

# *The Impact of Disjunction on Reasoning under Existential Rules: Research Summary*

MICHAEL MORAK

*University of Oxford, Department of Computer Science, OX1 3QD, United Kingdom*  
(e-mail: michael.morak@cs.ox.ac.uk)

*submitted 1 January 2003; revised 1 January 2003; accepted 1 January 2003*

---

## Abstract

Datalog<sup>±</sup> is a Datalog-based language family enhanced with existential quantification in rule heads, equalities and negative constraints. Query answering over databases with respect to a Datalog<sup>±</sup> theory is generally undecidable, however several syntactic restrictions have been proposed to remedy this fact. However, a useful and natural feature however is as of yet missing from Datalog<sup>±</sup>: The ability to express uncertain knowledge, or choices, using disjunction. It is the precise objective of the doctoral thesis herein discussed, to investigate the impact on the complexity of query answering, of adding disjunction to well-known decidable Datalog<sup>±</sup> fragments, namely guarded, sticky and weakly-acyclic Datalog<sup>±</sup> theories. For guarded theories with disjunction, we obtain a strong 2EXP lower bound in the combined complexity, even for very restricted formalisms like fixed sets of (disjunctive) inclusion dependencies. For sticky theories, the query answering problem becomes undecidable, even in the data complexity, and for weakly-acyclic query answering we see a reasonable and expected increase in complexity.

*A full version of a paper accepted to be presented at the Doctoral Consortium of the 30th International Conference on Logic Programming (ICLP 2014), July 19-22, Vienna, Austria*

**KEYWORDS:** Ontological Reasoning, Query Answering, Existential Rules, Logic, TGDs

---

## 1 Introduction and Problem Description

For the last thirty years, Datalog (see e.g., (Abiteboul et al. 1995)) has played an important role as a conceptual query language. Whilst not directly implemented in mainstream database management systems (DBMS), it did heavily influence the design of the SQL standard, which now also allows for recursive statements, as can be expressed in Datalog.

However in recent years it has become increasingly important to add ontological reasoning capabilities to the existing object-relational querying capabilities of traditional DBMS: A query is no longer just evaluated over the extensional relational database, but also over an ontological theory that, using rules and constraints, describes how to derive new (intensional) knowledge from the extensional data. By extending Datalog in such a way that existential quantification, the first-order logic constant *false* and equalities between variables are permitted in the rule heads, this behaviour can be expressed. Recently the Datalog<sup>±</sup> family of languages has been proposed in (Cali et al. 2011), that defines sensible restrictions on the structure of such an ontological theory. These restrictions are necessary as, depending on the structure of the ontological theory, an infinite amount of intensional knowledge might be derivable, rising the question of decidability of this type of reasoning. Also, as new values can be invented along the way, the domain can become infinite.

Despite these obstacles, commercial service providers have already started to integrate ontological reasoning engines into their database management systems (see e.g., (Oracle Inc. 2011; Microsoft Corp. 2011)), as there are several applications where such capabilities are desirable, such as data exchange, ontological reasoning (e.g., reasoning under description logics, or in the semantic web) and web data extraction.

*Problem Statement.* Given the fact that ontological reasoning is gaining mainstream acceptance and the fact that, as for example Answer Set Programming has proven, rule-based languages are well suited for knowledge representation and reasoning tasks, it is natural to ask how to enrich the languages we currently have with new, useful constructs. The construct that we want to focus on here is disjunction. Until now, Datalog<sup>±</sup> rules only allow us to express deterministic knowledge. But what about natural statements like “every person has a parent that is either male or female” or “every student is either an undergraduate or a graduate student”? Such statements are not captured by existing Datalog<sup>±</sup> languages. Seeing that disjunctive knowledge is an important feature in other logical languages like Answer Set Programming or Description Logics that allows users to intuitively formulate problems by, e.g., applying a guess-and-check approach, enriching Datalog<sup>±</sup> with disjunction is therefore a logical next step.

The objective of my doctoral studies is thus to introduce the language feature of disjunction to Datalog<sup>±</sup>, and investigate in-depth what the impact of doing so is w.r.t. decidability and complexity of reasoning, focussing on conjunctive query answering in particular.

## 2 Background and Literature Review

In the following subsections, we give a few basic preliminaries describing Datalog<sup>±</sup>, as well as an overview over the known results in the area.

### 2.1 Background

In this section the basic notions of conjunctive query evaluation under tuple generating dependencies (TGDs) are recalled, including a review of the chase procedure, an important algorithmic tool in the evaluation of queries under TGDs. Furthermore we briefly introduce the concept of stable models in the logic programming perspective. We assume that the reader is familiar with first-order logic as well as basic complexity theory. Good introductions to the former can be found in e.g. (Barwise 1977) and (Andrews 2002), for the latter we recommend (Papadimitriou 1994).

#### 2.1.1 Conjunctive Queries and the Relational Model

In order to define the semantics of conjunctive queries, we first need to introduce the relational data model. In the relational data model, the structure or *schema*  $\mathcal{S}$  of a database and its contents or *instance*  $D$  are distinct objects.

A schema  $\mathcal{S}$  consists of a finite number of *relation symbols* (also called *predicates*)  $r_i$ , that is,  $\mathcal{S} = \{r_1, \dots, r_n\}$ .

Such a relation symbol  $r_i \in \mathcal{S}$  (for any  $i$ ) consists of a finite number of *attributes*, such that each attribute has a *domain* of possible values. We consider here only the case that all predicates have a common domain  $\Gamma \cup \Gamma_N$ , where  $\Gamma$  is a set of constants and  $\Gamma_N$  is a set of labelled nulls (i.e.,

distinct null values, each with a unique name, comparable to skolem constants). The number of attributes of a relation symbol is called the *arity*, denoted  $\text{arity}(r_i)$ .

A relation  $R_i$  for predicate  $r_i$  is a set of *tuples* and each tuple is a mapping of each attribute in  $r_i$  to  $\Gamma \cup \Gamma_N$ . Such a tuple of  $R_i$  is denoted by  $r_i(x_1, \dots, x_k)$  (also referred to as an *atom*), where  $k = \text{arity}(r_i)$ .

An *instance*  $I$  for a schema  $\mathcal{S}$  consists of *relations*  $R_i$  for each  $r_i \in \mathcal{S}$ , that is,  $D = \{R_1, \dots, R_n\}$ . An instance in which no null values from  $\Gamma_N$  appear is referred to as a *database*, usually denoted  $D$ . Note that, when viewed as a first-order theory, we may simply interpret an instance as a conjunction of atoms.

A *conjunctive query*  $q$  over a database schema  $\mathcal{S}$  is an assertion of the form

$$q(\vec{X}) \leftarrow \exists \vec{Y} \varphi(\vec{X}, \vec{Y})$$

where  $\vec{X}$  and  $\vec{Y}$  are vectors of (first-order logic) variables,  $q(\vec{X})$  is called the *head*,  $\text{dimension}(\vec{X})$  is called the *arity* of  $q$  and  $\varphi(\vec{X}, \vec{Y})$  is called the *body*, where  $\varphi(\vec{X}, \vec{Y})$  is a first-order formula consisting of a conjunction of atoms of the form  $r_i(t_1, \dots, t_k)$  and equalities of the form  $t_1 = t_2$ , where  $r_i$  is a predicate of  $\mathcal{S}$  with arity  $k$  and each  $t_i$  is either a constant from  $\Gamma$  or a (first-order logic) variable. If the arity is 0 then  $q$  is called a *boolean conjunctive query*.

With every database  $D = \{R_1, \dots, R_n\}$  over a schema  $\mathcal{S}$ , we can now associate a finite first-order structure  $M_D = (U, R_1, \dots, R_n)$  with universe  $U = \Gamma$ . The evaluation of a conjunctive query  $q$  then comes down to checking satisfiability in first-order logic as follows:  $q$  has an answer over  $D$ , denoted  $D \models q$ , if and only if the set  $\{ \langle a_1, \dots, a_k \rangle \mid M_D \models q(a_1, \dots, a_k) \}$  is non-empty, with  $a_i \in \Gamma$ . This set is also called the set of *answers* to  $q$  over  $D$ , where  $k$  is the arity of  $q$ .

### 2.1.2 Dependencies

For reasoning tasks over databases, the need arises to express how new (*intensional*) knowledge can be derived from the data that is stored in the database (called the *extensional data*). An established way to do this is to introduce a set  $\Sigma$  of rules that describe the relation between intensional and extensional data. In this case for a database  $D$ , the logical theory  $D \cup \Sigma$ , i.e., the conjunction of the facts in the database with all the rules in  $\Sigma$ , is taken as a basis for conjunctive query evaluation.

Rules in  $\Sigma$  over a schema  $\mathcal{S}$  are of either one of the following two forms:

$$\forall \vec{X} (\varphi(\vec{X}) \rightarrow \exists \vec{Y} \psi(\vec{X}, \vec{Y})) \quad (1)$$

$$\forall \vec{X} (\varphi(\vec{X}) \rightarrow X_i = X_j) \quad (2)$$

where rules of the form of (1) are referred to as tuple generating dependencies (TGDs) and of (2) as equality generating dependencies (EGDs), with  $\varphi$  and  $\psi$  being conjunctions of predicates from  $\mathcal{S}$  (also called *atoms*) and  $X_i$  and  $X_j$  are the  $i$ -th and  $j$ -th position in vector  $\vec{X}$ .  $\varphi$  is also referred to as the body of the dependency and  $\psi$  or  $X_i = X_j$  as the head. TGDs where  $\psi = \perp$  are called *negative constraints*. For brevity, we will omit the universal quantifiers in front of TGDs and EGDs, and replace conjunctions in the body by commas.

Given an instance  $I$ , it is said to be satisfying a dependency  $\sigma \in \Sigma$ , that is,  $I \models \sigma$ , if the first-order sentence formed by a conjunction of the facts in  $I$  and  $\sigma$  is satisfiable. By extension,  $I$  satisfies  $\Sigma$  ( $I \models \Sigma$ ) iff it satisfies every  $\sigma \in \Sigma$ .

The models of a database  $D$  over a schema  $\mathcal{S}$  with respect to  $\Sigma$ , denoted  $\text{Mod}(D, \Sigma)$ , are all instances  $M$  that satisfy  $D \cup \Sigma$  (i.e.,  $I \supseteq D$  and  $I \models \Sigma$ ). When answering conjunctive queries we

use the certain answer semantics, i.e., we consider the query to be true only iff it is true under every model. The set of answers for a conjunctive query  $q$ , denoted  $ans(q, D, \Sigma)$ , thus equals the set

$$\{ \langle a_1, \dots, a_k \rangle \mid \forall M \in Mod(D, \Sigma) : M \models q(a_1, \dots, a_k) \}$$

For complexity analysis we focus on the decision version of this problem. This is the central problem when analyzing the complexity of databases, tuple and equality generating dependencies and therefore Datalog<sup>±</sup> complexity issues. Below it is formulated for boolean conjunctive queries, which we will focus on in this work:

#### BCQ-ANSWERING

**Instance:**  $\langle q, D, \Sigma \rangle$ :  $q$  a boolean conjunctive query,  $D$  a database and  $\Sigma$  a set of dependencies

**Question:**  $D \cup \Sigma \models q$ ?

Usually when dealing with query evaluation over databases the *data complexity* and the *combined complexity* are of interest. In this paper we follow the approach of (Vardi 1982) where for the former everything except the database  $D$  is considered fixed, i.e., the only input is the database. For the latter, the database  $D$ ,  $\Sigma$  and the query itself form the input.

Unfortunately, in general it holds that BCQ-ANSWERING is undecidable under unrestricted sets of TGDs, as has been shown in (Beeri and Vardi 1981). In (Cali et al. 2013; Baget et al. 2009; Baget et al. 2011) it has further been shown that even singleton sets of TGDs cause query answering to become undecidable.

These results clearly show that restrictions must be placed on the structure of  $\Sigma$  to ensure decidability. This is a non-trivial problem, as simple restrictions, like limiting the number of TGDs, are not enough.

### 2.1.3 The Chase

One of the fundamental tools to algorithmically check implication of dependencies is the *chase procedure*, introduced in (Maier et al. 1979), which was later adapted for checking query containment in (Johnson and Klug 1984), in the setting of databases with tuple and equality generating dependencies, or, more specifically, in the setting of databases with inclusion and functional dependencies. The chase algorithm tries to extend a given database instance in such a way that every TGD and EGD becomes satisfied. This is done by exhaustively (i.e., until a fix-point is reached) applying the *chase step*:

#### Definition 2.1

Let  $D$  be a database and  $\Sigma$  be a set of dependencies. A *chase step* is defined as follows:

**TGDs.** Let  $\Sigma$  contain a TGD  $\varphi(\vec{X}) \rightarrow \exists \vec{Y} \psi(\vec{X}, \vec{Y})$ , such that

- $D \models \varphi(\vec{a})$  for some assignment  $\vec{a}$  to  $\vec{X}$ , and
- $D \not\models \exists \vec{Y} \psi(\vec{a}, \vec{Y})$ .

Then extend  $D$  with facts  $\psi(\vec{a}, \vec{y})$ , where the elements of  $\vec{y}$  are fresh labelled nulls (i.e., values from  $\Gamma_N$  that have not been in use in  $D$  up to that point).

**EGDs.** Let  $\Sigma$  contain an EGD  $\varphi(\vec{X}) \rightarrow X_i = X_j$ , such that

- $D \models \varphi(\vec{a})$  for some assignment  $\vec{a}$  to  $\vec{X}$ , and
- $a_i \neq a_j$

If  $a_i$  is a labelled null, then replace every occurrence of  $a_i$  with  $a_j$  or vice-versa if  $a_j$  is a labelled null. If  $a_i$  and  $a_j$  are distinct constants, end the chase with failure.

### Definition 2.2

The *chase expansion* of a database instance  $D$  with respect to a set of dependencies  $\Sigma$  is a sequence  $D_0, D_1, \dots, D_m$ , such that  $D_0 = D$  and for  $i \geq 0$ ,  $D_{i+1}$  is obtained from  $D_i$  by applying a chase step. After exhaustively applying such chase steps, we obtain  $D_m$ , also denoted  $\text{chase}(D, \Sigma)$ .

The chase can have three different outcomes: Failure, non-terminating success or terminating success. In case of success the resulting instance  $D_m$  satisfies all dependencies in  $\Sigma$ . Note that if the chase does not terminate,  $m = \infty$  and the size of  $D_m$  is infinite.

We assume that the chase is fair, i.e., we exclude the possibility of a degenerated chase expansion by assuming that the chase expansion is constructed level by level, and after each application of a TGD, all applicable EGDs are applied. This ensures that every TGD that can be applied, is applied, and therefore we exclude the case that only a single infinite path in the chase expansion is ever expanded when in case the chase is infinite.

*Query Answering and the Chase* In case the chase succeeds, it computes a *universal solution* for  $\langle D, \Sigma \rangle$ . Every model  $M \in \text{Mod}(D, \Sigma)$  can then be obtained by appropriate instantiation of labelled nulls in  $\text{chase}(D, \Sigma)$  (i.e., for every model  $M$ , there exists a homomorphism mapping the universal solution to  $M$ ; cf. (Deutsch et al. 2008)). Using this property, the chase expansion of a database  $D$ , with respect to a set of dependencies  $\Sigma$ , can be used for answering conjunctive queries, as the following theorem shows:

### Theorem 2.3 ((Deutsch et al. 2008))

Given a boolean conjunctive query  $q$  over a schema  $\mathcal{S}$ , a database  $D$  of  $\mathcal{S}$  and a set of dependencies  $\Sigma$  over  $\mathcal{S}$ , then in cases where the chase does not fail, it holds that  $D \cup \Sigma \models q$  if and only if  $\text{chase}(D, \Sigma) \models q$ .

In case the chase fails, query answering is trivial: As there is no model, every boolean conjunctive query clearly is entailed by  $D \cup \Sigma$  (cf. the definition of certain answers in section 2.1.2).

## 2.2 Literature Review

In this section we discuss the different kinds of restrictions known to ensure decidability of query answering under sets of TGDs. The decidable classes of TGDs discussed below are defined by syntactic properties that either apply to single TGDs (local syntactic conditions) or to the set of all TGDs (global syntactic conditions). These properties can be checked in finite time using appropriate algorithms. Each subsection deals with a known syntactic condition that ensures decidability of query answering.

*Inclusion Dependencies* Inclusion dependencies (IDs), one of the simplest forms of dependencies, allow one to express that certain values occurring in a specific position in one relation, must also occur at (or be included in) a specific position in another relation. This allows for TGDs

that consist of one body and head atom only, and no variable may occur twice in the head or the body. The following is an example of an inclusion dependency, expressing that every student is a person:

$$student(X, Y) \rightarrow \exists Z person(X, Z)$$

The query answering problem was shown to be decidable, and in fact in  $AC_0$  (resp. PSPACE) in the data (resp. combined) complexity.

*Linear Tuple Generating Dependencies* This class is similar to IDs in that it allows for TGDs with only a single body atom, but generalizes them, because it allows repetition of variables in the body or head (e.g., the TGD  $r(X, Y, X) \rightarrow s(X, Y)$  is a linear TGD but not an ID).

Sets of linear TGDs enjoy the so-called *bounded derivation-depth property (BDDP)*, which roughly implies that only a finite initial part of the chase is required for query answering, thus ensuring decidability. As with inclusion dependencies, first-order rewritability (i.e., rewriting  $q$  and  $\Sigma$  into a first-order query  $q_\Sigma$ , such that  $D \cup \Sigma \models q$  iff  $D \models q_\Sigma$ ) is thus possible (cf. (Calì et al. 2009; Calì et al. 2010)). Therefore we get decidability, and query answering is in  $AC_0$  in the data complexity. Regarding combined complexity, results from inclusion dependencies carry over to linear TGDs, resulting in the PSPACE-completeness for query answering in the general case and NP-completeness in case of a fixed set of TGDs.

*Guarded Tuple Generating Dependencies* In (Calì et al. 2013), linear TGDs are extended to so-called *guarded* TGDs, that have a body atom that contains all variables occurring in the body, i.e., all universally quantified variables. This atom is called the *guard*. If there are multiple such atoms, the leftmost is taken as the guard. An example of a guarded TGD that says that if students are in their first semester, they have a tutor, is as follows. Note that it is not linear as it has multiple atoms in the body.

$$student(X, Y), firstsemester(X) \rightarrow \exists Z tutor(X, Z)$$

Linear TGDs and inclusion dependencies are trivially guarded, as they only have exactly one body atom. However, guarded TGDs are not first-order rewritable. This is shown by creating a database, query and a set of guarded TGDs in such a way that answering the query requires the computation of the transitive closure over a relation in the database. It is well known that this property cannot be expressed in a finite first-order query, and we cannot obtain decidability thusly. However, it can be shown that the universal model constructed by the chase, albeit possibly infinite, is of finite treewidth (i.e., it is tree-like and cannot be arbitrarily cyclic). From Courcelle's famous Theorem (cf. (Courcelle 1990)), which states that evaluating first-order sentences over structures of finite treewidth is decidable, we derive decidability for query answering under sets of guarded TGDs.

The complexity of query answering under guarded TGDs was investigated in (Calì et al. 2009), where it was established that, whenever a query is actually entailed by a database and a set of guarded TGDs, then all atoms needed to answer the query are derived in a finite, initial portion of the chase when restricted only to guards and atoms derived from them, whereby the size of this portion depends only on the query and the set of TGDs. Therefore, constructing this part of the chase and evaluating a boolean conjunctive query over it is enough to compute the answer. It is shown that this can be done in polynomial time in the data complexity, whereby P-membership follows. Hardness for P was shown in (Dantsin et al. 2001) by reduction to the implication problem for propositional logic programs.

The combined complexity is investigated in (Cali et al. 2013), proving the 2EXP-completeness for the general case and EXP-completeness in case of fixed arity. Also membership in NP was shown in case where the set of TGDs is fixed. NP-hardness follows from results in (Chandra and Merlin 1977), which show that NP-hardness holds even for the empty set of TGDs.

*Weakly-Guarded Sets of TGDs* In (Cali et al. 2013), guarded TGDs were extended to *weakly-guarded sets of TGDs*. Every TGD in such sets must have an atom in its body that contains all the variables where a null value may appear during the chase. The leftmost such atom is called the *weak-guard*. This class is the first class discussed here that is based on a global property. It is easy to see that, as guarded TGDs contain a body atom with all universally quantified variables, they are trivially weakly-guarded, as the guard is also a weak-guard.

It is implicit in (Cali et al. 2013) that it can be verified in polynomial time whether a set of TGDs is weakly-guarded or not: For a schema  $\mathcal{S}$  we first need to compute all the positions for each predicate where a null value can occur during the chase with respect to a set of TGDs  $\Sigma$ . These positions are called *affected* and computing them has been shown to be possible in polynomial time. Then we have to check for each TGD in  $\Sigma$  whether it contains a weak-guard, which, knowing the affected positions, is also possible in polynomial time.

It is then shown that weakly-guarded sets of TGDs enjoy the same favorable property as guarded TGDs, namely, the chase has finite treewidth. Given this fact, decidability of query answering is established as before. Regarding the complexity, in general the problem is 2EXP-complete, EXP-complete if the arity is fixed or the set of TGDs fixed, and it remains EXP-complete even if only the database is considered as input (data complexity).

*Weakly-Acyclic Sets of Tuple Generating Dependencies* The notion of *Weak Acyclicity* was established in the landmark paper (Fagin et al. 2005) as a syntactic condition to guarantee termination of the chase procedure. For this we first need to define the notion of a dependency graph.

A dependency graph  $G = (V, E)$  is constructed as follows:  $V$  is the set of attributes of all the relations occurring in  $\Sigma$ . We will denote the  $i$ th attribute of some relation  $r$  by  $r[i]$ . For each TGD  $\sigma = \varphi(\vec{X}) \rightarrow \exists \vec{Y} \psi(\vec{X}, \vec{Y})$  and each variable  $X \in \vec{X}$  shared between the relation attributes  $r[i]$  in  $\varphi$  and  $s[j]$  in  $\psi$ , we add an edge  $(r[i], s[j])$  to  $E$ . We add a special edge  $(r[i], p[k])$  to  $E$  for each attribute  $p[k]$  in  $\psi$  occupied by a variable  $Y \in \vec{Y}$ , and each attribute  $r[i]$  occurring in the body of  $\sigma$ .

A set  $\Sigma$  of TGDs is called *weakly-acyclic* if its dependency graph contains no cycles through special edges. The definition of weak acyclicity is a global property and can be decided in P, as the construction of the dependency graph and the cycle-check through a special edge are both feasible in P.

In (Fagin et al. 2005) it was shown that for weakly-acyclic sets of TGDs the chase always terminates. This is ensured by the fact that when cycles through special edges in the dependency graph are forbidden, no new null values can be added in a later chase step because of a null value added in an earlier chase step. Therefore we trivially get decidability: Simply compute the (finite) chase, and then answer the query on the obtained finite model.

Regarding complexity (cf. (Cali et al. 2013; Kolaitis et al. 2006)), in general the problem of BCQ-ANSWERING is 2EXP-complete for weakly-acyclic sets of TGDs. When the set of TGDs is fixed, the BCQ-ANSWERING problem is known to be NP-complete. P-completeness holds for the data complexity, following from the complexity of the fact inference problem for fixed Datalog programs (see (Dantsin et al. 2001)).

*Sticky Sets of Tuple Generating Dependencies* A recent addition to the set of syntactic conditions that ensure decidability and favourable complexity of conjunctive query evaluation is the paradigm of *stickiness*, introduced in (Cali et al. 2012). A survey of sticky classes can be found in (Cali et al. 2010). The class of *sticky sets of TGDs* is defined as follows: In a first step, a variable marking of all TGDs in a set  $\Sigma$  is computed by a procedure called SMarking. This is a two-step procedure:

1. *Initial marking*: For each  $\sigma \in \Sigma$ , if there exists a variable  $V$  in the body of  $\sigma$  and an atom without this variable exists in the head of  $\sigma$ , mark each occurrence of  $V$  in the body.
2. *Propagation step*: Until a fixpoint is reached, consider any pair  $\langle \sigma_1, \sigma_2 \rangle \in \Sigma \times \Sigma$ . If a universally quantified variable  $V$  occurs in  $head(\sigma_1)$  at positions  $\pi_1, \dots, \pi_m$  for  $m \geq 1$  and an atom in  $body(\sigma_2)$  exists where at each of these same positions a marked variable occurs, then mark each occurrence of  $V$  in  $body(\sigma_1)$ .

*Definition 2.4* ((Cali et al. 2012))

A set  $\Sigma$  of TGDs is called *sticky* if and only if there is no TGD in SMarking( $\Sigma$ ) such that a marked variable occurs in its body more than once.

The property of stickiness is incomparable to guardedness and weak acyclicity but strictly generalizes inclusion dependencies. In comparison to other discussed syntactic classes of TGDs, sticky sets of TGDs allow for a mildly restricted way to express joins. The following is an example of a sticky (singleton) set of TGDs, expressing the join between two tables, department and employee, to get a combined table of departments and their heads:

$$department(X, Y), employee(Y, Z) \rightarrow headofdept(X, Y, Z)$$

Note that the above TGD is not weakly-guarded.

### 3 Goal and Current Status of the Research

The goal, as already discussed in the introduction, is to introduce disjunction into Datalog<sup>±</sup> and investigate the impact of doing so on the decidability and complexity of query answering. We thus extend the definition of a TGD to allow for disjunction as follows:

A *disjunctive tuple-generating dependency (DTGD)*  $\sigma$  is a first-order formula  $\forall \vec{X} \varphi(\vec{X}) \rightarrow \bigvee_{i=1}^n \exists \vec{Y} \psi_i(\vec{X}, \vec{Y})$ , where  $n \geq 1$ ,  $\vec{X} \cup \vec{Y} \subset \Gamma_V$ , and  $\varphi, \psi_1, \dots, \psi_n$  are conjunctions of atoms;  $\varphi$  is the *body* of  $\sigma$ , denoted  $body(\sigma)$ , while  $\bigvee_{i=1}^n \psi_i$  is the *head* of  $\sigma$ , denoted  $head(\sigma)$ . If  $n = 1$ , then  $\sigma$  is a *tuple-generating dependency (TGD)*. Given a set  $\Sigma$  of DTGDs,  $schema(\Sigma)$  is the set of predicates occurring in  $\Sigma$ .

We employ the *disjunctive chase* introduced in (Deutsch and Tannen 2003) in order to answer queries. It is an extension of the chase procedure described in Section 2.1. Consider an instance  $I$ , and a DTGD  $\sigma = \varphi(\vec{X}) \rightarrow \bigvee_{i=1}^n \exists \vec{Y} \psi_i(\vec{X}, \vec{Y})$ . We say that  $\sigma$  is *applicable* to  $I$  if there exists a homomorphism  $h$  (i.e., a substitution of labelled nulls to either constants or other labelled nulls) such that  $h(\varphi(\vec{X})) \subseteq I$ , but there is no  $i \in \{1, \dots, n\}$  and a homomorphism  $h' \supseteq h$  such that  $h'(\psi_i(\vec{X}, \vec{Y})) \subseteq I$ . The result of applying  $\sigma$  to  $I$  with  $h$  is the set  $\{I_1, \dots, I_n\}$ , where  $I_i = I \cup h'(\psi_i(\vec{X}, \vec{Y}))$ , for each  $i \in \{1, \dots, n\}$ , and  $h' \supseteq h$  is such that  $h'(Y)$  is a “fresh” labelled null not occurring in  $I$ , for each  $Y \in \vec{Y}$ . For such an application of a DTGD, which defines a single DTGD *chase step*, we write  $I \langle \sigma, h \rangle \{I_1, \dots, I_n\}$ .

A *disjunctive chase tree* of a database  $D$  and a set  $\Sigma$  of DTGDs is a (possibly infinite) tree such

that the root is  $D$ , and for every node  $I$ , assuming that  $\{I_1, \dots, I_n\}$  are the children of  $I$ , there exists  $\sigma \in \Sigma$  and a homomorphism  $h$  such that  $I(\sigma, h)\{I_1, \dots, I_n\}$ . The disjunctive chase algorithm for  $D$  and  $\Sigma$  consists of an exhaustive application of DTGD chase steps in a fair fashion, which leads to a disjunctive chase tree  $T$  of  $D$  and  $\Sigma$ ; we denote by  $dchase(D, \Sigma)$  the set  $\{I \mid I \text{ is a leaf of } T\}$ . Note that each leaf of  $T$  is well-defined as the least fixpoint of a monotonic operator. By construction, each instance of  $dchase(D, \Sigma)$  is a model of  $D$  and  $\Sigma$ . Interestingly,  $dchase(D, \Sigma)$  is a *universal set model* of  $D$  and  $\Sigma$ , i.e., for each  $M \in Mod(D, \Sigma)$ , there exists  $I \in dchase(D, \Sigma)$  and a homomorphism  $h_I$  such that  $h_I(I) \subseteq M$  (Deutsch et al. 2008). This implies that w.r.t. certain answers, given a query  $q$ ,  $D \cup \Sigma \models q$  iff  $I \models q$ , for each  $I \in dchase(D, \Sigma)$ .

*Current Status.* Currently we have investigated and obtained results for all the decidable classes of TGDs. For the guarded-based classes, adding disjunction does not make the problem of query answering undecidable. However it does in certain cases increase the complexity of the problem by a significant amount. For the guarded-based classes of TGDs (i.e., IDs, linear, guarded and weakly-guarded), we have established all relevant complexity results when extending them to DTGDs.

In case of sticky TGDs, when adding disjunction the problem of query answering becomes undecidable. This was a very surprising result, given the fact that the complexity of query answering under sticky sets of TGDs is lower than under guarded TGDs.

In case of weakly-acyclic TGDs, data complexity results have been obtained, as well as certain lower bounds in the combined complexity, however a matching upper bound is still missing here. Decidability is assured in any case, because the disjunctive chase terminates, which follows from the definition of weak acyclicity.

#### 4 Preliminary Results

One classical work on disjunction in ontologies is (Calvanese et al. 2006), which immediately gives us coNP-hardness for conjunctive query answering over disjunctive ontologies, even if the query is fixed, and the ontology consists of a fixed, single rule of the form  $a(X) \rightarrow b(X) \vee c(X)$ . Without restricting the query language, there is thus no hope to get tractability results. However, for atomic queries, where the query consists only of a single atom, there are tractable data complexity cases to be found.

*Arbitrary queries.* In (Bourhis et al. 2013), we have investigated the complexity picture for answering arbitrary queries. The main results are as follows:

- 2EXP-completeness whenever the query is non fixed. This is shown by simulating a Büchi tree automaton, and it even holds for fixed sets of Disjunctive Inclusion Dependencies (DIDs) of arity at most three, or of non-fixed sets of the same with arity at most two.
- coNP-completeness in the data complexity for query answering under DIDs up to guarded DTGDs.
- EXP-completeness in the data complexity for query answering under weakly-guarded sets of DTGDs.

In case of (non-disjunctive, classical) TGDs, complexity results coincide in the data complexity, but vary from coNP-completeness to 2EXP-completeness for fixed sets of IDs to weakly-guarded sets of TGDs. It is thus interesting to note that adding disjunction to expressive languages

doesn't change the complexity in this case, but there is a high cost to add it to less expressive languages.

*Atomic queries.* In (Gottlob et al. 2012), we have investigated the complexity of answering single-atom queries. Here the complexity results vary considerably:

- 2EXP-completeness in the combined complexity for guarded DTGDs.
- EXP-completeness in the combined complexity for linear DTGDs.
- coNP-completeness in the data complexity for guarded DTGDs.
- Membership in  $AC_0$  in the data complexity for linear DTGDs.

In the case of atomic queries we do have a number of tractability results to offer, especially the highly parallelizable data complexity of  $AC_0$  in case of atomic query answering over sets of linear DTGDs (which captures the class of DIDs). For guarded DTGDs, most of the results follow directly from expressive fragments of first-order logic (the Guarded Fragment (Bárány et al. 2010; Grädel 1999), and Guarded-Negation First-Order Logic (Bárány et al. 2011; Bárány et al. 2012)). For linear, we develop novel machinery to obtain our respective bounds.

## 5 Open Issues and Expected Achievements

In addition to the published results, we would like to find answers to the following questions: What is the complexity of query answering under sets of

1. guarded-based DTGDs in case where the query is acyclic or of bounded (hyper)treewidth?
2. weakly-acyclic sets of DTGDs?
3. sticky sets of DTGDs?

Regarding the first item, we have already managed to obtain all the relevant results. In fact, for bounded (hyper)treewidth, the complexity table coincides with that of arbitrary queries. For acyclic queries, there are drops in complexity corresponding to the expressivity of the language considered. Papers containing these results have been submitted to this year's MFCS conference and DL workshop. It is our plan to subsequently publish these results, in addition to some extended work on arbitrary and atomic queries in a comprehensive journal paper, treating all the guarded-based classes of DTGDs, in the course of 2014.

Regarding weakly-acyclic, we already have answers to the complexity questions for data complexity and the cases of fixed sets and fixed arities. However, we are still missing the combined complexity results. Before submission of my thesis, we plan to close these open complexity questions as well.

Lastly, for sticky DTGDs, we have an undecidability proof, which is somewhat surprising as query answering under sticky TGDs is easier in terms of complexity than it is for guarded TGDs, yet the addition of disjunction doesn't cause a complexity increase in the latter. We have therefore focused on extending guarded DTGDs with cross-products (a form of join allowed in sticky TGDs). This again yields undecidability, however it becomes decidable if restricted to arity at most two, where binary predicates can never participate in a disjunction. For this case we are working on obtaining the relevant complexity results.

## References

- ABITEBOUL, S., HULL, R., AND VIANU, V. 1995. *Foundations of Databases*. Addison-Wesley.
- ANDREWS, P. B. 2002. *An Introduction to Mathematical Logic and Type Theory: To Truth Through Proof*. Academic Press Professional, Inc., San Diego, CA, USA.
- BAGET, J.-F., LECLÈRE, M., AND MUGNIER, M.-L. 2009. Walking the decidability line for rules with existential variables (long version). Tech. Rep. LIRMM RR-09030, Laboratoire d'Informatique et de Robotique et de Microélectronique de Montpellier, Université Montpellier, CNRS.
- BAGET, J.-F., LECLÈRE, M., MUGNIER, M.-L., AND SALVAT, E. 2011. On rules with existential variables: Walking the decidability line. *Artif. Intell.* 175, 9-10, 1620–1654.
- BÁRÁNY, V., GOTTLÖB, G., AND OTTO, M. 2010. Querying the guarded fragment. In *Proc. of LICS*. 1–10.
- BÁRÁNY, V., TEN CATE, B., AND OTTO, M. 2012. Queries with guarded negation. *PVLDB* 5, 11, 1328–1339.
- BÁRÁNY, V., TEN CATE, B., AND SEGOUFIN, L. 2011. Guarded negation. In *Proc. of ICALP*. 356–367.
- BARWISE, J. 1977. An introduction to first-order logic. In *Handbook of Mathematical Logic*, J. Barwise, Ed. North-Holland, Amsterdam, 5–46.
- BEERI, C. AND VARDI, M. Y. 1981. The implication problem for data dependencies. In *Proc. ICALP*. Lecture Notes in Computer Science, vol. 115. Springer, 73–85.
- BOURHIS, P., MORAK, M., AND PIERIS, A. 2013. The impact of disjunction on query answering under guarded-based existential rules. In *IJCAI*, F. Rossi, Ed. IJCAI/AAAI.
- CALÌ, A., GOTTLÖB, G., AND KIFER, M. 2013. Taming the infinite chase: Query answering under expressive relational constraints. *J. Artif. Intell. Res. (JAIR)* 48, 115–174.
- CALÌ, A., GOTTLÖB, G., AND LUKASIEWICZ, T. 2009. A general datalog-based framework for tractable query answering over ontologies. In *Proc. PODS*. ACM, 77–86.
- CALÌ, A., GOTTLÖB, G., LUKASIEWICZ, T., MARNETTE, B., AND PIERIS, A. 2010. Datalog<sup>±</sup>: A family of logical knowledge representation and query languages for new applications. In *Proc. LICS*. IEEE Computer Society, 228–242.
- CALÌ, A., GOTTLÖB, G., AND PIERIS, A. 2010. Query rewriting under non-guarded rules. In *Proc. AMW*. CEUR Workshop Proceedings, vol. 619. CEUR-WS.org.
- CALÌ, A., GOTTLÖB, G., AND PIERIS, A. 2011. New expressive languages for ontological query answering. In *Proc. AAI*. AAAI Press.
- CALÌ, A., GOTTLÖB, G., AND PIERIS, A. 2012. Towards more expressive ontology languages: The query answering problem. *Artif. Intell.* 193, 87–128.
- CALVANESE, D., GIACOMO, G. D., LEMBO, D., LENZERINI, M., AND ROSATI, R. 2006. Data complexity of query answering in description logics. In *Proc. KR*. AAAI Press, 260–270.
- CHANDRA, A. K. AND MERLIN, P. M. 1977. Optimal implementation of conjunctive queries in relational data bases. In *Proc. STOC*. ACM, 77–90.
- COURCELLE, B. 1990. The monadic second-order logic of graphs. i. recognizable sets of finite graphs. *Inf. Comput.* 85, 1, 12–75.
- DANTSIN, E., EITER, T., GOTTLÖB, G., AND VORONKOV, A. 2001. Complexity and expressive power of logic programming. *ACM Comput. Surv.* 33, 3, 374–425.
- DEUTSCH, A., NASH, A., AND REMMEL, J. B. 2008. The chase revisited. In *Proc. of PODS*. 149–158.
- DEUTSCH, A. AND TANNEN, V. 2003. Reformulation of XML queries and constraints. In *Proc. of ICDT*. 225–241.
- FAGIN, R., KOLAITIS, P. G., MILLER, R. J., AND POPA, L. 2005. Data exchange: Semantics and query answering. *Theor. Comput. Sci.* 336, 1, 89–124.
- GOTTLÖB, G., MANNA, M., MORAK, M., AND PIERIS, A. 2012. On the complexity of ontological reasoning under disjunctive existential rules. In *MFCS*, B. Rovan, V. Sassone, and P. Widmayer, Eds. Lecture Notes in Computer Science, vol. 7464. Springer, 1–18.

- GRÄDEL, E. 1999. On the restraining power of guards. *J. Symb. Log.* 64, 4, 1719–1742.
- JOHNSON, D. S. AND KLUG, A. C. 1984. Testing containment of conjunctive queries under functional and inclusion dependencies. *J. Comput. Syst. Sci.* 28, 1, 167–189.
- KOLAITIS, P. G., PANTTAJA, J., AND TAN, W. C. 2006. The complexity of data exchange. In *Proc. PODS*. ACM, 30–39.
- MAIER, D., MENDELZON, A. O., AND SAGIV, Y. 1979. Testing implications of data dependencies. *ACM Trans. Database Syst.* 4, 4, 455–469.
- MICROSOFT CORP. 2011. The microsoft connected services framework.
- ORACLE INC. 2011. Oracle database semantic technologies.
- PAPADIMITRIOU, C. M. 1994. *Computational Complexity*. Addison-Wesley, Reading, Massachusetts.
- VARDI, M. Y. 1982. The complexity of relational query languages (extended abstract). In *Proc. STOC*. ACM, 137–146.