Partial trees of the proof tree the encodes the branch  $q = q_0, \dots, q_n$  of the unfolding of  $q$  with  $\Sigma$  from Example 4.4.

tree, or an integer number), we write  $\sigma_o$  for the TGD obtained by renaming each variable  $x$  in  $\sigma$  into  $x_o$ . This is a simple mechanism for uniformly renaming the variables of a TGD in order to avoid undesirable clutter among variables during a resolution step.

**Definition 4.7 (Proof Tree).** Let  $q(\bar{x})$  be a CQ with  $\text{atoms}(q) = \{\alpha_1, \dots, \alpha_n\}$ , and  $\Sigma$  a set of TGDs. A *proof tree* of  $q$  w.r.t.  $\Sigma$  is a triple  $\mathcal{P} = (T, \lambda, \pi)$ , where  $T = (V, E)$  is a finite rooted tree,  $\lambda$  a labeling function that assigns a CQ to each node of  $T$ , and  $\pi$  a partition of  $\bar{x}$ , such that, for each node  $v \in V$ :

- (1) If  $v$  is the root node of  $T$ , then  $\lambda(v)$  is the CQ  $Q(\text{eq}_\pi(\bar{x})) \leftarrow \text{eq}_\pi(\alpha_1, \dots, \alpha_m)$ .
- (2) If  $v$  has only one child  $u$ ,  $\lambda(u)$  is an IDO  $\sigma_v$ -resolvent of  $\lambda(v)$  for some  $\sigma \in \Sigma$ , or a specialization of  $\lambda(v)$ .
- (3) If  $v$  has  $k > 1$  children  $u_1, \dots, u_k$ , then  $\{\lambda(u_1), \dots, \lambda(u_k)\}$  is a decomposition of  $\lambda(v)$ .

Assuming that  $v_1, \dots, v_m$  are the leaf nodes of  $T$ , the *CQ induced by  $\mathcal{P}$*  is defined as

$$Q(\text{eq}_\pi(\bar{x})) \leftarrow \alpha_1, \dots, \alpha_\ell,$$

where  $\{\alpha_1, \dots, \alpha_\ell\} = \bigcup_{i \in [m]} \text{atoms}(\lambda(v_i))$ . ■

The purpose of the partition  $\pi$  is to indicate that some output variables correspond to the same constant value – this is why variables in the same set of  $\pi$  are unified via the substitution  $\text{eq}_\pi$ . This unification step is crucial in order to safely use, in subsequent resolution steps, substitutions that are the identity on the output variables. If we omit this initial unification step, we may lose important resolution steps, and thus being incomplete for query answering purposes. The main result of this section, which exposes the connection between proof trees and CQ answering, follows. By abuse of notation, we write  $\mathcal{P}$  for the CQ induced by  $\mathcal{P}$ .

**THEOREM 4.8.** Consider a database  $D$ , a set  $\Sigma$  of TGDs, a CQ  $q(\bar{x})$ , and  $\bar{c} \in \text{dom}(D)^{|\bar{x}|}$ . The following are equivalent:

- (1)  $\bar{c} \in \text{cert}(q, D, \Sigma)$ .
- (2) There exists a proof tree  $\mathcal{P}$  of  $q$  w.r.t.  $\Sigma$  such that  $\bar{c} \in \mathcal{P}(D)$ .

The proof of the above result relies on the soundness and completeness of chunk-based resolution. Given a set  $\Sigma$  of TGDs and a CQ  $q(\bar{x})$ , by exhaustively applying chunk-based resolution, we can construct a (possibly infinite) union of CQs  $q_\Sigma$  such that, for every database  $D$ ,  $\text{cert}(q, D, \Sigma) = q_\Sigma(D)$ ; implicit in [15, 20]. In other words, given a tuple  $\bar{c} \in \text{dom}(D)^{|\bar{x}|}$ ,  $\bar{c} \in \text{cert}(q, D, \Sigma)$  iff there exists a CQ  $q'(\bar{x})$  in  $q_\Sigma$  such that  $\bar{c} \in q'(D)$ . It is now not difficult to show

that the latter statement is equivalent to the existence of a proof tree  $\mathcal{P}$  of  $q$  w.r.t.  $\Sigma$  such that  $\bar{c} \in \mathcal{P}(D)$ , and the claim follows.

## 4.2 Well-behaved Proof Trees

Theorem 4.8 states that checking whether a tuple  $\bar{c}$  is a certain answer boils down to deciding whether there exists a proof tree  $\mathcal{P}$  such that  $\bar{c}$  is an answer to the CQ induced by  $\mathcal{P}$  over the given database. Of course, the latter is an undecidable problem in general. However, if we focus on (piece-wise linear) warded sets of TGDs, it suffices to check for the existence of a well-behaved proof tree with certain syntactic properties, which in turn allows us to devise a decision procedure. We proceed to make this more precise. For technical clarity, we assume, w.l.o.g., TGDs with only one atom in the head since we can always convert a warded set of TGDs into one with single-atom heads, while certain answers are preserved; for the transformation see, e.g., [11].

**Piece-wise Linear Warded Sets of TGDs.** For piece-wise linear warded sets of TGDs, we can strengthen Theorem 4.8 by focussing on a certain class of proof trees that enjoy two syntactic properties: (i) they have a path-like structure, and (ii) the size of the CQs that label their nodes is bounded by a polynomial. The first property is formalized via linear proof trees. Let  $\mathcal{P} = (T, \lambda, \pi)$ , where  $T = (V, E)$ , be a proof tree of a CQ  $q$  w.r.t. a set  $\Sigma$  of TGDs. We call  $\mathcal{P}$  *linear* if, for each node  $v \in V$ , there exists at most one node  $u \in V$  such that  $(v, u) \in E$  and  $u$  is not a leaf in  $T$ , i.e.,  $v$  has at most one child that is not a leaf. For example, the proof tree given above, which consists of the partial trees depicted in Figure 1, is linear. The second property relies on the notion of node-width of a proof tree. Formally, the *node-width* of  $\mathcal{P}$  is

$$\text{nwd}(\mathcal{P}) := \max_{v \in V} \{|\lambda(v)|\},$$

i.e., the size of the largest CQ that labels a node of  $T$ .

Before we strengthen Theorem 4.8, let us define the polynomial that will allow us to bound the node-width of the linear proof trees that we need to consider. This polynomial relies on the notion of predicate level. Consider a set  $\Sigma$  of TGDs. For a predicate  $P \in \text{sch}(\Sigma)$ , we write  $\text{rec}(P)$  for the set of predicates of  $\text{sch}(\Sigma)$  that are mutually recursive to  $P$  according to  $\text{pg}(\Sigma) = (V, E)$ . Let  $\ell_\Sigma : \text{sch}(\Sigma) \rightarrow \mathbb{N}$  be the unique function that satisfies

$$\ell_\Sigma(P) = \max\{\ell_\Sigma(R) \mid (R, P) \in E, R \notin \text{rec}(P)\} + 1,$$

with  $\ell_\Sigma(P)$  being the *level* (w.r.t.  $\Sigma$ ) of  $P$ , for each  $P \in \text{sch}(\Sigma)$ . We can now define the polynomial

$$f_{\text{WARD} \cap \text{PWL}}(q, \Sigma) := (|q| + 1) \cdot \max_{P \in \text{sch}(\Sigma)} \{\ell_\Sigma(P)\} \cdot \max_{\sigma \in \Sigma} \{|\text{body}(\sigma)|\}.$$

We can now strengthen Theorem 4.8. But let us first clarify that, in the case of piece-wise linear warded sets of TGDs, apart from only one atom in the head, we also assume, w.l.o.g., that the level of a predicate in the body of TGD  $\sigma$  is  $k$  or  $k - 1$ , where  $k$  is the level of the predicate in the head of  $\sigma$ . The following holds:

**THEOREM 4.9.** *Consider a database  $D$ , a set  $\Sigma \in \text{WARD} \cap \text{PWL}$  of TGDs, a CQ  $q(\bar{x})$ , and  $\bar{c} \in \text{dom}(D)^{|\bar{x}|}$ . The following are equivalent:*

- (1)  $\bar{c} \in \text{cert}(q, D, \Sigma)$ .
- (2) *There is a linear proof tree  $\mathcal{P}$  of  $q$  w.r.t.  $\Sigma$  with  $\text{nwd}(\mathcal{P}) \leq f_{\text{WARD} \cap \text{PWL}}(q, \Sigma)$  such that  $\bar{c} \in \mathcal{P}(D)$ .*

**Warded sets of TGDs.** Now, in the case of arbitrary warded sets of TGDs, we cannot focus only on linear proof trees. Nevertheless, we can still bound the node-width of the proof trees that we need to consider by the following polynomial, which, unsurprisingly, does not rely anymore on the notion of predicate level:

$$f_{\text{WARD}}(q, \Sigma) := 2 \cdot \max \left\{ |q|, \max_{\sigma \in \Sigma} \{|\text{body}(\sigma)|\} \right\}.$$

Theorem 4.8 can be strengthened as follows:

**THEOREM 4.10.** *Consider a database  $D$ , a set  $\Sigma \in \text{WARD}$  of TGDs, a CQ  $q(\bar{x})$ , and  $\bar{c} \in \text{dom}(D)^{|\bar{x}|}$ . The following are equivalent:*

- (1)  $\bar{c} \in \text{cert}(q, D, \Sigma)$ .
- (2) *There exists a proof tree  $\mathcal{P}$  of  $q$  w.r.t.  $\Sigma$  with  $\text{nwd}(\mathcal{P}) \leq f_{\text{WARD}}(q, \Sigma)$  such that  $\bar{c} \in \mathcal{P}(D)$ .*

**A Proof Sketch.** Let us now provide some details on how Theorems 4.9 and 4.10 are shown. For both theorems, (2) implies (1) readily follows from Theorem 4.8. We thus focus on the other direction. The main ingredients of the proof can be described as follows:

- We introduce the auxiliary notion of chase tree, which can be seen as a concrete instantiation of a proof tree. It serves as an intermediate structure between proof trees and chase derivations, which allows us to use the chase as our underlying technical tool. Note that the notions of linearity and node-width can be naturally defined for chase trees.
- We then show that, if the given tuple of constants  $\bar{c}$  is a certain answer to the given CQ  $q$  w.r.t. the given database  $D$  and (piece-wise linear) warded set  $\Sigma$  of TGDs, then there exists a (linear) chase tree for the image of  $q$  to  $\text{chase}(D, \Sigma)$  such that its node-width respects the bounds given in the above theorems (Lemma 4.12).
- We finally show that the existence of a (linear) chase tree for the image of  $q$  to  $\text{chase}(D, \Sigma)$  with node-width at most  $m$  implies the existence of a (linear) proof tree  $\mathcal{P}$  of  $q$  w.r.t.  $\Sigma$  with node-width at most  $m$  such that  $\bar{c} \in \mathcal{P}(D)$  (Lemma 4.13).

Let us make the above description more formal. In order to introduce the notion of chase tree, we first need to recall the notion of chase graph, then introduce the notion of unravelling of the chase graph, and finally introduce the notions of unfolding and decomposition for sets of atoms in the unravelling of the chase graph.

Fix a chase sequence  $\delta = (I_i \langle \sigma_i, h_i \rangle I_{i+1})_{i \geq 0}$  for a database  $D$  under a set  $\Sigma$  of TGDs. The *chase graph* for  $D$  and  $\Sigma$  (w.r.t.  $\delta$ ) is a directed edge-labeled graph  $\mathcal{G}^{D, \Sigma} = (V, E, \lambda)$ , where  $V = \text{chase}(D, \Sigma)$ , and an edge  $(\alpha, \beta)$  labeled with  $(\sigma_k, h_k)$  belongs to  $E$  iff  $\alpha \in h_k(\text{body}(\sigma_k))$  and  $\beta \in I_{k+1} \setminus I_k$ , for some  $k \geq 0$ . In other words,  $\alpha$  has an edge to  $\beta$  if  $\beta$  is derived using  $\alpha$ , and if  $\beta$  is new in the sense that it has not been derived before. Notice that  $\mathcal{G}^{D, \Sigma}$  has no directed cycles. Notice also that  $\mathcal{G}^{D, \Sigma}$  depends on  $\delta$  – however, we can assume a fixed sequence  $\delta$  since, as discussed in Section 2, every chase sequence is equally useful for our purposes.

We now discuss the notion of unravelling of the chase graph; due to space reasons, we keep this discussion informal. Given a set  $\Theta \subseteq \text{chase}(D, \Sigma)$ , the *unravelling* of  $\mathcal{G}^{D, \Sigma}$  around  $\Theta$  is a directed node- and edge-labeled forest  $\mathcal{G}_\Theta^{D, \Sigma}$  that has a tree for each  $\alpha \in \Theta$  whose branches are backward-paths in  $\mathcal{G}$  from  $\alpha$  to a database atom. Intuitively,  $\mathcal{G}_\Theta^{D, \Sigma}$  is a forest-like reorganization of the atoms of  $\text{chase}(D, \Sigma)$  that are needed to derive  $\Theta$ . Due to its forest-like shape, it may contain multiple copies of atoms of  $\text{chase}(D, \Sigma)$ . The edges between nodes are labeled by pairs  $(\sigma, h)$  just like in  $\mathcal{G}^{D, \Sigma}$ , while the nodes are labeled by atoms and, importantly, the atoms along the paths in  $\mathcal{G}^{D, \Sigma}$  may be duplicated and labeled nulls are given new names. We write  $U(\mathcal{G}_\Theta^{D, \Sigma}, \Theta)$  for the set of all atoms that appear as labels in  $\mathcal{G}_\Theta^{D, \Sigma}$ , and  $\text{succ}_{\sigma, h}(\alpha)$  for the set of children of  $\alpha$  whose incoming edge is labeled with  $(\sigma, h)$ . It is important to say that there exists a homomorphism  $h_\Theta$  that maps  $\Theta$  to  $U(\mathcal{G}_\Theta^{D, \Sigma}, \Theta)$ .

Let us now introduce the notions of unfolding and decomposition. For sets  $\Gamma, \Gamma' \subseteq U(\mathcal{G}_\Theta^{D, \Sigma}, \Theta)$ , we say that  $\Gamma'$  is an *unfolding* of  $\Gamma$ , if there are  $\alpha \in \Gamma$  and  $\beta_1, \dots, \beta_k \in U(\mathcal{G}_\Theta^{D, \Sigma}, \Theta)$  such that

- (1)  $\text{succ}_{\sigma, h}(\alpha) = \{\beta_1, \dots, \beta_k\}$ , for some  $\sigma \in \Sigma$  and  $h$ ,
- (2) for every null that occurs in  $\alpha$ , either it does not appear in  $\Gamma \setminus \{\alpha\}$ , or it appears in  $\{\beta_1, \dots, \beta_k\}$ , and
- (3)  $\Gamma' = (\Gamma \setminus \{\alpha\}) \cup \{\beta_1, \dots, \beta_k\}$ .

Let  $\Gamma \subseteq U(\mathcal{G}_\Theta^{D, \Sigma}, \Theta)$  be a non-empty set. A *decomposition* of  $\Gamma$  is a set  $\{\Gamma_1, \dots, \Gamma_n\}$ , where  $n \geq 1$ , of non-empty subsets of  $\Gamma$  such that (i)  $\Gamma = \bigcup_{i \in [n]} \Gamma_i$ , and (ii)  $i \neq j$  implies that  $\Gamma_i$  and  $\Gamma_j$  do not share a labeled null. We can now define the key notion of chase tree:

**Definition 4.11 (Chase Tree).** Consider a database  $D$ , a set  $\Sigma$  of TGDs, and a set  $\Theta \subseteq \text{chase}(D, \Sigma)$ . A *chase tree* for  $\Gamma \subseteq U(\mathcal{G}_\Theta^{D, \Sigma}, \Theta)$  (w.r.t.  $\mathcal{G}_\Theta^{D, \Sigma}$ ) is a pair  $C = (T, \lambda)$ , where  $T = (V, E)$  is a finite rooted tree, and  $\lambda$  a labeling function that assigns a subset of  $U(\mathcal{G}_\Theta^{D, \Sigma}, \Theta)$  to each node of  $T$ , such that, for each  $v \in V$ :

- (1) If  $v$  is the root node of  $T$ , then  $\lambda(v) = \Gamma$ .
- (2) If  $v$  has only one child  $u$ , then  $\lambda(u)$  is an unfolding of  $\lambda(v)$ .
- (3) If  $v$  has  $k > 1$  children  $u_1, \dots, u_k$ , then  $\{\lambda(u_1), \dots, \lambda(u_k)\}$  is a decomposition of  $\lambda(v)$ .
- (4) If  $v$  is a leaf node, then  $\lambda(v) \subseteq D$ .

The *node-width* of  $C$  is  $\text{nwd}(C) := \max_{v \in V} \{|\lambda(v)|\}$ . Moreover, we say that  $C$  is *linear* if, for each node  $v \in V$ , there exists at most one  $u \in V$  such that  $(v, u) \in E$  and  $u$  is not a leaf. ■

We can now state our auxiliary technical lemmas. In what follows, fix a database  $D$ , and a set  $\Sigma$  of TGDs.

**LEMMA 4.12.** *Let  $\Theta \subseteq \text{chase}(D, \Sigma)$  and  $\Gamma \subseteq U(\mathcal{G}_\Theta^{D, \Sigma}, \Theta)$ . Then:*

- (1) *If  $\Sigma \in \text{WARD} \cap \text{PWL}$ , then there exists a linear chase tree  $C$  for  $\Gamma$  such that  $\text{nwd}(C) \leq f_{\text{WARD} \cap \text{PWL}}(\Gamma, \Sigma)$ .*



**Input:**  $D, \Sigma \in \text{WARD} \cap \text{PWL}$ ,  $q(\bar{x})$ ,  $\bar{c} \in \text{dom}(D)^{|\bar{x}|}$   
**Output:** Accept if  $\bar{c} \in \text{cert}(q, D, \Sigma)$ ; otherwise, Reject

$p := Q \leftarrow \alpha_1, \dots, \alpha_n$  with  $\text{atoms}(q(\bar{c})) = \{\alpha_1, \dots, \alpha_n\}$

**repeat**

- guess**  $op \in \{r, d, s\}$ 
  - if**  $op = r$  **then**
    - guess** a TGD  $\sigma \in \Sigma$
    - if**  $\text{mgcu}(p, \sigma) = \emptyset$  **then**
      - Reject
    - else**
      - guess**  $U \in \text{mgcu}(p, \sigma)$
      - if**  $|p[\sigma, U]| > f_{\text{WARD} \cap \text{PWL}}(q, \Sigma)$  **then**
        - Reject
      - else**
        - $p' := p[\sigma, U]$
  - if**  $op = d$  **then**
    - $p' := p[-D]$
  - if**  $op = s$  **then**
    - guess**  $V \subseteq \text{var}(p)$  and  $\gamma : V \rightarrow \text{dom}(D)$
    - $p' := \gamma(p)$

$p := p'$

**until**  $\text{atoms}(p) \subseteq D$ ;  
**return** Accept

(2) If  $\Sigma \in \text{WARD}$ , then there exists a chase tree  $C$  for  $\Gamma$  such that  $\text{nwd}(C) \leq f_{\text{WARD}}(\Gamma, \Sigma)$ .

The next technical lemma exposes the connection between chase trees and proof trees:

**LEMMA 4.13.** Consider a set  $\Theta \subseteq \text{chase}(D, \Sigma)$ , and let  $q'(\bar{x})$  be a CQ and  $\bar{c}$  a tuple of constants such that  $h'(\text{atoms}(q')) \subseteq U(\mathcal{G}^{D, \Sigma}, \Theta)$  and  $h'(\bar{x}) = \bar{c}$ , for some homomorphism  $h'$ . If there is a (linear) chase tree  $C$  for  $h'(\text{atoms}(q'))$  with  $\text{nwd}(C) \leq m$ , then there is a (linear) proof tree  $\mathcal{P}$  for  $q'$  w.r.t.  $\Sigma$  such that  $\text{nwd}(\mathcal{P}) \leq m$  and  $\bar{c} \in \mathcal{P}(D)$ .

We now show Theorem 4.9, while Theorem 4.10 can be shown analogously. Consider a CQ  $q(\bar{x})$  and a tuple  $\bar{c} \in \text{dom}(D)^{|\bar{x}|}$  such that  $\bar{c} \in \text{cert}(q, D, \Sigma)$ . We need to show that if  $\Sigma \in \text{WARD} \cap \text{PWL}$ , then there exists a linear proof tree  $\mathcal{P}$  of  $q$  w.r.t.  $\Sigma$  with  $\text{nwd}(\mathcal{P}) \leq f_{\text{WARD} \cap \text{PWL}}(q, \Sigma)$  such that  $\bar{c} \in \mathcal{P}(D)$ . By hypothesis, there is a homomorphism  $h$  such that  $h(\text{atoms}(q)) \subseteq \text{chase}(D, \Sigma)$  and  $h(\bar{x}) = \bar{c}$ . Let  $\Theta_q$  be the set of atoms  $h(\text{atoms}(q))$ . Recall that there is a homomorphism  $h_{\Theta_q}$  that maps  $\Theta_q$  to  $U(\mathcal{G}^{D, \Sigma}, \Theta_q)$ . Thus, the homomorphism  $h' = h_{\Theta_q} \circ h$  is such that  $h'(\text{atoms}(q)) \subseteq U(\mathcal{G}^{D, \Sigma}, \Theta_q)$  and  $h'(\bar{x}) = \bar{c}$ . By Lemma 4.12, there exists a chase tree  $C$  for  $h'(\text{atoms}(q))$  with  $\text{nwd}(C) \leq f_{\text{WARD} \cap \text{PWL}}(h'(\text{atoms}(q)), \Sigma)$ . By Lemma 4.13, there exists a linear proof tree  $\mathcal{P}$  of  $q$  w.r.t.  $\Sigma$  with  $\text{nwd}(\mathcal{P}) \leq f_{\text{WARD} \cap \text{PWL}}(h'(\text{atoms}(q)), \Sigma) \leq f_{\text{WARD} \cap \text{PWL}}(q, \Sigma)$  such that  $\bar{c} \in \mathcal{P}(D)$ , and the claim follows.

### 4.3 Complexity Analysis

We now have all the tools for showing that CQ answering under piece-wise linear warded sets of TGDs is in PSPACE in combined complexity, and in NLOGSPACE in data complexity, and also for

re-establishing the complexity of warded sets of TGDs (see Proposition 3.2) in a more transparent way than the approach of [16].

**The Case of CQAns(WARD  $\cap$  PWL).** Given a database  $D$ , a set  $\Sigma \in \text{WARD} \cap \text{PWL}$  of TGDs, a CQ  $q(\bar{x})$ , and a tuple  $\bar{c} \in \text{dom}(D)^{|\bar{x}|}$ , by Theorem 4.9, our problem boils down to checking whether there exists a linear proof tree  $\mathcal{P}$  of  $q$  w.r.t.  $\Sigma$  with  $\text{nwd}(\mathcal{P}) \leq f_{\text{WARD} \cap \text{PWL}}(q, \Sigma)$  such that  $\bar{c} \in \mathcal{P}(D)$ . This can be easily checked via a space-bounded algorithm that is trying to build such a proof tree in a level-by-level fashion. Essentially, the algorithm builds the  $i$ -th level from the  $(i - 1)$ -th level of the proof tree by non-deterministically applying the operations introduced above, i.e., resolution, decomposition and specialization.

The algorithm is depicted in the box above. Here is a semi-formal description of it. The first step is to store in  $p$  the Boolean CQ obtained after instantiating the output variables of  $q$  with  $\bar{c}$ . The rest of the algorithm is an iterative procedure that non-deterministically constructs  $p'$  (the  $i$ -th level) from  $p$  (the  $(i - 1)$ -th level) until it reaches a level that is a subset of the database  $D$ . Notice that  $p$  and  $p'$  always hold one CQ since at each level of a linear proof tree only one node has a child, while all the other nodes are leaves, which essentially means that their atoms appear in the database  $D$ . At each iteration, the algorithm constructs  $p'$  from  $p$  by applying resolution (r), decomposition (d), or specialization (s):

**Resolution.** It guesses a TGD  $\sigma \in \Sigma$ . If the set  $\text{mgcu}(p, \sigma)$ , i.e., the set of all MGCU of  $p$  with  $\sigma$ , is empty, then rejects; otherwise, it guesses  $U \in \text{mgcu}(p, \sigma)$ . If the size of the  $\sigma$ -resolvent of  $p$  obtained via  $U$ , denoted  $p[\sigma, U]$ , does not exceed the bound given by Theorem 4.9, then it assigns  $p[\sigma, U]$  to  $p'$ ; otherwise, it rejects. Recall that during a resolution step we need to rename variables in order to avoid undesirable clutter. However, we cannot blindly use new variables at each step since this will explode the space used by algorithm. Instead, we should reuse variables that have been lost due to their unification with an existentially quantified variable. A simple analysis shows that we only need polynomially many variables, while this polynomial depends only on  $q$  and  $\Sigma$ .

**Decomposition.** It deletes from  $p$  the atoms that occur in  $D$ , and it assigns the obtained CQ  $p[-D]$  to  $p'$ . Notice that  $p[-D]$  may be empty in case  $\text{atoms}(p) \subseteq D$ . Essentially, the algorithm decomposes  $p$  in such a way that the subquery of  $p$  consisting of  $\text{atoms}(p) \cap D$  forms a child of  $p$  that is a leaf, while the subquery consisting of  $\text{atoms}(p) \setminus D$  is the non-leaf child.

**Specialization.** It assigns to  $p'$  a specialized version of  $p$ , where some variables are instantiated by constants of  $\text{dom}(D)$ . Notice that the convention that output variables correspond to constants is implemented by directly instantiating them with actual constants from  $\text{dom}(D)$ .

After constructing  $p'$ , the algorithm assigns it to  $p$ , and this ends one iteration. If  $\text{atoms}(p) \subseteq D$ , then a linear proof tree  $\mathcal{P}$  such that  $\bar{c} \in \mathcal{P}(D)$  has been found, and the algorithm accepts; otherwise, it proceeds with the next iteration.

It is easy to see that the algorithm uses polynomial space in general. Moreover, in case the set of TGDs and the CQ are fixed, the algorithm uses logarithmic space, which is the space needed for representing constantly many elements of  $\text{dom}(D)$ ; each element

of  $\text{dom}(D)$  can be represented using logarithmically many bits. The desired upper bounds claimed in Theorem 4.2 follow.

**The Case of CQAns(WARD).** The non-deterministic algorithm discussed above cannot be directly used for warded sets of TGDs since it is not enough to search for a linear proof tree as in the case of piece-wise linear warded sets of TGDs. However, by Theorem 4.10, we can search for a proof tree that has bounded node-width. This allows us to devise a space-bounded algorithm, which is similar in spirit as the one presented above, with the crucial difference that it constructs in a level-by-level fashion the branches of the proof tree in parallel universal computations using alternation. Since this alternating algorithm uses polynomial space in general, and logarithmic space when the set of TGDs and the CQ are fixed, we immediately get an EXPTIME upper bound in combined, and a PTIME upper bound in data complexity. This confirms Proposition 3.2 established in [16]. However, our new algorithm is significantly simpler than the one employed in [16], while Theorem 4.10 reveals the main property of warded sets of TGDs that leads to the desirable complexity upper bounds.

## 5 A JUSTIFIED COMBINATION

It is interesting to observe that the class of piece-wise linear warded sets of TGDs generalizes the class of *intensionally linear* sets of TGDs, denoted IL, where each TGD has at most one body atom whose predicate is intensional. Therefore, Theorem 4.2 immediately implies that CQAns(IL) is PSPACE-complete in combined complexity, and NLOGSPACE-complete in data complexity. Notice that IL generalizes linear Datalog, which is also PSPACE-complete in combined complexity, and NLOGSPACE-complete in data complexity. Thus, we can extend linear Datalog by allowing existentially quantified variables in rule heads, which essentially leads to IL, without affecting the complexity of query answering.

At this point, one maybe tempted to think that the same holds for piece-wise linear Datalog, i.e., we can extend it with existentially quantified variables in rule heads, which leads to PWL, without affecting the complexity of query answering, that is, PSPACE-complete in combined, and NLOGSPACE-complete in data complexity. However, if this is the case, then wardedness becomes redundant since the formalism that we are looking for is the class of piece-wise linear sets of TGDs, without the wardedness condition. It turned out that this is not the case. To our surprise, the following holds:

**THEOREM 5.1.** *CQAns(PWL) is undecidable in data complexity.*

To show the above result we exploit an undecidable tiling problem [8]. A *tiling system* is a tuple  $\mathbb{T} = (T, L, R, H, V, a, b)$ , where  $T$  is a finite set of tiles,  $L, R \subseteq T$  are special sets of left and right border tiles, respectively, with  $L \cap R = \emptyset$ ,  $H, V \subseteq T^2$  are the horizontal and vertical constraints, and  $a, b$  are distinguished tiles of  $T$  called the start and the finish tile, respectively. A *tiling* for  $\mathbb{T}$  is a function  $f: [n] \times [m] \rightarrow T$ , for some  $n, m > 0$ , such that  $f(1, 1) = a$ ,  $f(1, m) = b$ ,  $f(1, i) \in L$  and  $f(n, i) \in R$ , for every  $i \in [m]$ , and  $f$  respects the horizontal and vertical constraints. In other words, the first and the last rows of a tiling for  $\mathbb{T}$  start with  $a$  and  $b$ , respectively, while the leftmost and rightmost columns contain only tiles from  $L$  and  $R$ , respectively. We reduce from the UnboundedTiling problem, that is, given a tiling system  $\mathbb{T}$ , decide whether there is a tiling for

$\mathbb{T}$ . Given a tiling system  $\mathbb{T} = (T, L, R, H, V, a, b)$ , the goal is to construct in polynomial time a database  $D_{\mathbb{T}}$ , a set of TGDs  $\Sigma \in \text{PWL}$ , and a Boolean CQ  $q$ , such that  $\mathbb{T}$  has a tiling iff  $() \in \text{cert}(q, D_{\mathbb{T}}, \Sigma)$ ;  $()$  is the empty tuple. Note that  $\Sigma$  and  $q$  should not depend on  $\mathbb{T}$ .

**The Database  $D_{\mathbb{T}}$ .** It simply stores the tiling system  $\mathbb{T}$ :

$$\begin{aligned} & \{\text{Tile}(t) \mid t \in T\} \cup \{\text{Left}(t) \mid t \in L\} \cup \{\text{Right}(t) \mid t \in R\} \\ & \cup \{H(t, t') \mid (t, t') \in H\} \cup \{V(t, t') \mid (t, t') \in V\} \\ & \cup \{\text{Start}(a), \text{Finish}(b)\}. \end{aligned}$$

**The Set of TGDs  $\Sigma$ .** It is responsible for generating all the candidate tilings for  $\mathbb{T}$ , i.e., tilings without the condition  $f(1, m) = b$ , of arbitrary width and depth. Whether there exists a candidate tiling for  $\mathbb{T}$  that satisfies the condition  $f(1, m) = b$  will be checked by the CQ  $q$ . The set  $\Sigma$  essentially implements the following idea: construct rows of size  $\ell$  from rows of size  $\ell - 1$ , for  $\ell > 1$ , that respect the horizontal constraints, and then construct all the candidate tilings by combining compatible rows, i.e., rows that respect the vertical constraints. A row  $r$  is encoded as an atom  $\text{Row}(p, c, s, e)$ , where  $p$  is the id of the row from which  $r$  has been obtained, i.e., the previous one,  $c$  is the id of  $r$ , i.e., the current one,  $s$  is the starting tile of  $r$ , and  $e$  is the ending tile of  $r$ . We write  $\text{Row}(c, c, s, s)$  for rows consisting of a single tile, which do not have a previous row (hence the id of the previous row coincides with the id of the current row), and the starting tile is the same as the ending tile. The following two TGDs construct all the rows that respect the horizontal constraints:

$$\begin{aligned} \text{Tile}(x) & \rightarrow \exists z \text{Row}(z, z, x, x), \\ \text{Row}(\_, x, y, z), H(z, w) & \rightarrow \exists u \text{Row}(x, u, y, w). \end{aligned}$$

Analogously to Prolog, we write “ $\_$ ” for a “don’t-care” variable that occurs only once in the TGD. The next set of TGDs constructs all the pairs of compatible rows, i.e., pairs of rows  $(r_1, r_2)$  such that we can place  $r_2$  below  $r_1$  without violating the vertical constraints. This is done again inductively as follows:

$$\begin{aligned} \text{Row}(x, x, y, y), \text{Row}(x', x', y', y'), V(y, y') & \rightarrow \text{Comp}(x, x'), \\ \text{Row}(x, y, \_, z), \text{Row}(x', y', \_, z'), & \\ \text{Comp}(x, x'), V(z, z') & \rightarrow \text{Comp}(y, y'). \end{aligned}$$

We finally compute all the candidate tilings, together with their bottom-left tile, using the following two TGDs:

$$\begin{aligned} \text{Row}(\_, x, y, z), \text{Start}(y), \text{Right}(z) & \rightarrow \text{CTiling}(x, y), \\ \text{CTiling}(x, \_), \text{Row}(\_, y, z, w), \text{Comp}(x, y), & \\ \text{Left}(z), \text{Right}(w) & \rightarrow \text{CTiling}(y, z). \end{aligned}$$

This concludes the definition of  $\Sigma$ .

**The Boolean CQ  $q$ .** Recall that  $q$  is responsible for checking whether there exists a candidate tiling such that its bottom-left tile is  $b$ . This can be easily done via the query

$$Q \leftarrow \text{CTiling}(x, y), \text{Finish}(y).$$

By construction, the set  $\Sigma$  of TGDs belongs to PWL. Moreover, there is a tiling for  $\mathbb{T}$  iff  $() \in \text{cert}(q, D_{\mathbb{T}}, \Sigma)$ , and Theorem 5.1 follows.

## 6 EXPRESSIVE POWER

A class of TGDs naturally gives rise to a declarative database query language. More precisely, we consider queries of the form  $(\Sigma, q)$ , where  $\Sigma$  is a set of TGDs, and  $q$  a CQ over  $\text{sch}(\Sigma)$ . The *extensional (database) schema* of  $\Sigma$ , denoted  $\text{edb}(\Sigma)$ , is the set of extensional predicates of  $\text{sch}(\Sigma)$ , i.e., the predicates that do not occur in the head of a TGD of  $\Sigma$ . Given a query  $Q = (\Sigma, q)$  and a database  $D$  over  $\text{edb}(\Sigma)$ , the *evaluation* of  $Q$  over  $D$ , denoted  $Q(D)$ , is defined as  $\text{cert}(q, D, \Sigma)$ . We write  $(C, \text{CQ})$  for the query language consisting of all the queries  $(\Sigma, q)$ , where  $\Sigma \in C$ , and  $q$  is a CQ. The evaluation problem for such a query language, dubbed  $\text{Eval}(C, \text{CQ})$ , is defined in the usual way. By definition,  $\bar{c} \in Q(D)$  iff  $\bar{c} \in \text{cert}(q, D, \Sigma)$ . Therefore, the complexity of  $\text{Eval}(C, \text{CQ})$  when  $C = \text{WARD} \cap \text{PWL}$  and  $C = \text{WARD}$  is immediately inherited from Theorem 4.2 and Proposition 3.2, respectively:

**THEOREM 6.1.** *The following statements hold:*

- (1)  $\text{Eval}(\text{WARD} \cap \text{PWL}, \text{CQ})$  is PSPACE-complete in combined complexity, and NLOGSPACE-complete in data complexity.
- (2)  $\text{Eval}(\text{WARD}, \text{CQ})$  is EXPTIME-complete in combined complexity, and PTIME-complete in data complexity.

The main goal of this section is to understand the relative expressive power of  $(\text{WARD} \cap \text{PWL}, \text{CQ})$  and  $(\text{WARD}, \text{CQ})$ . To this end, we are going to adopt two different notions of expressive power namely the classical one, which we call combined expressive power since it considers the set of TGDs and the CQ as one composite query, and the program expressive power, which aims at the decoupling of the set of TGDs from the actual CQ. We proceed with the details starting with the combined expressive power.

### 6.1 Combined Expressive Power

Consider a query  $Q = (\Sigma, q)$ , where  $\Sigma$  is a set of TGDs and  $q(\bar{x})$  a CQ over  $\text{sch}(\Sigma)$ . The *expressive power* of  $Q$ , denoted  $\text{ep}(Q)$ , is the set of pairs  $(D, \bar{c})$ , where  $D$  is a database over  $\text{edb}(\Sigma)$ , and  $\bar{c} \in \text{dom}(D)^{|\bar{x}|}$ , such that  $\bar{c} \in Q(D)$ . The *combined expressive power* of a query language  $(C, \text{CQ})$ , where  $C$  is a class of TGDs, is defined as the set

$$\text{cep}(C, \text{CQ}) = \{\text{ep}(Q) \mid Q \in (C, \text{CQ})\}.$$

Given two query languages  $Q_1, Q_2$ , we say that  $Q_2$  is *more expressive* (w.r.t. the combined expressive power) than  $Q_1$ , written  $Q_1 \leq_{\text{cep}} Q_2$ , if  $\text{cep}(Q_1) \subseteq \text{cep}(Q_2)$ . Moreover, we say that  $Q_1$  and  $Q_2$  are *equally expressive* (w.r.t. the combined expressive power), written  $Q_1 =_{\text{cep}} Q_2$ , if  $Q_1 \leq_{\text{cep}} Q_2$  and  $Q_2 \leq_{\text{cep}} Q_1$ .

The next easy lemma states that  $Q_1 =_{\text{cep}} Q_2$  is equivalent to say that every query of  $Q_1$  can be equivalently rewritten as a query of  $Q_2$ , and vice versa. Given two query languages  $Q_1$  and  $Q_2$ , we write  $Q_1 \leq Q_2$  if, for every  $Q = (\Sigma, q) \in Q_1$ , there exists  $Q' = (\Sigma', q') \in Q_2$  such that, for every  $D$  over  $\text{edb}(\Sigma)$ ,  $Q(D) = Q'(D)$ .

**LEMMA 6.2.** *Consider two query languages  $Q_1$  and  $Q_2$ . It holds that  $Q_1 \leq_{\text{cep}} Q_2$  iff  $Q_1 \leq Q_2$ .*

We are now ready to state the main result of this section, which reveals the expressiveness of  $(\text{WARD} \cap \text{PWL}, \text{CQ})$  and  $(\text{WARD}, \text{CQ})$  relative to Datalog. Let us clarify that a Datalog query is essentially a pair  $(\Sigma, q)$ , where  $\Sigma$  is a Datalog program, or a set of full TGDs, i.e., TGDs without existentially quantified variables, that

have only one head atom, and  $q$  a CQ. We write  $\text{FULL}_1$  for the above class of TGDs. In other words, piece-wise linear Datalog, denoted  $\text{PWL-DATALOG}$ , is the language  $(\text{FULL}_1 \cap \text{PWL}, \text{CQ})$ , while Datalog, denoted  $\text{DATALOG}$ , is the language  $(\text{FULL}_1, \text{CQ})$ , and thus we can refer to their combined expressive power.

**THEOREM 6.3.** *The following statements hold:*

- (1)  $\text{PWL-DATALOG} =_{\text{cep}} (\text{WARD} \cap \text{PWL}, \text{CQ})$ .
- (2)  $\text{DATALOG} =_{\text{cep}} (\text{WARD}, \text{CQ})$ .

Let us explain how (1) is shown; the proof for (2) is similar. We need to show that: (a)  $\text{PWL-DATALOG} \leq_{\text{cep}} (\text{WARD} \cap \text{PWL}, \text{CQ})$ , and (b)  $(\text{WARD} \cap \text{PWL}, \text{CQ}) \leq_{\text{cep}} \text{PWL-DATALOG}$ . By definition,  $\text{FULL}_1 \cap \text{PWL} \subseteq \text{WARD} \cap \text{PWL}$ . Thus,  $(\text{FULL}_1 \cap \text{PWL}, \text{CQ}) \leq (\text{WARD} \cap \text{PWL}, \text{CQ})$ , which, together with Lemma 6.2, implies (a). For showing (b), by Lemma 6.2, it suffices to show that:

**LEMMA 6.4.**  $(\text{WARD} \cap \text{PWL}, \text{CQ}) \leq \text{PWL-DATALOG}$ .

The key idea underlying the above lemma is to convert a linear proof tree  $\mathcal{P}$  of a CQ  $q(\bar{x})$  w.r.t. a set  $\Sigma \in \text{WARD} \cap \text{PWL}$  of TGDs into a piece-wise linear Datalog query  $Q = (\Sigma', q'(\bar{x}))$  such that, for every database  $D$  over  $\text{edb}(\Sigma)$ ,  $\mathcal{P}(D) = Q(D)$ . Roughly, each node of  $\mathcal{P}$  together with its children, is converted into a full TGD that is added to  $\Sigma'$ . Assume that the node  $v$  has the children  $u_1, \dots, u_k$  in  $\mathcal{P}$ , where  $v$  is labeled by  $p_0(\bar{x}_0)$  and, for  $i \in [k]$ ,  $u_i$  is labeled by the CQ  $p_i(\bar{x}_i)$  with  $\bar{x} \subseteq \bar{x}_i$ . In this case, we add to  $\Sigma'$  the full TGD

$$C_{[p_1]}(\bar{x}_1), \dots, C_{[p_k]}(\bar{x}_k) \rightarrow C_{[p_0]}(\bar{x}_0),$$

where  $C_{[p_i]}$  is a predicate that corresponds to the CQ  $p_i$ , while  $[p_i]$  refers to a *canonical renaming* of  $p_i$ . The intention underlying such a canonical renaming is the following: if  $p_i$  and  $p_j$  are the same up to variable renaming, then  $[p_i] = [p_j]$ . We also add to  $\Sigma'$  a full TGD

$$R(x_1, \dots, x_n) \rightarrow C_{[p_R]}(x_1, \dots, x_n)$$

for each  $n$ -ary predicate  $R \in \text{edb}(\Sigma)$ , where  $p_R(x_1, \dots, x_n)$  is the atomic query consisting of the atom  $R(x_1, \dots, x_n)$ . Since in  $\mathcal{P}$  we may have several CQs that are the same up to variables renaming, the set  $\Sigma'$  is recursive, but due to the linearity of  $\mathcal{P}$ , the employed recursion is piece-wise linear, i.e.,  $\Sigma' \in \text{FULL}_1 \cap \text{PWL}$ . The CQ  $q'(\bar{x})$  is simply the atomic query  $C_{[q]}(\bar{x})$ . It should not be difficult to see that indeed  $\mathcal{P}(D) = Q(D)$ , for every database  $D$  over  $\text{edb}(D)$ .

Having the above transformation of a linear proof tree into a piece-wise linear Datalog query in place, we can easily rewrite every query  $Q = (\Sigma, q) \in (\text{WARD} \cap \text{PWL}, \text{CQ})$  into an equivalent query that falls in  $\text{PWL-DATALOG}$ . We exhaustively convert each linear proof tree  $\mathcal{P}$  of  $q$  w.r.t.  $\Sigma$  such that  $\text{nwd}(\mathcal{P}) \leq f_{\text{WARD} \cap \text{PWL}}(q, \Sigma)$  into a piece-wise linear Datalog query  $Q_{\mathcal{P}}$ , and then we take the union of all those queries. Since we consider the canonical renaming of the CQs occurring in a proof tree, and since the size of those CQs is bounded by  $f_{\text{WARD} \cap \text{PWL}}(q, \Sigma)$ , we immediately conclude that we need to explore finitely many CQs. Thus, the above iterative procedure will eventually terminate and construct a finite piece-wise linear Datalog query that is equivalent to  $Q$ , as needed.

### 6.2 Program Expressive Power

The *expressive power* of a set  $\Sigma$  of TGDs, denoted  $\text{ep}(\Sigma)$ , is the set of triples  $(D, q(\bar{x}), \bar{c})$ , where  $D$  is a database over  $\text{edb}(\Sigma)$ ,  $q(\bar{x})$  is a CQ over  $\text{sch}(\Sigma)$ , and  $\bar{c} \in \text{dom}(D)^{|\bar{x}|}$ , such that  $\bar{c} \in \text{cert}(q, D, \Sigma)$ . The

program expressive power of a query language  $(C, CQ)$ , where  $C$  is a class of TGDs, is defined as the set

$$\text{pep}(C, CQ) = \{\text{ep}(\Sigma) \mid \Sigma \in C\}.$$

Given two query languages  $Q_1, Q_2$ , we say that  $Q_2$  is *more expressive* (w.r.t. program expressive power) than  $Q_1$ , written  $Q_1 \leq_{\text{pep}} Q_2$ , if  $\text{pep}(Q_1) \subseteq \text{pep}(Q_2)$ . Moreover, we say that  $Q_2$  is *strictly more expressive* (w.r.t. the program expressive power) than  $Q_1$ , written  $Q_1 <_{\text{pep}} Q_2$ , if  $Q_1 \leq_{\text{pep}} Q_2$  and  $Q_2 \not\leq_{\text{pep}} Q_1$ .

Let us now establish a useful lemma, analogous to Lemma 6.2, which reveals the essence of the program expressive power. For brevity, given two classes of TGDs  $C_1$  and  $C_2$ , we write  $C_1 \leq C_2$  if, for every  $\Sigma \in C_1$ , there exists  $\Sigma' \in C_2$  such that, for every  $D$  over  $\text{edb}(\Sigma)$ , and  $CQ \ q$  over  $\text{sch}(\Sigma)$ ,  $Q(D) = Q'(D)$ , where  $Q = (\Sigma, q)$  and  $Q' = (\Sigma', q)$ . The following holds:

LEMMA 6.5. *Consider two query languages  $Q_1 = (C_1, CQ)$  and  $Q_2 = (C_2, CQ)$ . Then,  $Q_1 \leq_{\text{pep}} Q_2$  iff  $C_1 \leq C_2$ .*

We are now ready to study the expressiveness (w.r.t. the program expressive power) of  $(\text{WARD} \cap \text{PWL}, CQ)$  and  $(\text{WARD}, CQ)$  relative to Datalog. In particular, we show that:

THEOREM 6.6. *The following statements hold:*

- (1)  $\text{PWL-DATALOG} <_{\text{pep}} (\text{WARD} \cap \text{PWL}, CQ)$ .
- (2)  $\text{DATALOG} <_{\text{pep}} (\text{WARD}, CQ)$ .

Let us explain how (1) is shown; the proof for (2) is similar. We need to show that: (a)  $\text{PWL-DATALOG} \leq_{\text{pep}} (\text{WARD} \cap \text{PWL}, CQ)$ , and (b)  $(\text{WARD} \cap \text{PWL}, CQ) \not\leq_{\text{pep}} \text{PWL-DATALOG}$ . Since, by definition,  $\text{FULL}_1 \cap \text{PWL} \subseteq \text{WARD} \cap \text{PWL}$ , we immediately get that  $\text{FULL}_1 \cap \text{PWL} \leq \text{WARD} \cap \text{PWL}$ , and thus, by Lemma 6.5, (a) follows. For showing (b), by Lemma 6.5, it suffices to show that:

LEMMA 6.7.  $\text{WARD} \cap \text{PWL} \not\leq \text{FULL}_1 \cap \text{PWL}$ .

By contradiction, assume the opposite. We define the set of TGDs  $\Sigma = \{P(x) \rightarrow \exists y R(x, y)\}$ , the database  $D = \{P(c)\}$ , and the CQs  $q_1 = Q \leftarrow R(x, y)$  and  $q_2 = Q \leftarrow R(x, y), P(y)$ . By hypothesis, there exists  $\Sigma' \in \text{FULL}_1 \cap \text{PWL}$  such that  $Q_1(D) = Q'_1(D)$  and  $Q_2(D) = Q'_2(D)$ , where  $Q_i = (\Sigma, q_i)$  and  $Q'_i = (\Sigma', q_i)$ , for  $i \in \{1, 2\}$ . Clearly,  $Q_1(D) \neq \emptyset$  and  $Q_2(D) = \emptyset$ , which implies that  $Q'_1(D) \neq \emptyset$  and  $Q'_2(D) = \emptyset$ . However, it is easy to see that  $Q'_1(D) \neq \emptyset$  implies  $Q'_2(D) \neq \emptyset$ , which is a contradiction, and the claim follows.

## 7 IMPLEMENTATION AND FUTURE WORK

The Vadalog system is currently optimized for piece-wise linear warded sets of TGDs in three ways: (i) the first one is related to the way that existential quantifiers interact with recursion; (ii) the second one is related to the optimizer, which detects and uses piece-wise linearity for the purpose of join ordering; (iii) the third way is related to the architecture of the system. Here are some promising directions for future research:

- (1) As said in Section 1,  $\text{NLOGSPACE}$  is contained in the class  $\text{NC}_2$  of highly parallelizable problems. This means that reasoning under piece-wise linear warded sets of TGDs is principally parallelizable, unlike warded sets of TGDs. We plan to exploit this for the parallel execution of reasoning tasks in both core settings and in the map-reduce model.

- (2) Reasoning with piece-wise linear warded sets of TGDs is  $\text{LOGSPACE}$ -equivalent to reachability in directed graphs. Reachability in very large graphs has been well-studied and many algorithms and heuristics have been designed that work well in practice; see, e.g., [12, 17, 19]. We are confident that several of these algorithms can be adapted for our purposes.
- (3) Reachability in directed graphs is known to be in the *dynamic* parallel complexity class  $\text{DYN-FO}$  [13, 25]. We plan to analyze whether reasoning under piece-wise linear warded sets of TGDs, or relevant subclasses thereof, can be shown to be in  $\text{DYN-FO}$ .

## REFERENCES

- [1] Foto N. Afrati, Manolis Gergatsoulis, and Francesca Toni. 2003. Linearisability on datalog programs. *Theor. Comput. Sci.* 308, 1-3 (2003), 199–226.
- [2] Marcelo Arenas, Georg Gottlob, and Andreas Pieris. 2014. Expressive languages for querying the semantic web. In *PODS*. 14–26.
- [3] Patricia C. Arocena, Boris Glavic, Radu Ciucanu, and Renée J. Miller. 2015. The iBench Integration Metadata Generator. *PVLDB* 9, 3 (2015), 108–119.
- [4] Catriel Beeri and Moshe Y. Vardi. 1981. The Implication Problem for Data Dependencies. In *ICALP*. 73–85.
- [5] Luigi Bellomarini, Georg Gottlob, Andreas Pieris, and Emanuel Sallinger. 2017. Swift Logic for Big Data and Knowledge Graphs. In *IJCAI*. 2–10.
- [6] Luigi Bellomarini, Emanuel Sallinger, and Georg Gottlob. 2018. The Vadalog System: Datalog-based Reasoning for Knowledge Graphs. *PVLDB* 11, 9 (2018), 975–987.
- [7] Michael Benedikt, George Konstantinidis, Giansalvatore Mecca, Boris Motik, Paolo Papotti, Donatello Santoro, and Efthymia Tsamoura. 2017. Benchmarking the Chase. In *PODS*. 37–52.
- [8] Peter Van Emde Boas. 1997. The Convenience of Tilings. In *Complexity, Logic, and Recursion Theory*. 331–363.
- [9] Andrea Cali, Georg Gottlob, and Michael Kifer. 2013. Taming the Infinite Chase: Query Answering under Expressive Relational Constraints. *J. Artif. Intell. Res.* 48 (2013), 115–174.
- [10] Andrea Cali, Georg Gottlob, Thomas Lukasiewicz, Bruno Marnette, and Andreas Pieris. 2010. Datalog+/-: A Family of Logical Knowledge Representation and Query Languages for New Applications. In *LICS*. 228–242.
- [11] Andrea Cali, Georg Gottlob, and Andreas Pieris. 2012. Towards more expressive ontology languages: The query answering problem. *Artif. Intell.* 193 (2012), 87–128.
- [12] Edith Cohen, Eran Halperin, Haim Kaplan, and Uri Zwick. 2003. Reachability and Distance Queries via 2-Hop Labels. *SIAM J. Comput.* 32, 5 (2003), 1338–1355.
- [13] Samir Datta, Raghav Kulkarni, Anish Mukherjee, Thomas Schwentick, and Thomas Zeume. 2015. Reachability is in DynFO. In *ICALP*. 159–170.
- [14] Ronald Fagin, Phokion G. Kolaitis, Renée J. Miller, and Lucian Popa. 2005. Data exchange: semantics and query answering. *Theor. Comput. Sci.* 336, 1 (2005), 89–124.
- [15] Georg Gottlob, Giorgio Orsi, and Andreas Pieris. 2014. Query Rewriting and Optimization for Ontological Databases. *ACM Trans. Database Syst.* 39, 3 (2014), 25:1–25:46.
- [16] Georg Gottlob and Andreas Pieris. 2015. Beyond SPARQL under OWL 2 QL Entailment Regime: Rules to the Rescue. In *IJCAI*. 2999–3007.
- [17] Ruoming Jin, Yang Xiang, Ning Ruan, and Haixun Wang. 2008. Efficiently answering reachability queries on very large directed graphs. In *SIGMOD*. 595–608.
- [18] David S. Johnson and Anthony C. Klug. 1984. Testing Containment of Conjunctive Queries under Functional and Inclusion Dependencies. *J. Comput. Syst. Sci.* 28, 1 (1984), 167–189.
- [19] Valerie King. 1999. Fully Dynamic Algorithms for Maintaining All-Pairs Shortest Paths and Transitive Closure in Digraphs. In *FOCS*. 81–91.
- [20] Mélanie König, Michel Leclère, Marie-Laure Mugnier, and Michaël Thomazo. 2015. Sound, complete and minimal UCQ-rewriting for existential rules. *Semantic Web* 6, 5 (2015), 451–475.
- [21] Roman Kontchakov, Martin Rezk, Mariano Rodriguez-Muro, Guohui Xiao, and Michael Zakharyashev. 2014. Answering SPARQL Queries over Databases under OWL 2 QL Entailment Regime. In *ISWC*. 552–567.
- [22] David Maier, Alberto O. Mendelzon, and Yehoshua Sagiv. 1979. Testing Implications of Data Dependencies. *ACM Trans. Database Syst.* 4, 4 (1979), 455–469.
- [23] Jeffrey F. Naughton. 1986. Data Independent Recursion in Deductive Databases. In *PODS*. 267–279.
- [24] Jeffrey F. Naughton and Yehoshua Sagiv. 1987. A Decidable Class of Bounded Recursions. In *PODS*. 227–236.
- [25] Sushant Patnaik and Neil Immerman. 1997. Dyn-FO: A Parallel, Dynamic Complexity Class. *J. Comput. Syst. Sci.* 55, 2 (1997), 199–209.