

Static Analysis of Schema-Mappings Ensuring Oblivious Termination

Bruno Marnette
Oxford University Computing Laboratory
bruno.marnette@comlab.ox.ac.uk

Floris Geerts
University of Edinburgh
fgeerts@inf.ed.ac.uk

ABSTRACT

A schema-mapping is a high level specification of a data-exchange setting where a set of source-to-target dependencies is used to realize basic operations from source to target relations (such as copy, selection, join or union) while the target schema is subject to a set of target constraints (such as inclusion dependencies or key constraints). In this paper, we consider *strong* schema-mappings that allow for additional constraints such as *source dependencies* on the source schema and *target-to-source dependencies* from the target relations back to the source. Furthermore, strong schema-mappings may include disjunctive dependencies. We argue that this extension is desirable when the source instance is to provide both a lower and upper bound on the information that a target instance can have.

We first focus on the *implication problem* for strong schema-mappings which is to determine whether a given constraint δ is logically implied by the set Σ of constraints (denoted by $\Sigma \models \delta$). After providing complete characterizations for this problem in terms of universal solutions (while supporting equality constraints), we introduce criteria of termination, denoted by TOC, DTOC and MTOC, that allow the efficient computation of universal solutions for standard constraints, disjunctive constraints, and when the source instance is assumed to be immutable (i.e., it is master data), respectively. We obtain decision procedures for the implication problem, provided that Σ satisfies these termination conditions, and give the corresponding complexity bounds. As an immediate application we revisit the problems of determinacy, relative information completeness and variations thereof, all for strong schema-mappings. Indeed, by viewing them as implication problems we obtain efficient decision procedures when the relevant termination conditions are satisfied.

We then focus on the problem of deciding whether source-to-target constraints in a strong schema-mapping are already implied by the embedded (standard) schema-mapping. This problem is important if one wants to use target-to-source constraints in standard data-exchange tools. Since no

such constraints are logically implied by standard schema-mappings (and hence the results established earlier are of no use), we provide an alternative semantics for implication. More specifically, we want the constraint to be satisfied by every solution corresponding to the output of a standard data-exchange tool. We consider three semantics based on universal solutions, cores and CWA-solutions, respectively. Decidability of the implication of general (resp. safe) target-to-source constraints is shown for the CWA-based semantics (resp. core-semantics).

Categories and Subject Descriptors: H.2.5 [Heterogeneous Databases]: Data translation; H.2.4 [Systems]: Relational databases.

General Terms: Algorithms, Theory.

Keywords: Schema mapping, data exchange, data integration, relative completeness, determinacy.

1. INTRODUCTION

A schema-mapping is a high-level specification that describes the relationship between two database schemas. As schema-mapping typically consists of a tuple $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma_{st}, \Sigma_t)$ where \mathbf{S} is a source schema, \mathbf{T} is a target schema, Σ_{st} is a set of *source-to-target dependencies* (tgds) that describe the relationship between the two schemas and Σ_t is a set of *target dependencies* (tgds and egds) that describe constraints on the target instance. Since schema-mappings are crucial for data interoperability tasks (see e.g., [14]), an extensive investigation into the foundations of schema-mappings has been carried out in recent years. In particular, one main line of investigations has focused on the computation of *solutions* of schema-mappings [6, 7, 11, 23]. Here, given \mathcal{M} and a source instance I of \mathbf{S} , a target instance J of \mathbf{T} is a solution for \mathcal{M} and I if $(I, J) \models \Sigma$ for $\Sigma = \Sigma_{st} \cup \Sigma_t$. These solutions play a prominent role throughout this paper.

In spite of its success, schema-mappings are limited in the sense that constraints in Σ only provide a *lower bound* on the set of “true” facts in the solutions. Similarly, a solution is understood as a *sound view* of the schema \mathbf{T} , stating that some facts must hold, but without providing any upper-bound.

EXAMPLE 1.1. Consider the following schema-mapping with source schema

```
{Emp(name, salary, area); Client(name, phone, country)}
```

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

ICDT 2010, March 22–25, 2010, Lausanne, Switzerland.

Copyright 2010 ACM 978-1-60558-947-3/10/0003 ...\$10.00

which stores employees and clients in a world-wide company. The UK branch of the company has target schema

$\{\text{Emp}_{\text{UK}}(\text{name}, \text{salary}); \text{Client}_{\text{UK}}(\text{name}, \text{phone})\}$.

The following set Σ_{st} of source-to-target dependencies can be used to (1) export the list of employees already assigned to the UK while authorizing the branch to adapt their salary and (2) export the list of clients located in the UK:

$$\begin{aligned} \text{Emp}(e, s, \text{UK}) &\rightarrow \exists s' \text{Emp}_{\text{UK}}(e, s') \\ \text{Client}(n, p, \text{UK}) &\rightarrow \text{Client}_{\text{UK}}(n, p) \end{aligned}$$

If we consider the source instance I below (left) then a solution J for \mathcal{M} and I is given below (right).

I	$\text{Emp}(\text{alice}, 3000, \text{UK})$ $\text{Emp}(\text{bob}, 2500, \text{USA})$ $\text{Emp}(\text{cedric}, 3500, \text{UK})$ $\text{Client}(\text{denise}, 01234098, \text{UK})$ $\text{Client}(\text{edward}, 01899785, \text{AU})$ $\text{Client}(\text{fred}, 06947647, \text{US})$	J
		$\text{Emp}_{\text{UK}}(\text{alice}, \#_a)$ $\text{Emp}_{\text{UK}}(\text{cedric}, \#_c)$ $\text{Client}_{\text{UK}}(\text{denise}, 01234098)$

Any other solutions is obtained by instantiating the *nulls* $\#_a$ and $\#_c$ in J arbitrarily. Hence, no upper bound on neither those values nor the cardinality of solutions is provided. \square

In many practical settings, however, it is desirable to also provide an ‘‘upper bound’’ on the solutions of a schema-mapping.

EXAMPLE 1.2. Continuing the previous example, suppose that a list Sal of possible salaries for employees in the UK branch is available. One could then enforce the nulls to take values from this list, provided *disjunctive target* constraints are allowed. Indeed, adding

$$\forall xy \text{Emp}_{\text{UK}}(x, y) \rightarrow \bigvee_{s \in \text{Sal}} y = s$$

to Σ_t would provide the desired upper bound. Similarly, suppose that one wants to encode that every employee of the UK branch has to be registered in the main company, or equivalently, that the branch cannot hire a new employee without notifying the main company (which would lead to updating the source instance). For this, *target-to-source dependencies* are needed. Indeed, the dependency

$$\delta : \text{Emp}_{\text{UK}}(e, s) \rightarrow \exists s', l' \text{Emp}(e, s', l')$$

suffices for this. Furthermore, one needs to distinguish between the case when the source instance I is assumed to be immutable (i.e., I is so-called *master data* using the terminology from [8]), and when I can be updated. In the former case, J can be updated to J' below (where the salaries have been fixed and where clients have been added) but cannot be updated to J'' below (where a new employee has been added). When I can be updated, both solutions are feasible.

J'	$\text{Emp}_{\text{UK}}(\text{alice}, 3000)$ $\text{Emp}_{\text{UK}}(\text{cedric}, 3500)$ $\text{Client}_{\text{UK}}(\text{denise}, 01234098)$ $\text{Client}_{\text{UK}}(\text{georg}, 01993707)$ $\text{Client}_{\text{UK}}(\text{isabel}, 01993707)$	J''
		$\text{Emp}_{\text{UK}}(\text{alice}, 3000)$ $\text{Emp}_{\text{UK}}(\text{cedric}, 3500)$ $\text{Emp}_{\text{UK}}(\text{heather}, 3500)$ $\text{Client}_{\text{UK}}(\text{denise}, 01234098)$ $\text{Client}_{\text{UK}}(\text{georg}, 01993707)$

Motivated by the previous examples, it therefore seems natural to extend schema-mappings to *strong* schema-mappings $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$ where Σ now consists of Σ_t (target dependencies), Σ_{st} (source-to-target dependencies), Σ_s (source dependencies) and finally Σ_{ts} (target-to-source dependencies).

This more general setting, has been previously studied by Deutsch et al [4]. In fact, it is shown there that most concepts and techniques from standard schema-mappings generalize to the strong schema-mapping setting. In particular, they studied the *implication problem* which is to decide, given a strong schema-mapping \mathcal{M} and a dependency δ over $\mathbf{S} \cup \mathbf{T}$, whether Σ *logically implies* δ , denoted by $\Sigma \models \delta$. In other words, it is to decide whether for any source instance I of \mathbf{S} and *any* solution J for \mathcal{M} and I , $(I, J) \models \delta$. It is shown in [4] that when Σ consists of *tgds only, universal* solutions can be used to determine $\Sigma \models \delta$, provided that these solutions can be computed. However, no termination conditions of the chase procedure (that returns universal solutions) are established in [4] when Σ consists of disjunctive constraints or when the source instance is assumed to be immutable.

We therefore believe that the implication problem for (strong) schema-mappings needs revisiting and this is indeed the first main focus of this paper. Our first contributions (Section 3) are the following: Let $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$ be a strong schema-mapping and δ a constraint over $\mathbf{S} \cup \mathbf{T}$.

- We provide characterizations for when $\Sigma \models \delta$ holds in terms of universal solutions when Σ possibly contains egds and/or disjunctive constraints, hereby extending the work of [4].
- Termination conditions of the *oblivious chase* [17] are given for standard constraints, disjunctive constraints and for the case when the source instance is assumed to be immutable (master data). The corresponding classes of strong schema-mappings are denoted by TOC, DTOC and MTOC, respectively.
- Finally, complexity bounds for $\Sigma \models \delta$ are established under the assumption that the constraints belong to TOC, DTOC or MTOC. In particular, the implication problem is in NP when $\mathcal{M} \in \text{TOC}$, and in Π_2^P when $\mathcal{M} \in \text{DTOC}$ or $\mathcal{M} \in \text{MTOC}$. The complexity bounds carry over to the class of weakly acyclic [6] or super-weakly acyclic [17] schema-mappings.

As an immediate application of these results, we revisit the problems of relative information completeness [8] and determinacy [21, 19] in the strong schema-mapping setting in Section 4. Indeed, we show that the two problems can be viewed as an implication problem and hence, by leveraging the results in Section 3, we establish complexity bounds for both these problems under some termination conditions.

The second main focus of this paper addresses the important question whether the target-to-source constraints in a strong schema-mapping are already implied by the embedded (standard) schema-mapping. More precisely, given a strong schema-mapping $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$, where as before Σ consists of Σ_{st} , Σ_{ts} , Σ_s and Σ_t , does $\mathcal{M}' = (\mathbf{S}, \mathbf{T}, \Sigma_{st}, \Sigma_t)$, already imply the constraints in Σ_{ts} ? This question is important for the reason that current schema-mapping management systems (such as e.g., [12]) do not support target-

to-source constraints. Consequently, only target-to-source constraints can be included that are “harmless” (i.e., already implied by standard constraints).

Note that the results in Section 3 are of no use for this problem. Indeed, standard schema-mappings can only logically entail source-to-target or target constraints. We are therefore faced with the problem of properly defining the semantics of entailment for target-to-source constraints in the (standard) schema-mapping setting, while at the same time ensuring decidability of the implication problem for large classes of constraints. Let $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$ be a strong schema-mapping and assume that $\Sigma_{ts} = \{\delta\}$. We explore under which semantics of solutions (denoted by **Sem**) we can be sure that, for all source instances I and every J in $\mathbf{Sem}(\mathcal{M}, I)$ we have $(I, J) \models \delta$.

Our contributions (Section 5) consists of the following:

- ▶ We show that the semantics based on *universal solutions* (**USol**) is not restrictive enough while the semantics based on *cores of universal solutions* (**Core**) is on the contrary too restrictive and leads very easily to undecidability, even for weakly-acyclic schema-mappings.
- ▶ We propose an intermediate semantics, which in the case of tgds only, coincides with the notion of *closed-world assumption* semantics [15, 13].
- ▶ We show that when the embedded standard schema-mapping $\mathcal{M}' \in \mathbf{TOC}$, the CWA-based implication problem is decidable for arbitrary target-to-source dependencies. The core-based implication problem is decidable for *safe* target-to-source constraints only.

Organization. In Section 2 we define schema-mappings and universal solutions. We define strong schema-mappings and study its corresponding implication problem in Section 3, and show as an application its use for the problems of completeness and determinacy in Section 4. The implication problem for target-to-source constraints by the embedded standard schema-mapping is studied in Section 5. Finally, Section 6 discusses related work.

2. PRELIMINARIES

Terms. We fix a universal schema \mathbf{U} consisting of an infinite set of predicate symbols R of fixed and finite arity a_R . We define a schema σ as a subset of \mathbf{U} . We also fix an infinite set \mathbf{Cst} of constants, an infinite set \mathbf{Var} of variables, and a set \mathbf{Funcs} of function symbols f of fixed and finite arity a_f . We define the set **Terms** of terms as the minimal set containing $\mathbf{Cst} \cup \mathbf{Var}$ and such that $f\langle t_1, \dots, t_n \rangle$ is a term whenever $f \in \mathbf{Funcs}$, $n = a_f$ and $t_1, \dots, t_n \in \mathbf{Terms}$. We inductively define the *size* $|t|$ of a term t as $|t| = 1$ if $t \in \mathbf{Cst} \cup \mathbf{Var}$ and $|t| = 1 + |t_1| + \dots + |t_n|$ if $t = f\langle t_1, \dots, t_n \rangle$. We let \mathbf{Dom} be the set of terms in which no variables occur, and we define the set \mathbf{Nulls} of nulls as the set of non-constant terms in \mathbf{Dom} so that $\mathbf{Dom} = \mathbf{Cst} \uplus \mathbf{Nulls}$. We say that a null is *flat* if it is of the form $f\langle \rangle$ for some $f \in \mathbf{Funcs}$ of arity zero. In this paper, we use either integers $1, 2, \dots$ or (bold) letters $\mathbf{a}, \mathbf{b}, \dots$ to denote constants, we use x, y, z, u, v, \dots to denote variables, and we use $\#_1, \#_2, \dots$ to denote flat nulls.

Instances and Morphisms. We define a *fact* as an object $R\langle t_1, \dots, t_n \rangle$ where $R \in \mathbf{U}$, $n = a_R$ and $t_1, \dots, t_n \in \mathbf{Dom}$. We define an *instance* as a finite or infinite set of facts.

(The semantics of finite instances is discussed in Section 6.) Given an instance I , we let $\mathbf{Cst}(I)$ and $\mathbf{Nulls}(I)$ be the set of constants and nulls occurring in I and we define the active domain of I as $\mathbf{Dom}(I) = \mathbf{Cst}(I) \cup \mathbf{Nulls}(I)$. We say that an instance I is *ground* iff $\mathbf{Dom}(I) = \mathbf{Cst}(I)$ and $\mathbf{Nulls}(I) = \emptyset$. We denote by σ_I the set of predicate symbols occurring in an instance I and say that I is *over* a schema σ iff $\sigma_I \subseteq \sigma$. Given a mapping $h : \mathbf{Dom}(I) \rightarrow \mathbf{Dom}$, we denote by $h(I)$ the instance containing the fact $R\langle h(t_1), \dots, h(t_n) \rangle$, for each fact $R\langle t_1, \dots, t_n \rangle$ in I . Given two instances I and J , an *homomorphism* h from I to J is a mapping $h : \mathbf{Dom}(I) \rightarrow \mathbf{Dom}(J)$ such that $h(I) \subseteq J$ and $h(c) = c$ for all constant $c \in \mathbf{Cst}(I)$. We write $I \xrightarrow{\text{hom}} J$ iff there exists a homomorphism from I to J .

Tgds and Egds. An *atom* is an object $R\langle t_1, \dots, t_n \rangle$ where $R \in \mathbf{U}$, $n = a_R$ and $t_1, \dots, t_n \in \mathbf{Var} \cup \mathbf{Cst}$ (note the difference with the definition of a fact given earlier). Given a set of atoms we let $\mathbf{Var}(\phi)$ and $\mathbf{Cst}(\phi)$ be the sets of variables and constants occurring in ϕ . We define a *tuple-generating dependency*, or *tgd* for short, as a first-order formula r of the form

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

where ϕ and ψ are two finite sets of atoms called respectively the *body* and the *head* of r and \bar{x}, \bar{y} and \bar{z} are three disjoint tuples of distinct variables such that $\{\bar{x}\} = \mathbf{Var}(\phi) \cap \mathbf{Var}(\psi)$; $\{\bar{y}\} = \mathbf{Var}(\phi) \setminus \{\bar{x}\}$ and $\{\bar{z}\} = \mathbf{Var}(\psi) \setminus \{\bar{x}\}$. We define an *equality-generating dependency*, or *egd* for short, as a first-order formula s of the form

$$\forall \bar{x}, \phi(\bar{x}) \rightarrow \alpha(\bar{x}) = \beta(\bar{x})$$

where ϕ is a finite set of atoms called the *body* of s , and $\alpha(\bar{x})$ and $\beta(\bar{x})$ are two terms consisting either of a constant or a variable of \bar{x} .

Disjunctive dependencies. We define a *disjunctive tgd* as a first-order formula of the form

$$\forall \bar{x}, \bar{y}, \phi(\bar{x}, \bar{y}) \rightarrow (\exists \bar{z}_1, \psi_1(\bar{x}, \bar{z}_1)) \vee \dots \vee (\exists \bar{z}_k, \psi_k(\bar{x}, \bar{z}_k))$$

where for each $i \in [1, k]$, $\forall \bar{x}, \bar{y}, \phi(\bar{x}, \bar{y}) \rightarrow \exists \bar{z}_i, \psi(\bar{x}, \bar{z}_i)$ is a tgd. We define a *disjunctive egd* as a first-order formula of the form

$$\forall \bar{x}, \phi(\bar{x}) \rightarrow (e_{1,1} \wedge \dots \wedge e_{1,p_1}) \vee \dots \vee (e_{k,1} \wedge \dots \wedge e_{k,p_k})$$

where for each $i \in [1, k]$ and each $j \in [1, p_k]$, $e_{i,j}$ denotes an equality $\alpha_{i,j}(\bar{x}) = \beta_{i,j}(\bar{x})$ and $\forall \bar{x} \phi(\bar{x}) \rightarrow e_{i,j}$ is an egd.

Given a set Σ of (disjunctive) tgds and egds we denote by $\mathbf{Cst}(\Sigma)$ the set of constants occurring in Σ and by σ_Σ the set of predicate symbols occurring in Σ . Given an instance I we write $I \models \Sigma$ iff I is a model of Σ according to the standard semantics of FO. Given two sets Σ and Σ' of constraints we write $\Sigma \models \Sigma'$ iff for all (possibly infinite) instance I we have that $I \models \Sigma$ implies $I \models \Sigma'$. In the following sections, we often omit to write explicitly the universal quantification in front of the (disjunctive) tgds and egds.

Schema-Mappings. A schema-mapping is a tuple $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma_{st}, \Sigma_t)$ where: (i) \mathbf{S} and \mathbf{T} are two disjoint finite schemas respectively called *source schema* and *target schema*; (ii) Σ_{st} is a finite set of tgds called *source-to-target tgds* of the form $\phi \rightarrow \psi$ for some ϕ over \mathbf{S} and some ψ over \mathbf{T} ; and (iii) Σ_t is a finite set of tgds over \mathbf{T} called *target tgds* and egds over \mathbf{T} called *target egds*. Given a schema-mapping \mathcal{M} ,

a *source instance* for \mathcal{M} is a ground instance over \mathbf{S} (possibly infinite, with constants only) while a *target instance* for \mathcal{M} is an instance over \mathbf{T} (possibly infinite and with nulls).

Universal Solutions and Cores. Given a schema-mapping $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma_{st}, \Sigma_t)$ and a source instance I , a target instance J is a *universal solution* and we write $J \in \text{USol}(\mathcal{M}, I)$ iff (i) J is a solution, meaning that $I \cup J \models (\Sigma_{st} \cup \Sigma_t)$, and (ii) J is a universal, meaning that for all solutions K we have $J \xrightarrow{\text{hom}} K$. A universal solution $J \in \text{USol}(\mathcal{M}, I)$ is a *core of universal solutions* and we write $J \in \text{Core}(\mathcal{M}, I)$ iff J is finite, and there exists no strict subset $J' \subset J$ such that $J' \in \text{USol}(\mathcal{M}, I)$.

Universal models. In general, one may consider any finite set Σ of tgds and egds (not necessarily restricted to the schema-mapping format). Similar to the notion of universal solution, and following the terminology of [4], an instance J is a *universal model* of Σ and an instance I if (i) $I \xrightarrow{\text{hom}} J$; (ii) $J \models \Sigma$; and (iii) for every instance K , if $K \models \Sigma$ and $I \xrightarrow{\text{hom}} K$, then $J \xrightarrow{\text{hom}} K$. We denote by $\text{UMod}(\Sigma, I)$ the set of universal models of Σ and I .

3. LOGICAL ENTAILMENT

The *implication problem* for a class \mathcal{C} of dependencies is the following: given a finite set Σ of dependencies, $\Sigma \subseteq \mathcal{C}$, and a dependency $\delta \in \Sigma$, does $\Sigma \models \delta$? In this section, we study the implication problem for (disjunctive) tgds and egds, a problem which is also known as the problem of *logical entailment* in the context of schema-mappings. In order to capture the generality promised in the introduction, we first define a notion of *strong* schema-mappings.

DEFINITION 3.1. We define a *strong schema-mapping* \mathcal{M} as a triple $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$ where $\Sigma = \Sigma_s \uplus \Sigma_{st} \uplus \Sigma_{ts} \uplus \Sigma_t$ is a set of disjunctive tgds and egds decomposed as follows:

- Σ_{st} is a set of disjunctive tgds from \mathbf{S} to \mathbf{T} ;
- Σ_t is a set of disjunctive tgds and egds over \mathbf{T} ;
- Σ_{ts} is a set of disjunctive tgds from \mathbf{T} to \mathbf{S} ;
- Σ_s is a set of disjunctive tgds and egds over \mathbf{S} .

In this definition, Σ_{st} corresponds to the source-to-target dependencies of a standard schema-mapping, Σ_t and Σ_s captures the dependencies that often come with a schema (e.g., tgds encoding inclusion dependencies, egds encoding key constraints, and disjunctive egds encoding domain constraints), and Σ_{ts} is a set of target-to-source dependencies allowing to encode complete views (see also Section 5).

Given a strong schema-mapping $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$ and a source instance I , we denote by $\text{Sol}(\mathcal{M}, I)$ the set of target instances J such that $(I, J) \models \Sigma$. The following observation is straightforward.

OBSERVATION 3.1. Given a strong schema-mapping $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$ and a formula δ over $\mathbf{S} \cup \mathbf{T}$, the following statements are equivalent:

- $\Sigma \models \delta$
- for every source instance I and every target instance $J \in \text{Sol}(\mathcal{M}, I)$ we have $(I, J) \models \delta$.

To obtain decidability for the implication problem, we provide a complete characterization that (unlike previous approaches) supports egds. More specifically, we first recall the characterization given in [4] and show that it is a complete characterization in the case of tgds only. We then show how the notion of *substitution-free simulation* of egds introduced in [11, 17] can be used to also provide a complete characterization in the presence of egds.

3.1 Implication under Tgds and Egds

We first recall a well-known (partial) characterization of the implication problem for tgds and egds. The constraints in this and the following section are non-disjunctive; disjunctive constraints are considered in Section 3.3

PROPOSITION 3.1 ([4]). *Given a finite set Σ of tgds and egds, given a tgd δ of the form $\phi(\bar{x}, \bar{y}) \rightarrow \exists z \psi(\bar{x}, \bar{z})$, given two tuples of fresh constants $\bar{a} \in \text{Dom}^{|\bar{x}|}$ and $\bar{b} \in \text{Dom}^{|\bar{y}|}$, and given a universal model $K \in \text{UMod}(\Sigma, \phi(\bar{a}, \bar{b}))$ we have $\Sigma \models \delta$ iff $K \models \exists \bar{z}, \psi(\bar{a}, \bar{z})$.*

When Σ consists only of tgds, it can be seen that a universal model K necessarily exists (even though it might be infinite) and therefore the characterization offered by Proposition 3.1 is complete in the case without egds. More precisely, using the notations of Proposition 3.1 the following statements are all equivalent when Σ is a set of tgds only:

- $\Sigma \models \delta$
- $\exists K \in \text{UMod}(\Sigma, \phi(\bar{a}, \bar{b})), K \models \exists \bar{z}, \psi(\bar{a}, \bar{z})$
- $\forall K \in \text{UMod}(\Sigma, \phi(\bar{a}, \bar{b})) K \models \exists \bar{z}, \psi(\bar{a}, \bar{z})$

An important observation here is that the above characterization is not complete in the presence of egds, simply because universal models may not exist. This is illustrated in the following example.

EXAMPLE 3.1. Consider $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$ with $\mathbf{S} = \emptyset$, $\mathbf{T} = \{A, B, C\}$ and $\Sigma = \Sigma_t, \delta_1$ and δ_2 defined as:

$$\Sigma_t = \left\{ \begin{array}{l} A(x, y) \rightarrow x = y \\ A(x, x) \wedge B(y) \rightarrow x = y \\ A(x, x) \wedge B(x) \rightarrow C(x) \end{array} \right\}$$

$$\delta_1 = A(x, y) \rightarrow C(x)$$

$$\delta_2 = A(x, y) \wedge B(z) \rightarrow C(y)$$

Clearly, $\Sigma \not\models \delta_1$ and $\Sigma \models \delta_2$. It is easily verified that no universal model exists in neither $\text{UMod}(\Sigma, \{A(a, b)\})$ (for δ_1) nor $\text{UMod}(\Sigma, \{A(a, b), B(c)\})$ (for δ_2). Hence, Proposition 3.1 is of no use to decide whether δ_1 or δ_2 are implied.

We next show how egds can be incorporated by means of a transformation process (see [11, 17]) that transforms a set Σ consisting of tgds and egds into a set $s(\Sigma)$ consisting of tgds only.

Simulation $\epsilon(\Sigma)$. Given a set Σ of tgds and egds, we let $\epsilon(\Sigma)$ be the set of tgds obtained from Σ by: (i) replacing each equality atom $(x = y)$ by $E(x, y)$, where E is a fresh binary predicate symbol; (ii) adding the tgds $E(x, y) \rightarrow E(y, x)$ and $E(x, y) \wedge E(y, z) \rightarrow E(x, z)$ to encode the symmetry and transitivity of E ; and adding finally for all predicates R occurring in Σ the tgd $R(x_1, \dots, x_{a_R}) \rightarrow E(x_1, x_1) \wedge \dots \wedge E(x_{a_R}, x_{a_R})$ to encode the reflexivity of E .

Singularisation. Given a finite set of atoms ϕ over a schema σ (that does not contain E) we define a *singularisation* of ϕ as a set of atoms $\phi^\sigma \cup \phi^E$ obtained, starting from $\phi^\sigma := \phi$ and $\phi^E := \emptyset$, and by repeating the following operations until no constant occur in ϕ^σ and no variable occurs more than once in ϕ^σ :

- if a constant \mathbf{c} occurs in ϕ^σ , introduce a fresh variable u , replace every occurrence of \mathbf{c} by u in ϕ^σ and add the atom $E(u, \mathbf{c})$ to ϕ^E ;
- if a variable x occurs more than once in ϕ^σ , introduce a fresh variable x' , replace one occurrence of x by x' in ϕ^σ and add the atom $E(x, x')$ to ϕ^E .

Note that there may be an exponential number of possible non-equivalent singularizations for a given set ϕ of atoms, but to simplify the discussion, we assume a fixed algorithm **sing** that, given a set ϕ of atoms, returns a singularisation **sing**(ϕ) of ϕ . Given a set Σ of tgds and egds, we then define the *substitution-free simulation* $s(\Sigma)$ of Σ as the set of tgds obtained from $\epsilon(\Sigma)$ by replacing each tgd ($\phi \rightarrow \psi$) in $\epsilon(\Sigma)$ by (**sing**(ϕ) \rightarrow ψ).

EXAMPLE 3.2. Consider again the target constraints of Example 3.1. The simulation $\epsilon(\Sigma)$ consists of

$$\epsilon(\Sigma) = \left\{ \begin{array}{ll} A(x, y) \rightarrow E(x, y) & (r_1) \\ A(x, x) \wedge B(y) \rightarrow E(x, y) & (r_2) \\ A(x, x) \wedge B(x) \rightarrow C(x) & (r_3) \\ E(x, y) \wedge E(y, z) \rightarrow E(x, z) & (\text{transitivity}) \\ E(x, y) \rightarrow E(y, z) & (\text{symmetry}) \\ A(x, y) \rightarrow E(x, x) \wedge E(y, y) & (\text{reflexivity}) \\ B(y) \rightarrow E(y, y) & (\text{reflexivity}) \\ C(y) \rightarrow E(y, y) & (\text{reflexivity}) \end{array} \right\}$$

In this case, the substitution-free simulation $s(\Sigma)$ only replaces r_2 and r_3 by

$$\begin{array}{l} r'_2 : A(x, x') \wedge B(y) \wedge E(x, x') \rightarrow E(x, y) \\ r'_3 : A(x, x') \wedge B(x'') \wedge E(x, x') \wedge E(x', x'') \rightarrow C(x) \quad \square \end{array}$$

We are now almost ready to provide a complete characterization for $\Sigma \models \delta$ in the presence of egds. Given an instance I where the predicate symbol E is used to encode equality constraints, we define I^+ as the instance obtained by applying to I the substitution axioms of equality, that is, I^+ is the set of atoms $R(t_1, \dots, t_n)$ such that I contains $\{R(t'_1, \dots, t'_n); E(t_1, t'_1); \dots; E(t_n, t'_n)\}$ for some terms $t'_1, \dots, t'_n \in \text{Dom}(I)$. Given a tgd δ of the form $\phi(\bar{x}, \bar{y}) \rightarrow \exists z \psi(\bar{x}, \bar{z})$, and considering two tuples of fresh constants $\bar{a} \in \text{Dom}^{|\bar{x}|}$ and $\bar{b} \in \text{Dom}^{|\bar{y}|}$ we let $B_\delta = \phi(\bar{a}, \bar{b})$ and $H_\delta = \exists \bar{z}, \phi(\bar{a}, \bar{z})$. Given an egd δ of the form $\phi(\bar{x}) \rightarrow \alpha(\bar{x}) = \beta(\bar{x})$ and considering a tuple $\bar{a} \in \text{Dom}^{|\bar{x}|}$ of fresh constants, we let $B_\delta = \phi(\bar{a})$ and $H_\delta = E(\alpha(\bar{a}), \beta(\bar{a}))$.

PROPOSITION 3.2. *Given a finite set Σ of tgds and egds, given a tgd or egd δ , and given a universal model $K \in \text{UMod}(s(\Sigma), B_\delta)$, we have $\Sigma \models \delta$ iff one of the following statements hold:*

- $K^+ \models H_\delta$, or
- $\exists c, c' \in \text{Cst}(\Sigma \cup \delta)$, $c \neq c' \wedge E(c, c') \in K^+$.

The following examples shed light on the two conditions in Proposition 3.2.

EXAMPLE 3.3. Consider $s(\Sigma)$ given in Example 3.2 and δ_2 from Example 3.1. The following instance K is a universal model for $s(\Sigma)$ and $B_{\delta_2} = \{A(a, b), B(c)\}$:

$$K = \{ \begin{array}{l} A(a, b); B(c); E(a, b); E(b, a); E(b, c); \\ E(c, b); E(a, c); E(c, a); C(a) \end{array} \}$$

Since $C(b) \in K^+$, we have $\Sigma \models \delta_2$. \square

EXAMPLE 3.4. Consider $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$ with $\mathbf{S} = \emptyset$, $\mathbf{T} = \{A, B\}$ and $\Sigma = \{A(x) \rightarrow x = 1\}$ and $\delta = A(0) \rightarrow B(0)$. Here $s(\Sigma)$ contains $A(x) \rightarrow E(x, 1)$ and a universal solution K for $s(\Sigma)$ and $B_\delta = \{A(0)\}$ clearly exists. However, K must contain $E(0, 1)$ which means that two different constants must be identified. Since these constants originate from Σ and δ , this not allowed. Note that this precisely corresponds to the fact that $\Sigma \wedge \{A(0)\}$ is unsatisfiable. Hence, in this case, δ is vacuously implied by Σ , a situation that did not occur in the tgd case only. \square

3.2 Oblivious Termination

In order to obtain decidability for the implication problem, Proposition 3.2 tells that it suffices to identify classes of dependencies allowing (after substitution-free simulation) the computation of universal models. In this section, we first recall the notion of *oblivious termination* introduced in [17] for tgds and egds before extending this notion to disjunctive dependencies in the next section.

Skolemization P_Σ . Given a set of tgds Σ we define the *Skolemization* P_Σ of Σ as the logic program obtained by replacing each tgd of the form

$$\phi(\bar{x}, \bar{y}) \rightarrow \exists z_1, \dots, z_n \psi(\bar{x}, z_1, \dots, z_n)$$

by a rule of the form

$$\phi(\bar{x}, \bar{y}) \rightarrow \psi(\bar{x}, f_1(\bar{x}), \dots, f_n(\bar{x}))$$

where each f_i is a fresh function symbol in **Funcs** of arity $|\bar{x}|$. Given a logic program P_Σ and an instance I we denote by $P_\Sigma(I)$ the minimal Herbrand model of $P_\Sigma(I)$, or equivalently, the least-fixed point of I by P_Σ . We say that P_Σ *terminates* on I iff $P_\Sigma(I)$ is finite.

PROPOSITION 3.3 ([17]). *For any instance I and set Σ of tgds we have that $P_\Sigma(I) \in \text{UMod}(\Sigma, I)$.*

TOC. Given a finite set of tgds Σ , we say that Σ *ensures oblivious termination* and write $\Sigma \in \text{TOC}$ iff P_Σ terminates on every finite ground instance. Given a set Σ of tgds and egds, we write $\Sigma \in \text{TOC}^e$ iff $s(\Sigma) \in \text{TOC}$. It can easily be observed that a set Σ of tgds (without egds) is in **TOC** iff $s(\Sigma)$ is in **TOC**^e and we can therefore drop the “e” in the notation **TOC**^e without causing any ambiguity. Note that the class **TOC** was originally defined in [17] for schema-mappings $(\mathbf{S}, \mathbf{T}, \Sigma)$ where $\Sigma = \Sigma_{st} \cup \Sigma_t$ and was only requiring the termination of P_Σ (or $P_{s(\Sigma)}$) for every finite ground source instance I . This requirement is however not sufficient in the context of logical entailment because P_Σ may terminate on every finite source instance (over \mathbf{S}) without necessarily terminating on every finite instance (over $\mathbf{S} \cup \mathbf{T}$).

Critical Instance. Let Σ be a set of dependencies and let $*$ be a fresh constant ($*$ $\notin \text{Cst}(\Sigma)$). We define the critical instance I_Σ of Σ as the set of all the facts $R(t_1, \dots, t_n)$ such that R is a predicate symbol occurring in Σ (i.e., $R \in \sigma_\Sigma$) and $\{t_1, \dots, t_n\} \subseteq \text{Cst}(\Sigma) \cup \{*\}$.

PROPOSITION 3.4 ([17]). *For every finite set of tgds Σ , the following statements are equivalent:*

- Σ ensures oblivious termination ($\Sigma \in \text{TOC}$);
- P_Σ terminates on the critical instance I_Σ ;
- there exists a finite bound k depending only on Σ such that, for any instance I over $\mathbf{S} \cup \mathbf{T}$ and any term t occurring in $P_\Sigma(I)$, the size $|t|$ of t is at most k .

It follows from this proposition that, given a set of tgds and egds $\Sigma \in \text{TOC}$, and given an instance I , we can compute the universal model $P_\Sigma(I)$ in time $O(|\text{Dom}(I)|^k)$ for some k depending only on Σ . Combined with Proposition 3.2 we have the following theorem.

THEOREM 3.1. *For every strong schema-mapping $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$ without disjunction, and more generally, for every finite set Σ of tgds and egds, if $\Sigma \in \text{TOC}$ the following problem is in NP: given a tgd δ , does $\Sigma \models \delta$? Moreover, the problem becomes PTIME if δ is an egd or if δ is a tgd with a bounded number of atoms in its head.*

Since containment of CQ queries is NP-complete [1] and can be viewed as an implication problem, NP-hardness is also easily verified.

3.3 Supporting Disjunction

In this section, we extend our framework to also support disjunctive tgds and egds. In the case of disjunctive tgds (and egds) it as been shown in [4] that there is generally no (single) universal model and the more appropriate notion of *universal model set* was introduced in [4]. Given an instance I and a set of disjunctive tgds Σ , a universal model set is a set of instances $\mathbf{J} = \{J_0, J_1, J_2, \dots\}$ such that (i) for all $i \geq 0$ we have $J_i \models \Sigma$ and $I \xrightarrow{\text{hom}} J_i$, and (ii) for every instance J' such that $J' \models \Sigma$ and $I \xrightarrow{\text{hom}} J'$ there exists some $i \geq 0$ such that $J_i \xrightarrow{\text{hom}} J'$. We denote by $\text{UModS}(\Sigma, I)$ the universal model sets for Σ and I .

We now generalize Proposition 3.2 to the case of disjunctive tgds and egds by considering universal models sets and extending the definitions of $s(\Sigma)$, B_δ and H_δ in a natural way. More precisely: the singularisation of a disjunctive tgd $\phi \rightarrow \bigvee_i \phi_i$ is simply defined as $\text{sing}(\phi) \rightarrow \bigvee_i \psi_i$; B_δ is defined exactly as in the case of tgds and egds; and H_δ is the union of boolean conjunctive queries corresponding to the head of δ (while replacing “=” by “E” in the case of disjunctive egds).

PROPOSITION 3.5. *Given a finite set Σ of disjunctive tgds and egds, given a disjunctive tgd or egd δ , and given a universal model set $\mathbf{K} \in \text{UModS}(s(\Sigma), B_\delta)$, we have $\Sigma \models \delta$ iff for every $K \in \mathbf{K}$ one of the following statements hold:*

- $K^+ \models H_\delta$, or
- $\exists c, c' \in \text{Cst}(\Sigma \cup \delta)$, $c \neq c' \wedge E(c, c') \in K^+$.

We next generalize the class TOC, defined only for tgds and egds, to the case of disjunctive dependencies. For this, we first generalize the process of Skolemization. Given a set Σ of disjunctive tgds we define the skolemization P_Σ of Σ as the set of rules of the form $B^s \rightarrow H_1^s \vee \dots \vee H_k^s$ such that $(B \rightarrow H_1 \vee \dots \vee H_k) \in \Sigma$ and each $(B^s \rightarrow H_i^s)$ corresponds to the standard skolemization of the tgd $(B \rightarrow H_i)$. Given a logic program P_Σ as above and an instance I we define a P_Σ -derivation of I as a (possibly infinite) series $(I_0, I_1, \dots, I_i, \dots)$ such that $I_0 = I$ and for all $i > 0$ there exists a justification (r_i, \bar{c}_i, j_i) of I_i such that:

- r_i is a rule $B^s(\bar{x}) \rightarrow H_1^s(\bar{x}) \vee \dots \vee H_k^s(\bar{x})$ in P_Σ ;
- $\bar{c}_i \in \text{Dom}(I_{i-1})^{|\bar{x}|}$; and $B^s(\bar{c}_i) \subseteq I_{i-1}$;
- $j_i \in \{1, \dots, k\}$ and $I_i = I_{i-1} \cup H_{j_i}^s(\bar{c}_i)$
- there is no $i' < i$ such that $(r_{i'}, \bar{c}_{i'}) = (r_i, \bar{c}_i)$

Given a finite set Σ of disjunctive tgds and an instance I we say that P_Σ terminates on I iff all the P_Σ -derivations of I are finite. We say that a finite P_Σ derivation (I_0, I_1, \dots, I_n) of I leads to I_n iff there is no $I_{n+1} \supset I_n$ such that $(I_0, I_1, \dots, I_n, I_{n+1})$ is a P_Σ -derivation of I . We finally denote by $P_\Sigma(I)$ the sets of instances J such that some P_Σ -derivation of I leads to J .

PROPOSITION 3.6. *For every finite set Σ of disjunctive tgds and every instance I such that P_Σ terminates on I , we have $P_\Sigma(I) \in \text{UModS}(\Sigma, I)$.*

DEFINITION 3.2. *We say that a finite set Σ of disjunctive tgds and egds ensures oblivious termination, denoted by $\Sigma \in \text{DTC}$, iff $P_{s(\Sigma)}$ terminates on every finite instance I .*

PROPOSITION 3.7. *For any finite sets Σ of disjunctive tgds, the following statements are equivalent:*

- Σ ensures oblivious termination ($\Sigma \in \text{DTC}$);
- P_Σ terminates on the critical instance I_Σ ;
- there exists a bound k , depending only on Σ such that, for any instances I and any term t occurring in a P_Σ -derivation of I , the size $|t|$ of t is at most k .

From this, the generalization of Proposition 3.1 to disjunctive constraints easily follows:

THEOREM 3.2. *For every strong schema-mappings $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$, and more generally for every set Σ of disjunctive tgds and egds, if $\Sigma \in \text{DTC}$ the following problem is in Π_2^P : given a disjunctive tgd δ , does $\Sigma \models \delta$? Moreover, the problem is in coNP if δ is a disjunctive egd or a disjunctive tgd where each head has a bounded number of atoms.*

Matching lower bounds can be obtained by a reduction from the containment of UCQ queries, which is Π_2^P -complete [20], and the tautology problem, which is coNP-complete.

3.4 Supporting Master Data

In this section we consider the problem of implication in the case where a source instance I is immutable while the target instance J can be any solution. In the context of strong schema-mappings, this setting corresponds to a situation of *semi-dynamic* data-exchange where the target instance J

can evolve in time after its creation while the initial source instance I is not subject to evolution. This situation occurs, e.g., when I is so-called *master data* [8] which is assumed to be complete and consistent. As illustrated in the introduction, master data together with target-to-source constraints allow to provide an upper bound on the set of true facts in solutions. Note that in this section the distinction between the source schema and the target schema is crucial.

DEFINITION 3.3. *Given a strong schema-mapping $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$, a source instance I , and a constraint δ over \mathbf{T} , we write $\mathcal{M} \stackrel{\text{sol}[I]}{=} \delta$ iff for every target instance J such that $(I, J) \models \Sigma$ (that is, every $J \in \text{Sol}(\mathcal{M}, I)$) we have $J \models \delta$.*

We first illustrate on an example the additional source of difficulty brought by the introduction of (immutable) master data.

EXAMPLE 3.5. Consider the following strong schema-mapping $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$ where $\Sigma \in \text{TOC}$ and the following source instance I .

$$\begin{aligned} \mathbf{S} &= \{R, S\} \text{ and } \mathbf{T} = \{A, B, C\} \\ \Sigma &= \left\{ \begin{array}{l} R(x, y) \rightarrow A(x, y) \wedge B(y) \\ A(x, y) \wedge B(y) \rightarrow \exists u, v C(x, u) \wedge A(u, v) \\ A(x, y) \rightarrow S(y) \end{array} \right\} \\ I &= \{R(0, 0); S(0)\} \leftarrow \text{(immutable) master data.} \end{aligned}$$

The set of solutions $\text{Sol}(\mathcal{M}, I)$ is characterized by the following set of tgds and egds

$$\begin{aligned} \Sigma' &\equiv \left\{ \begin{array}{l} \rightarrow A(0, 0) \wedge B(0) \\ A(x, y) \wedge B(y) \rightarrow \exists u, v C(x, u) \wedge A(u, v) \\ A(x, y) \rightarrow y = 0 \end{array} \right\} \\ &\equiv \left\{ \begin{array}{l} \rightarrow A(0, 0) \wedge B(0) \\ A(x, 0) \rightarrow \exists u C(x, u) \wedge A(u, 0) \\ A(u, v) \rightarrow v = 0 \end{array} \right\} \end{aligned}$$

but the set of constraints Σ' ensures no form of termination because the universal models of Σ' are generally infinite. In particular, $\text{UMod}(s(\Sigma'), \emptyset)$ contains no finite instance. \square

Worse still, the following is readily verified:

PROPOSITION 3.8. *The following problem is undecidable: given a strong schema-mapping $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$, $\Sigma \in \text{TOC}$, a source instance I and a target tgd δ , does $\mathcal{M} \stackrel{\text{sol}[I]}{=} \delta$?*

Including the Source Instance. Given a strong schema-mapping $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$ where $\Sigma = \Sigma_{st} \uplus \Sigma_{ts} \uplus \Sigma_s \uplus \Sigma_t$ and given a source instance I for \mathbf{S} , we denote by $\Sigma[I]$ the set of disjunctive tgds and egds over \mathbf{T} such that:

- $\Sigma[I]$ contains Σ_t
- $\Sigma[I]$ contains the set $\Sigma_{st}[I]$ of tgds with an empty body (corresponding to boolean conjunctive queries) of the form $\rightarrow \exists \bar{z} \psi(\bar{c}, \bar{z})$ such that there exists a tgd $\phi(\bar{x}, \bar{y}) \rightarrow \exists \bar{z}, \psi(\bar{x}, \bar{z})$ in Σ_{st} and a tuple of constants $\bar{c} \in \text{Dom}(I)^{|\bar{x}|}$ satisfying $I \models \exists \bar{y}, \phi(\bar{c}, \bar{y})$
- $\Sigma[I]$ contains the set $\Sigma_{ts}[I]$ of disjunctive egds

$$\begin{aligned} \phi(x_1, \dots, x_n, \bar{y}) \rightarrow & (x_1 = c_1^1 \wedge \dots \wedge x_n = c_n^1) \\ & \vee \dots \vee (x_1 = c_1^p \wedge \dots \wedge x_n = c_n^p) \end{aligned}$$

such that there exists a disjunctive tgd

$$\phi(x_1, \dots, x_n, \bar{y}) \rightarrow \bigvee_{i=1}^k (\exists \bar{z}_i \psi_i(x_1, \dots, x_n, \bar{z}_i))$$

in Σ_{ts} such that $\{(c_1^j, \dots, c_n^j) \mid j \leq p\}$ corresponds to the set of tuples (c_1, \dots, c_n) satisfying

$$I \models \bigvee_{i=1}^k (\exists \bar{z}_i \psi_i(c_1, \dots, c_n, \bar{z}_i))$$

PROPOSITION 3.9. *For any strong schema-mapping $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$, any source instance I and any target constraint δ we have $\mathcal{M} \stackrel{\text{sol}[I]}{=} \delta$ iff $\Sigma[I] \models \delta$.*

It follows easily from Theorem 3.2 that $\mathcal{M} \stackrel{\text{sol}[I]}{=} \delta$ becomes decidable under the assumption that $\Sigma[I] \in \text{DTC}$. A very natural question however remains: how to ensure decidability (and reasonable data-complexity) when the source instance I is not known in advance? The key contribution of this section is the following observation: it is possible to test the worst case scenario to ensure not only that $\Sigma[I] \in \text{DTC}$ for every source instance I , but also that $\Sigma[I]$ ensures fairly fast termination, even for a large source instance I . In the following definition we denote by $I_{\mathcal{M}}$ the *critical source instance* of $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$ consisting of all the facts $R(t_1, \dots, t_n)$ such that $R \in \mathbf{S}$ and $\{t_1, \dots, t_n\} \subseteq \text{Cst}(\Sigma) \cup \{*\}$.

DEFINITION 3.4. *Given a strong schema-mapping $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$ we write $\mathcal{M} \in \text{MTOC}$ iff $\Sigma[I_{\mathcal{M}}] \in \text{DTC}$.*

PROPOSITION 3.10. *For any strong schema-mapping $\mathcal{M} \in \text{MTOC}$ there exists a finite k , depending only on \mathcal{M} such that, for every source instance I , every instance J , and every term t occurring in a $P_{s(\Sigma[I])}$ -derivation of J , the size $|t|$ of t is bounded by k .*

THEOREM 3.3. *For any strong schema-mapping $\mathcal{M} \in \text{MTOC}$ the following problem is in Π_2^P : given a source instance I and a target disjunctive tgd δ , does $\mathcal{M} \stackrel{\text{sol}[I]}{=} \delta$? The problem is in coNP if δ is a disjunctive egd or a disjunctive tgd with a bounded number of atoms in each head.*

Matching lower bounds can be obtained in a similar way as for Theorem 3.2.

4. DETERMINACY AND COMPLETENESS

The goal of this section is three-fold: First, we rephrase the notions of *relative information completeness* [8] and *determinacy* [21, 19] in the setting of strong schema-mappings; Second, we introduce interesting variations of these problems; and finally, we establish complexity results for all of these problems, hereby leveraging the results of Section 3.

4.1 Relative Information Completeness

In [8], the problem whether an instance J has complete information to answer a query Q was investigated. More precisely, J is said to be complete for Q if adding tuples to J

does not change the query result $Q(J)$. This problem becomes even more challenging when constraints are in place that limit which tuples can be added to J . In [8], these extensions of J should adhere to certain containment constraints of the form $q(J) \subseteq \pi_X(I)$, where q is a CQ query and $\pi_X(I)$ is a projection of a fixed (immutable) *master data*. As mentioned in the introduction, master data represents an instance which contains complete and consistent information. It is used in this context to provide an “upper bound” for the information extracted by $q(J)$. An instance J is then said to be *complete for Q relative to a set of containment constraints Σ and master data I* , if for any $J' \supseteq J$ such that $J' \models \Sigma$, we have that $Q(J) = Q(J')$. In other words, adding tuples to J either violates the constraints or does not change the query result $Q(J)$.

The following definition generalizes the notion of relative completeness, by incorporating **(i)** more general constraints from J to I ; and **(ii)** allowing constraints from I to J ; and **(iii)** accommodating for tgds and egds on the schema of J .

DEFINITION 4.1. *Given a strong schema-mapping \mathcal{M} , a source instance I , a solution $J \in \text{Sol}(\mathcal{M}, I)$ and a target query Q we say that J is complete for Q relative to (\mathcal{M}, I) if for any other solution $J' \in \text{Sol}(\mathcal{M}, I)$ such that $J \subseteq J'$, we have that $Q(J) = Q(J')$.*

It is easily verified (because of the target dependencies) that this problem is undecidable in general. We obtain decidability, however, when restricting \mathcal{M} to the class MTOC.

THEOREM 4.1. *Let \mathcal{M} be a strong schema-mapping such that $\mathcal{M} \in \text{MTOC}$, instances I and J and target query Q , deciding whether J is complete for Q relative to (\mathcal{M}, I) is in Π_2^P . The problem is in coNP when Q has a bounded number of atoms.*

This generalizes the Π_2^P -completeness result reported in [8]. Indeed, it is easily verified that the containment constraints studied there belong to MTOC. A matching coNP lower bound follows from a reduction from the tautology problem.

We are also interested in the following two stronger notions of completeness:

DEFINITION 4.2. *Given a strong schema-mapping \mathcal{M} , a source instance I and a target query Q we say that (1) Q is complete relative to (\mathcal{M}, I) if any $J \in \text{Sol}(\mathcal{M}, I)$ is complete for Q relative to (\mathcal{M}, I) ; (2) Q is universally complete relative to \mathcal{M} iff Q is complete relative to (\mathcal{M}, I) for every source instance I .*

It is again easily verified that these two problems are undecidable in general. In order to obtain decidability, we proceed as follows: we first reduce the above two problems to an implication problem; this requires the modification of the input schema-mapping \mathcal{M} into a new schema-mapping $\mathcal{M}^{2, \subseteq}$; the decidability then follows by assuming that $\mathcal{M}^{2, \subseteq}$ belongs to one of the classes introduced in Section 3.

Given a strong schema-mapping $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$ we define $\mathcal{M}^{2, \subseteq} = (\mathbf{S}, \mathbf{T}_1 \cup \mathbf{T}_2, \Sigma_1 \cup \Sigma_2 \cup \Sigma_{\subseteq})$, where for $i = 1, 2$:

- \mathbf{T}_i contains a predicate R_i for every $R \in \mathbf{T}$;
- Σ_i contains, for every $r \in \Sigma$, the constraint obtained from r by replacing every atom $R(\bar{t})$ such that $R \in \mathbf{T}$ by the atom $R_i(\bar{t})$.
- Σ_{\subseteq} contains a tgd $R_1(x_1, \dots, x_n) \rightarrow R_2(x_1, \dots, x_n)$ for every $R \in \mathbf{T}$.

Intuitively, \mathbf{T}_1 represents $J \in \text{Sol}(\mathcal{M}, I)$ while \mathbf{T}_2 represents $J' \in \text{Sol}(\mathcal{M}, I)$ such that $J' \supseteq J$. In view of the monotonicity of CQ queries, it suffices to check whether $Q(J') \subseteq Q(J)$. The key observation is that this containment test reduces to checking whether a certain target constraint holds. From the results in Section 3, we then obtain:

THEOREM 4.2. *For every strong schema-mapping \mathcal{M} ,*

1. *If $\mathcal{M}^{2, \subseteq} \in \text{TOC}$ (no disjunction), then deciding whether Q is universally complete relative to \mathcal{M} is in NP. The problem is in PTIME if Q has a bounded number of atoms.*
2. *If $\mathcal{M}^{2, \subseteq} \in \text{DTC}$, then deciding whether Q is universally complete relative to \mathcal{M} is in Π_2^P . The problem is in coNP if Q has a bounded number of atoms.*
3. *If $\mathcal{M}^{2, \subseteq} \in \text{MTOC}$, then deciding whether Q is complete relative to (\mathcal{M}, I) is in Π_2^P . The problem is in coNP if Q has a number of atoms.*

4.2 Determinacy

A similar strategy can be followed for the determinacy problem [21, 19]. Recall that determinacy concerns the question whether a query Q can be answered using a set \mathbf{V} of views. More specifically, one says that \mathbf{V} *determines* Q iff $\mathbf{V}(J) = \mathbf{V}(J')$ implies $Q(J) = Q(J')$, for any two instances J, J' . In other words, determinacy says that \mathbf{V} provides enough information to uniquely determine the answer to Q .

We propose the following generalization of determinacy:

DEFINITION 4.3. *Given a strong schema-mapping \mathcal{M} , a source instance I , and a target query Q we say that (1) Q is determined by (\mathcal{M}, I) iff for all $J_1, J_2 \in \text{Sol}(\mathcal{M}, I)$ we have $Q(J_1) = Q(J_2)$; and (2) Q is determined by \mathcal{M} iff Q is determined in (\mathcal{M}, I') for every source instance I' .*

When the views \mathbf{V} are CQ queries, one can find an \mathcal{M} such that $J \in \text{Sol}(\mathcal{M}, I)$ iff $I = \mathbf{V}(J)$ (and similarly when \mathbf{V} are UCQ queries). Hence, the above definition indeed generalizes the standard notion of determinacy.

It is readily verified that it is undecidable to check whether Q is determined by \mathcal{M} or (\mathcal{M}, I) . Similar to the completeness problem, we can see determinacy as an implication problem. Indeed, for a strong schema-mapping $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$ we let \mathcal{M}^2 be the strong schema-mapping $\mathcal{M}^2 = (\mathbf{S}, \mathbf{T}_1 \cup \mathbf{T}_2, \Sigma_1 \cup \Sigma_2)$ where for $i = 1, 2$, \mathbf{T}_i and Σ_i are defined as in $\mathcal{M}^{2, \subseteq}$. Again, \mathbf{T}_1 represents $J \in \text{Sol}(\mathcal{M}, I)$ while \mathbf{T}_2 represents $J' \in \text{Sol}(\mathcal{M}, I)$. In contrast to the completeness problem, no further constraints between J and J' are enforced. Checking determinacy now pours down to checking whether $Q(J) \subseteq Q(J')$ and $Q(J') \subseteq Q(J)$. That is, one needs to check the implication of two target constraints.

THEOREM 4.3. For every strong schema-mapping \mathcal{M} ,

1. If $\mathcal{M}^2 \in \text{TOC}$ (no disjunction), then deciding whether Q is determined by \mathcal{M} is in NP. The problem is in PTIME if Q has a bounded number of atoms.
2. If $\mathcal{M}^2 \in \text{DTC}$, then deciding whether Q is determined by \mathcal{M} is in Π_2^P . The problem is in coNP if Q has a bounded number of atoms.
3. If $\mathcal{M}^2 \in \text{MTC}$, then deciding whether Q is determined by (\mathcal{M}, I) is in Π_2^P . The problem is in coNP if Q has a bounded number of atoms.

4.3 Example

In this section, we illustrate determinacy and universal completeness. Consider the following strong schema-mapping $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma_{st} \cup \Sigma_t \cup \Sigma_{ts} \cup \Sigma_s)$,

$$\begin{aligned} \mathbf{S} &= \{ \text{Phone}(\text{name}, \text{phone-nb}), \text{Age}(\text{name}, \text{age}) \} \\ \mathbf{T} &= \{ \text{Employee}(\text{name}, \text{phone-nb}, \text{age}) \} \\ \Sigma_{st} &= \{ \text{Phone}(n, p) \rightarrow \exists a, \text{Employee}(n, p, a) \\ &\quad \text{Age}(n, a) \rightarrow \exists p, \text{Employee}(n, p, a) \} \\ \Sigma_t &= \{ \text{Employee}(n, p, a) \wedge \text{Employee}(n, p', a') \rightarrow p = p' \\ &\quad \text{Employee}(n, p, a) \wedge \text{Employee}(n, p', a') \rightarrow a = a' \} \\ \Sigma_{ts} &= \{ \text{Employee}(n, x, y) \rightarrow \\ &\quad (\exists a, \text{Age}(n, a)) \vee (\exists p, \text{Phone}(n, p)) \} \\ \Sigma_s &= \{ \text{Age}(n, a) \wedge \text{Age}(n, a') \rightarrow a = a' \\ &\quad \text{Phone}(n, p) \wedge \text{Phone}(n, p') \rightarrow p = p' \} \end{aligned}$$

We can check that \mathcal{M}^2 and $\mathcal{M}^{2 \subseteq}$ both belong to DTC (by evaluating the corresponding logical programs on the critical instance). Consequently, the algorithms behind Theorems 4.2 and 4.3 allow us to verify that

$$Q_1 = \{(x), \exists p, a, \text{Employee}(x, p, a)\}$$

is determined by \mathcal{M} , while the query

$$Q_2 = \{(p), \exists n, a, \text{Employee}(n, p, a)\}$$

is universally complete relative to \mathcal{M} without being determined by \mathcal{M} .

To be more precise, Q_1 is determined in \mathcal{M} because Q_1 is equivalent to $\{(x), (\exists a, \text{Age}(x, a)) \vee (\exists p, \text{Phone}(x, p))\}$. The query Q_2 is universally complete relative to \mathcal{M} for the following reason: Let $J \in \text{Sol}(\mathcal{M}, I)$ for some I . If we add an atom $\text{Employee}(k, \ell, m)$ in J for some new $\ell \notin Q_2(J)$, then by Σ_{st} either $\text{Age}(k, m)$ or $\text{Phone}(k, \ell)$ must already exist in I . This in turn implies that Employee already contains a tuple $\text{Employee}(k, \ell', m')$. However, Σ_t enforces that $\ell = \ell'$ and $m = m'$. In other words, no new tuple can be added to Employee . The reason why Q_2 is not determined is the following: for the solution $(\{\text{Age}(n, a)\}, \{\text{Employee}(n, \#_0, a)\})$ where $n, a \in \text{Cst}$ and $\#_0 \in \text{Nulls}$ the value of $\#_0$ is not determined. For instance $(\{\text{Age}(n, a)\}, \{\text{Employee}(n, \#_1, a)\})$ is another solution. \square

5. TARGET-TO-SOURCE CONSTRAINTS

So far, we have considered the implication problem in the context of *strong* schema-mappings which (possibly) contain target-to-source dependencies. As discussed in the introduction, enriching standard schema-mappings with target-to-source dependencies requires some care and raises several questions.

In particular, in the context of dynamic data-exchange, and when assuming that a set Σ_{ts} of some target-to-source dependencies are part of the input, it is indeed important to make sure that these target-to-source dependencies, meant to constrain the evolution of the target instances, are satisfied at the initial stage of the exchange, when a first target instance J is computed from the initial source instance I and with a standard schema-mapping management system (SMMS).

More formally, given a strong schema-mapping $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$ where $\Sigma = \Sigma_{st} \cup \Sigma_t \cup \Sigma_{ts}$ and given two instances I and J such that

- I is a source instance of \mathbf{S} ;
- J is a universal solution in $\text{USol}(\mathcal{M}, I)$ computed by some SMMS using only the standard schema-mapping $\mathcal{M}' = (\mathbf{S}, \mathbf{T}, \Sigma_{st}, \Sigma_t)$, without taking Σ_{ts} into account,

we want to make sure that $(I, J) \models \Sigma_{ts}$ will be true, before knowing the source instance I , and without necessarily knowing precisely which SMMS is to be used. In other words, we want to decide whether Σ_{ts} is “implied” by the (standard) schema-mapping embedded in the strong schema-mapping.

Note that here we ignore the source dependencies Σ_s . Indeed, it is assumed that every given source instance I already satisfies Σ_s .

As already mentioned in the introduction, the results of Section 3 are of no use in this case since no target-to-source constraints can be logically implied by a standard schema-mapping. We therefore set out to define alternative semantics of entailment with decidable properties.

5.1 Universal Solutions and Cores

In this section, we first consider the two most popular notions of semantics for schema-mappings based $\text{USol}(\mathcal{M}, I)$ (the set of *universal solutions*) and $\text{Core}(\mathcal{M}, I)$ (the set of *cores of universal solutions*). We refer to Section 2 for the definitions of these notions.

DEFINITION 5.1. Given a schema-mapping \mathcal{M} and a target-to-source dependency δ we write

- $\mathcal{M} \stackrel{\text{USol}}{\models} \delta$ iff for every source instance I and every $J \in \text{USol}(\mathcal{M}, I)$ we have $(I, J) \models \delta$, and
- $\mathcal{M} \stackrel{\text{Core}}{\models} \delta$ iff for every source instance I and every $J \in \text{Core}(\mathcal{M}, I)$ we have $(I, J) \models \delta$.

In the light of the following example (similar to examples used in [16, 13]), it becomes clear that the notion of universal solution is not restrictive enough because (i) some universal solutions are neither natural nor realistic; and (ii) USol leads to a notion of entailment $\stackrel{\text{USol}}{\models}$ which is not general enough for target-to-source dependencies (intuitively, $\mathcal{M} \stackrel{\text{USol}}{\models} \delta$ only holds for a trivial schema-mappings \mathcal{M} or a trivial target-to-source δ).

¹Note that for the sake of simplicity (and tractability) we ignore here the set of source dependencies Σ_s that may come with the source schema.

EXAMPLE 5.1. Consider $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma_{st}, \Sigma_t)$ with $\mathbf{S} = \{A\}$, $\mathbf{T} = \{B\}$, $\Sigma_t = \emptyset$ and where $\Sigma_{st} = \{A(x) \rightarrow B(x)\}$ copies A to B . Let δ be the source-to-target dependency $B(x) \rightarrow A(x)$. The “natural” solution for \mathcal{M} and $I_n = \{A(1), \dots, A(n)\}$ would be the target instance $J_n = \{B(1), \dots, B(n)\}$ which satisfies $(I_n, J_n) \models \delta$. We can observe, however, that the target instance $J'_n = \{B(1), \dots, B(n), B(\#_0)\}$ where $\#_0 \in \text{Nulls}$ is also a universal solution. Since $(I_n, J'_n) \not\models \delta$, we have $\mathcal{M} \not\models^{\text{core}} \delta$. In contrast, $\text{Core}(\mathcal{M}, I_n)$ captures precisely the desired semantics since $\text{Core}(\mathcal{M}, I_n) = \{J_n\}$ and we have $\mathcal{M} \models^{\text{core}} \delta$. \square

It can easily be argued that \models^{core} is the most natural semantics for schema-mappings and also the most relevant semantics in the context of target-to-source dependencies. However, the following proposition shows that general and unrestricted use of \models^{core} easily leads to undecidability. This is even the case for schema-mappings that ensures oblivious termination ($\Sigma_{st} \cup \Sigma_t \in \text{TOC}$), are without egds, without disjunction in δ , and under the assumption that Σ_t satisfies the syntactic criteria of *weak acyclicity* mentioned in the introduction and further discussed in Section 6.

PROPOSITION 5.1. *The following problem is undecidable: given a schema-mapping $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma_{st}, \Sigma_t)$ where Σ_t consists only of a weakly acyclic set of target tgds and given a target-to-source tgd δ , does $\mathcal{M} \models^{\text{core}} \delta$?*

5.2 Core-Ground and Core-Safe Queries

In this section, we identify more precisely the origin of the undecidability stated in Proposition 5.1. In particular, we show that decidability can be obtained by requiring δ to satisfy either a very general and high-level property (δ is *core-ground*) or a less general but tractable property (δ is *core-safe*).

Core-ground queries. Given a schema-mapping $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma_{st}, \Sigma_t)$ and a conjunctive query Q over \mathbf{T} we say Q is core-ground (in \mathcal{M}) iff for every source instance I and every $J \in \text{Core}(\mathcal{M}, I)$ there is no null in $Q(J)$ (we say that $Q(J)$ is ground). Given a target-to-source disjunctive tgd δ of the form $\phi(\bar{x}, \bar{y}) \rightarrow \bigvee_i H_i$, we then say that δ is core-ground iff the conjunctive query $Q_\phi = \{\bar{x} \mid \exists \bar{y}, \phi(\bar{x}, \bar{y})\}$ is core-ground.

PROPOSITION 5.2. *The following problem is undecidable: given a schema-mapping $\mathcal{M} \in \text{TOC}$ and a target query of the form $Q = \{\{\bar{x}\} \mid R(x)\}$ for some unary predicate R , is Q core-ground in \mathcal{M} ?*

THEOREM 5.1. *The following problem is decidable: given a schema-mapping $\mathcal{M} \in \text{TOC}$ and a given a target-to-source disjunctive tgds δ such that δ is core-ground in \mathcal{M} , does $\mathcal{M} \models^{\text{core}} \delta$?*

We next identify a decidable syntactic criterion ensuring that a target query is core-ground in a schema-mapping.

Core-safe Queries. Let $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma_{st}, \Sigma_t)$ be a schema-mapping. We define a *position* as a pair (R, i) preferably written R^i where $R \in \mathbf{T}$ and $i \in \{1, \dots, a_R\}$. Given a set of atoms ϕ over \mathbf{T} and a variable x we denote by $\mathcal{P}_{\phi, x}$ the set of positions R^i such that ϕ contains an atom $R(t_1, \dots, t_n)$

where $t_i = x$. We say that a position p of \mathbf{T} is *initially affected* in \mathcal{M} iff there is a tgd $\phi(\bar{x}, \bar{y}) \rightarrow \exists \bar{z}, \psi(\bar{x}, \bar{z})$ in $\Sigma_{st} \cup \Sigma_t$ such that $p \in \bigcup_{z \in \bar{z}} \mathcal{P}_{\psi, z}$. We let $\text{IA}(\mathcal{M}) = \{\{p_1\}, \dots, \{p_m\}\}$ where $\{p_1, \dots, p_m\}$ is the set of initially affected positions.

Let $P = \{P_1, P_2, \dots, P_m\}$ be a finite set where each P_i is a set of positions. Given a target tgd $\phi(\bar{x}, \bar{y}) \rightarrow \exists \bar{z}, \psi(\bar{x}, \bar{z})$ in Σ and given a universal quantified variable $x \in \bar{x}$ such that $\mathcal{P}_{\phi, x} \subseteq P_1$ we write $P \sim_{\mathcal{M}}^{\text{tgd}} \{P_1 \cup \mathcal{P}_{\psi, x}, P_2, \dots, P_m\}$. Given a target egd of the form $\phi(x_1, x_2, \bar{y}) \rightarrow x_1 = x_2$ in Σ such that $\mathcal{P}_{\phi, x_1} \subseteq P_1$ and $\mathcal{P}_{\phi, x_2} \subseteq P_2$ we write $P \sim_{\mathcal{M}}^{\text{egd}} \{P_1 \cup P_2, P_3, \dots, P_m\}$. We finally denote by $\sim_{\mathcal{M}}$ the reflexive and transitive closure of the binary relation $\sim_{\mathcal{M}}^{\text{tgd}} \cup \sim_{\mathcal{M}}^{\text{egd}}$.

Given a conjunctive query $Q = \{\bar{x} \mid \exists \bar{y}, \phi(\bar{x}, \bar{y})\}$ over \mathbf{T} , we say that Q is unsafe in \mathcal{M} iff there exists an $x \in \bar{x}$ and a set $\{P_1, \dots, P_m\}$ such that (i) $\text{IA}(\mathcal{M}) \sim_{\mathcal{M}}^* \{P_1, \dots, P_m\}$; and (ii) $\mathcal{P}_{\phi, x} \subseteq P_1$. We say that Q is *core-safe* in \mathcal{M} iff Q is not unsafe in \mathcal{M} . We say that a disjunctive tgd $\phi(\bar{x}, \bar{y}) \rightarrow \bigvee_i \exists \bar{z}_i, \psi_i(\bar{x}, \bar{z}_i)$ is core-safe in \mathcal{M} iff Q_ϕ is core-safe in \mathcal{M} .

EXAMPLE 5.2. Consider the two schema mappings $\mathcal{M}_a = (\mathbf{S}, \mathbf{T}, \{r_1; r_2; r_3\})$ and $\mathcal{M}_b = (\mathbf{S}, \mathbf{T}, \{r_1; r_2; r_3; e_1\})$, with $\mathbf{S} = \{A; B\}$, $\mathbf{T} = \{R; S\}$ and constraints:

$$\begin{aligned} r_1 &= A(x) \rightarrow \exists y, z, R(x, y, z) \\ r_2 &= B(x, y, z) \rightarrow R(x, y, z) \\ r_3 &= R(x, y, y) \rightarrow S(x, y) \\ e_1 &= R(x, y, z) \rightarrow y = z \end{aligned}$$

Let $Q = \{x, y \mid S(x, y)\}$. Then, $\text{IA}(\mathcal{M}_a) = \{\{R^1\}; \{R^2\}\}$ which is a fixed point for $\sim_{\mathcal{M}}$ and therefore Q is core-safe in \mathcal{M}_a . In contrast, Q is not core-safe in \mathcal{M}_b . Indeed,

$$\text{IA}(\mathcal{M}_b) = \{\{R^1\}; \{R^2\}\} \sim_{\mathcal{M}}^{\text{egd}} \{\{R^1; R^2\}\} \sim_{\mathcal{M}}^{\text{tgd}} \{\{R^1; R^2; S^2\}\}$$

but $\mathcal{P}_{Q, y} = \{S^2\}$ is included in $\{R^1; R^2; S^2\}$. \square

PROPOSITION 5.3.

- For all schema-mappings \mathcal{M} and all target-to-source disjunctive tgds δ , if δ is core-safe in \mathcal{M} , then δ is core-ground in \mathcal{M} (and Theorem 5.1 is applicable).
- Moreover, the following problem is in PTIME: given \mathcal{M} and δ as above, is δ core-safe in \mathcal{M} ?

5.3 Closed-World Assumption

In this section, we propose another semantics which is more flexible than \models^{core} and captures an intuitive notion of *closed-world semantics*, in a style similar to that of [15, 13]. Furthermore, under this semantics we obtain decidability for arbitrary target-to-source dependencies (in contrast to the safety assumption in the previous section).

Canonical Instance. Let $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma_{st}, \Sigma_t)$ be a schema-mapping. In the case of tgds only (without target egds) and given a source instance I we define the *canonical solution* of I as the target instance $\Pi_{\mathbf{T}}(P_{\Sigma}(I))$ where P_{Σ} is the fixed point of I by the Skolemisation P_{Σ} of I and $\Pi_{\mathbf{T}}$ is the operation selecting the atoms $R(\bar{i})$ such that $R \in \mathbf{T}$. From Proposition 3.3 we know that the canonical solution is indeed a solution (i.e., it belongs to $\text{USol}(\mathcal{M}, I)$).

In the presence of egds, we use the substitution-free simulation defined in Section 3.1 to define a *canonical instance* $\mathcal{M}^{\text{can}}(I)$ of I . That is, we define $\mathcal{M}^{\text{can}}(I) = \Pi_{\mathbf{T}}(P_{s(\Sigma)}(I)^+)$ where: $s(\Sigma)$ is the substitution-free simulation of Σ ; $^+$ is the operation that applies the substitution axioms; and $\Pi_{\mathbf{T}}$ selects the facts that are over the target schema (dropping in particular the facts $E(t, t')$). It can be seen that $\mathcal{M}^{\text{can}}(I)$ coincides precisely with the canonical solution of I when there is no egd in Σ . We use, however, the expression *canonical instance* because in the presence of egds, $\mathcal{M}^{\text{can}}(I)$ is not necessarily a solution. This is illustrated in the following example.

EXAMPLE 5.3. Consider the schema-mapping \mathcal{M} where $\mathbf{S} = \{A, B\}$, $\mathbf{T} = \{R, S, T\}$ and (Σ_{st}, Σ_t) is as follows.

$$\Sigma_{st} = \left\{ \begin{array}{l} A(x) \rightarrow \exists u, v R(x, u) \wedge R(v, x) \wedge S(u, v) \\ B(x) \rightarrow \exists R(x, x) \end{array} \right\}$$

$$\Sigma_t = \left\{ \begin{array}{l} S(x, y) \rightarrow x = y \\ R(x, x) \rightarrow T(x) \end{array} \right\}$$

In this setting, assuming that $P_{s(\Sigma)}$ is equal to

$$\left\{ \begin{array}{l} A(x) \rightarrow R(x, f_u(x)) \wedge R(f_v(x), x) \wedge S(f_u(x), f_v(x)) \\ B(x) \rightarrow R(x, x) \\ S(x, y) \rightarrow E(x, y) \\ R(x, x') \wedge E(x, y') \rightarrow T(x) \\ (+\text{axioms of } E) \end{array} \right\}$$

we observe that for $I = \{A(0), B(1)\}$ the canonical instance

$$\mathcal{M}^{\text{can}}(I) = \{ R(0, f_u(0)); R(f_v(0), 0); S(f_u(0), f_v(0)); \\ R(0, f_v(0)); R(f_u(0), 0); S(f_v(0), f_u(0)); \\ R(1, 1); T(1) \}$$

violates Σ_t and is not a solution. \square

From Theorem 3.4 in [15], it follows that the definition below coincides with the notion of CWA-solutions introduced in [15] when only source-to-target tgds are present.

DEFINITION 5.2. *Given a schema-mapping \mathcal{M} and a source instance I , we say that a universal solution $J \in \text{USol}(\mathcal{M}, I)$ is CwA iff there exists a homomorphism h from $\mathcal{M}^{\text{can}}(I)$ to J such that $J = h(\mathcal{M}^{\text{can}}(I))$. We denote by $\text{CwA}(\mathcal{M}, I)$ the set of CwA universal solutions in $\text{USol}(\mathcal{M}, I)$ and, given a dependency δ we write $\mathcal{M} \stackrel{\text{cwa}}{\models} \delta$ iff for every source instance I and every $J \in \text{CwA}(\mathcal{M}, I)$ we have $(I, J) \stackrel{\text{cwa}}{\models} \delta$.*

We can first observe that $\stackrel{\text{cwa}}{\models}$ is a good candidate for decidability in the sense that it is more restrictive than $\stackrel{\text{usol}}{\models}$ while being less restrictive than $\stackrel{\text{core}}{\models}$.

PROPOSITION 5.4. $\stackrel{\text{usol}}{\models} \subsetneq \stackrel{\text{cwa}}{\models} \subsetneq \stackrel{\text{core}}{\models}$, that is, for all schema-mappings \mathcal{M} and all source instances I we have

$$\text{Core}(\mathcal{M}, I) \subseteq \text{CwA}(\mathcal{M}, I) \subseteq \text{USol}(\mathcal{M}, I)$$

and both inclusions can be strict.

We are now ready to prove our final theorem.

THEOREM 5.2. *The following problem is decidable: given a schema-mapping $\mathcal{M} \in \text{TOC}$ and given a target-to-source disjunctive tgd δ , does $\mathcal{M} \stackrel{\text{cwa}}{\models} \delta$?*

The following example shows that Theorem 5.2 can be used in scenarios that are not covered by the previous section.

EXAMPLE 5.4. Consider the schema-mapping \mathcal{M} from Example 5.3 and consider the target-to-source tgd

$$\delta : T(x) \rightarrow B(x)$$

It is clear that $\mathcal{M} \in \text{TOC}$ but it can be seen that the query $\{x | T(x)\}$ is not core-safe (and Theorem 5.3 does not apply). We can nonetheless use the algorithm corresponding to Theorem 5.2 to check that we have $\mathcal{M} \stackrel{\text{cwa}}{\models} \delta$ (and therefore $\mathcal{M} \stackrel{\text{core}}{\models} \delta$). \square

6. RELATED WORK

Finite Models. The different semantics considered in this paper rely on *infinite models*. In particular, for a strong schema-mapping $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$, we write $\Sigma \models \delta$ whenever, for every (possibly infinite) instances I and J over the source schema and target schema, respectively, we have $(I, J) \models \Sigma$ implies $(I, J) \models \delta$. As shown in [4], in the case of strong schema-mappings this notion of entailment is not equivalent to the one where only *finite* instances are considered. The reason why we did not discuss this distinction in the previous sections is the following observation (similar to Proposition 5 in [17]): the assumptions of termination considered in the Theorems of this paper all ensure both tractability and *finite controllability* (meaning that finite-model semantics and infinite-model semantics coincide). For instance, our notion of determinacy has been shown in [19] to differ from *finite* determinacy (for every *finite* solutions $J_1, J_2 \in \text{Sol}(\mathcal{M}, I)$ we have $Q(J_1) = Q(J_2)$). However, if for a given \mathcal{M} , $\mathcal{M}^2 \in \text{DTC}$, then a target query Q is determined by \mathcal{M} iff it is finitely determined by \mathcal{M} .

Weak Acyclicity. It has been shown in [17] that the class TOC is a strict generalisation of the widely-used notion of *weak acyclicity* [6, 5], even in the presence of arbitrary egds. It can be seen that MTOC is similarly a strict generalization of *weak acyclicity*. More precisely, if we consider a strong schema-mapping $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$ where Σ consists of tgds and egds only (no disjunction), then we can check that the following holds: if Σ_t is the union of a weakly-acyclic set of tgds and an arbitrary set of egds, then $\mathcal{M} \in \text{MTOC}$. Note that MTOC is at the same time a different kind of criterion than *weak-acyclicity*. Indeed, while *weak-acyclicity* of schema-mappings can be decided in PTIME, membership in MTOC is in general undecidable. Recall, however, that MTOC (just like TOC and DTC) can be tested on a critical instance which is typically of small size and provided that “efficient” schema-mappings are given, the MTOC class appears feasible in practice. Finally, tractable subclasses of MTOC (and DTC) can be defined in a similar way as in [17].

Peer Data Exchange. Target-to-source dependencies have been considered in [9]. There, a specific class of strong schema-mappings called Peer-Data-Exchange settings (PDE) was introduced. In particular, a PDE setting is a strong schema-mapping $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$ with no source constraints and no disjunction (where $\Sigma = \Sigma_{st} \cup \Sigma_{ts} \cup \Sigma_t$). Static analysis of the PDE setting was not the focus of [9]. Instead,

a very interesting decision problem, denoted $\text{SOL}(\mathcal{M})$ was considered: given a source instance I and a target instance J for a fixed PDE setting $\mathcal{M} = (\mathbf{S}, \mathbf{T}, \Sigma)$, is there a target instance J^+ (called a *solution* in [9]) such that $J \subseteq J^+$ and $(I, J^+) \models \Sigma$? This problem $\text{SOL}(\mathcal{M})$ has been proven in [9] to be NP-complete under the assumption of weak acyclicity, while it is clearly undecidable in general. Our results in Section 3 can be used to strengthen this result: $\text{SOL}(\mathcal{M})$ is in NP whenever $\mathcal{M} \in \text{MTOC}$. Note that other results (such as Theorem 2 in [9]) obtained under the assumption of weak acyclicity can similarly be generalized to the class MTOC.

Restricted Chase. The assumptions of termination in this paper rely on the oblivious chase, which differs from the well-known procedure called *the restricted chase* (often simply referred to as *the chase*). With the restricted chase, an inconsistency is resolved (intuitively, by adding the head of a tgds) only when necessary. We can observe that the restricted chase terminates whenever the oblivious chase terminates (while the converse does not hold). It can be seen however that the restricted chase (and similarly, the *core chase* introduced in [4]) is problematic for several reasons: (i) termination cannot be tested on the critical source instance, (ii) applying a rule is more costly than applying a rule of the oblivious chase, and (iii) many questions remain open (e.g. does termination of the restricted chase implies *polynomial* termination?). Note however that several decidable criteria such as the very involved *inductive restriction* in [18] (which is incomparable with TOC) have been designed to ensure termination of (only) the restricted chase. We intend to show in the future that our results can be extended to this class by using a rewriting algorithm, that given an inductively restricted set Σ of tgds and egds, returns a set of tgds and egds $\Sigma' \in \text{TOC}$.

Guarded Tgds. Another setting in which entailment of tgds and egds has been shown decidable (with a good complexity) is when all the tgds are *guarded* (and the egds are *separable*) [2]. A tgd is called guarded iff there is an atom in the body covering all the universal variables. For instance, the following tgd δ_1 is guarded (by the atom $A(x, y)$) while the tgd δ_2 is not guarded.

$$\begin{aligned} \delta_1 : & \quad A(x, y) \wedge C(y) \rightarrow R(x, y) \\ \delta_2 : & \quad A(x, y) \wedge B(y, z) \rightarrow R(x, z) \end{aligned}$$

The main advantage of this setting is that it neither requires the existence of a *finite* universal model (in contrast with [4]), nor any assumption of termination. It is also interesting to observe that the target tgds used in real-life schema-mappings usually consist of inclusion dependencies, that is, of tgds which have only one atom in the body and are therefore guarded. One limit of this framework, however, is the fact that it does not apply to strong schema-mappings with non-guarded source-to-target tgds and target-to-source tgds. To be more precise, as soon as Σ_{st} contains a single non-guarded tgd such as δ_2 above (used to materialize the join of A and B), the results of [2] cannot be applied directly.

Hypertree-width. We have shown that several complexity results can be significantly improved when assuming a fixed bound on the number of atoms, either in the head H of

the constraints δ considered in Section 3 (Theorems 3.1, 3.2 and 3.3), or in the conjunctive queries Q considered in Section 4 (Theorems 4.1, 4.2 and 4.3). It can be seen that the same complexity results hold under the much more general assumption of *bounded hypertree-width* [10] which strictly covers by itself the cases where H (or Q) has a bounded number of existential variables, is *acyclic* or has a bounded *tree-width* [10].

7. CONCLUSION

We have studied the implication problem for strong schema-mappings which are extensions of standard schema-mappings with (disjunctive) target-to-source and source dependencies. As an application, we revisited the problems of relative information completeness and determinacy, both instances of implication problems. We also addressed the question whether target-to-source constraints are already implied by the embedded standard schema-mapping, a problem that cannot be solved using the standard logical implication techniques mentioned above.

Interesting open problems concern the precise complexity of the implication problem of target-to-source constraints under both the core- and CwA-implication semantics. Other questions involve the identification of tractable subclasses of TOC, DTOC and MTOC (similar to those in [17]) and the generalisation to classes of dependencies ensuring (only) the termination of the restricted chase.

Acknowledgment. The first author was funded by EPSRC EP/E010865/1 and partially supported by ERC Advanced Grant Webdam on Foundation of Web data management. The second author was supported by EPSRC EP/E029213/1. We thank Serge Abiteboul and Alin Deutsch for helpful discussions.

8. REFERENCES

- [1] S. Abiteboul, R. Hull, and V. Vianu. *Foundations of Databases*. 1995.
- [2] A. Cali, G. Gottlob, and T. Lukasiewicz. A general datalog-based framework for tractable query answering over ontologies. In *PODS*, pages 77–86, 2009.
- [3] S. Chaudhuri and M. Y. Vardi. On the equivalence of recursive and nonrecursive datalog programs. *J. Comput. Syst. Sci.*, 54(1):61–78, 1997.
- [4] A. Deutsch, A. Nash, and J. B. Remmel. The chase revisited. In *PODS*, pages 149–158, 2008.
- [5] A. Deutsch and V. Tannen. Xml queries and constraints, containment and reformulation. *Theor. Comput. Sci.*, 336(1):57–87, 2005.
- [6] R. Fagin, P. G. Kolaitis, R. J. Miller, and L. Popa. Data exchange: semantics and query answering. *Theor. Comput. Sci.*, 336(1):89–124, 2005.
- [7] R. Fagin, P. G. Kolaitis, and L. Popa. Data exchange: getting to the core. *ACM Trans. Database Syst.*, 30(1):174–210, 2005.
- [8] W. Fan and F. Geerts. Relative information completeness. In *PODS*, pages 97–106, 2009.

- [9] A. Fuxman, P. G. Kolaitis, R. J. Miller, and W. C. Tan. Peer data exchange. *ACM Trans. Database Syst.*, 31(4):1454–1498, 2006.
- [10] G. Gottlob, N. Leone, and F. Scarcello. Hypertree decompositions and tractable queries. *J. Comput. Syst. Sci.*, 64(3):579–627, 2002.
- [11] G. Gottlob and A. Nash. Efficient core computation in data exchange. *J. ACM*, 55(2), 2008.
- [12] L. M. Haas, M. A. Hernández, H. Ho, L. Popa, and M. Roth. Clio grows up: from research prototype to industrial tool. In *SIGMOD*, pages 805–810, 2005.
- [13] A. Hernich and N. Schweikardt. Cwa-solutions for data exchange settings with target dependencies. In *PODS*, pages 113–122, 2007.
- [14] P. G. Kolaitis. Schema mappings, data exchange, and metadata management. In *PODS*, pages 61–75, 2005.
- [15] L. Libkin. Data exchange and incomplete information. In *PODS*, pages 60–69, 2006.
- [16] L. Libkin and C. Sirangelo. Data exchange and schema mappings in open and closed worlds. In *PODS*, pages 139–148, 2008.
- [17] B. Marnette. Generalized schema-mappings: from termination to tractability. In *PODS*, pages 13–22, 2009.
- [18] M. Meier, M. Schmidt, and G. Lausen. On chase termination beyond stratification. *PVLDB*, 2(1):970–981, 2009.
- [19] A. Nash, L. Segoufin, and V. Vianu. Determinacy and rewriting of conjunctive queries using views: A progress report. In *ICDT*, pages 59–73, 2007.
- [20] Y. Sagiv and M. Yannakakis. Equivalences among relational expressions with the union and difference operators. *J. ACM*, 27(4):633–655, 1980.
- [21] L. Segoufin and V. Vianu. Views and queries: determinacy and rewriting. In *PODS*, pages 49–60, 2005.
- [22] O. Shmueli. Equivalence of datalog queries is undecidable. *J. Log. Program.*, 15(3):231–241, 1993.
- [23] B. ten Cate, L. Chiticariu, P. G. Kolaitis, and W. C. Tan. Laconic schema mappings: Computing the core with sql queries. *PVLDB*, 2(1):1006–1017, 2009.