

Normalization and optimization of schema mappings

Georg Gottlob · Reinhard Pichler · Vadim Savenkov

Received: 23 August 2010 / Accepted: 8 February 2011 / Published online: 16 March 2011
© Springer-Verlag 2011

Abstract Schema mappings are high-level specifications that describe the relationship between database schemas. They are an important tool in several areas of database research, notably in data integration and data exchange. However, a concrete theory of schema mapping optimization including the formulation of optimality criteria and the construction of algorithms for computing optimal schema mappings is completely lacking to date. The goal of this work is to fill this gap. We start by presenting a system of rewrite rules to minimize sets of source-to-target tuple-generating dependencies. Moreover, we show that the result of this minimization is unique up to variable renaming. Hence, our optimization also yields a schema mapping normalization. By appropriately extending our rewrite rule system, we also provide a normalization of schema mappings containing equality-generating target dependencies. An important application of such a normalization is in the area of defining a set Σ of dependencies. These dependencies express conditions that instances S and T must fulfill. In data exchange, several definitions in this area depend on the concrete syntactic representation of the mappings. This is, in particular, the case for queries with negated atoms and for aggregate queries. The normalization of schema mappings allows us to eliminate the effect of the concrete syntactic representation of the mapping from the semantics of query answering. We discuss in detail how our results can be fruitfully applied to aggregate queries.

Keywords Data integration Data exchange Schema mappings optimization

Schema mappings are high-level specifications that describe the relationship between two database schemas. They play an important role in data integration [4, 19] and data exchange [9]. A schema mapping is usually given in the form $M = (S, T, \Sigma)$, indicating the two database schemas S and T plus a set Σ of dependencies. These dependencies express conditions that instances S and T must fulfill. In data exchange, several definitions in this area depend on the concrete syntactic representation of the mappings. This is, in particular, the case for queries with negated atoms and for aggregate queries. The normalization of schema mappings allows us to eliminate the effect of the concrete syntactic representation of the mapping from the semantics of query answering. We discuss in detail how our results can be fruitfully applied to aggregate queries.

G. Gottlob
Computing Laboratory, Oxford University, Room 358,
Wolfson Building, Parks Road, Oxford OX1 3QD,
United Kingdom
e-mail: georg.gottlob@comlab.ox.ac.uk

R. Pichler · V. Savenkov (✉)
Database and Artificial Intelligence Group,
Institute of Information Systems,
Vienna University of Technology, Favoritenstrasse 9,
1040 Vienna, Austria
e-mail: savenkov@dbai.tuwien.ac.at

R. Pichler
e-mail: pichler@dbai.tuwien.ac.at

Over the past years, schema mappings have been extensively studied (see [6, 18] for numerous pointers to the literature). However, only recently, the questions of schema mapping optimization has been raised. In [10], the foundation for optimization has been laid by defining various forms of equivalence of schema mappings and by proving important properties of the resulting notions. However, a concrete theory of schema mapping optimization including the formulation of optimality criteria and the construction of algorithms for computing optimal schema mappings is completely lacking to date. The goal of this work is to fill this gap. Below,

we illustrate the basic ideas of our approach by a series of simple examples, where it is clear “at a glance” what the optimal form of the schema mappings should look like. In fact, one would expect that a human user designs these mappings as mentioned above, GAV mappings are only a special case of schema mappings given by s-t tgds which, in the general case, may have existentially quantified variables and conjunctions of atoms in the conclusion. Note that the existentially quantified variables are used to represent incomplete data (in the form of marked nulls) in the target instance. Hence, as an additional optimization goal, we would like to minimize the number of existentially quantified variables in each s-t tgd. Moreover, we would now also like to minimize the number of atoms in the CQ of the conclusion.

As mentioned above, GAV mappings are only a special case of schema mappings given by s-t tgds which, in the general case, may have existentially quantified variables and conjunctions of atoms in the conclusion. Note that the existentially quantified variables are used to represent incomplete data (in the form of marked nulls) in the target instance. Hence, as an additional optimization goal, we would like to minimize the number of existentially quantified variables in each s-t tgd. Moreover, we would now also like to minimize the number of atoms in the CQ of the conclusion.

For the most common form of schema mappings considered in the literature, the dependencies are source-to-target tuple-generating dependencies (or s-t tgds, for short) of the form $x(\varphi(x) \rightarrow y \psi(x, y))$, where the antecedent φ is a conjunctive query (CQ) over S and the conclusion ψ is a CQ over T . The universal quantification is usually not denoted explicitly. Instead, it is assumed implicitly for all variables in $\varphi(x)$.

Example 1 Consider a schema mapping $M = \langle S, T, \Sigma \rangle$ with $S = \{L(\cdot, \cdot, \cdot), P(\cdot, \cdot)\}$ and $T = \{C(\cdot, \cdot)\}$, where L, P , and C are abbreviations for the relational schema L (lecture(title, year, prof), Prof(name, area), and Course(title, prof-area), respectively. Moreover, suppose that S consists of two rules expressing the following constraints: If any lecture is specified in the source instance, then the title of all lectures for 3rd year students as well as the area of the professor giving this lecture should be present in the Course relation of the target instance. Moreover, T contains a specific rule which takes care of the lectures given by professors from the database area. We get the following set of s-t tgds:

$$\begin{aligned} L(x_1, x_2, x_3) \wedge L(x_4, 3, x_5) \wedge P(x_5, x_6) \wedge C(x_4, x_6) \\ L(x_1, 3, x_2) \wedge P(x_2, db) \wedge C(x_1, db) \end{aligned}$$

The above schema mapping has a specific form called GAV (global-as-view) [19], i.e., we only have s-t tgds of the form $A(x) \rightarrow B(x)$, where the conclusion is a single atom $B(x)$ without existentially quantified variables. In this special case, we see a close relationship of schema mappings with unions of conjunctive queries (UCQs). Indeed, given a source instance over S , the tuples that have to be present in any legal target instance J according to the above schema mapping are precisely the tuples in the result of the following UCQ:

$$\begin{aligned} \text{ans}(x_4, x_6) : \check{S} L(x_1, x_2, x_3) \wedge L(x_4, 3, x_5) \wedge P(x_5, x_6) \\ \text{ans}(x_1, db) : \check{S} L(x_1, 3, x_2) \wedge P(x_2, db) \end{aligned}$$

The goal of UCQ optimization is usually twofold [17, 25], namely to minimize the number of CQs and to minimize the number of atoms in each CQ. In the above UCQ, we would thus delete the second CQ and, moreover, eliminate the first atom from the body of the first CQ. In total, the above UCQ can be replaced by a single CQ $\text{ans}(x_4, x_6) :-$

$$\Sigma = \{ L(x_4, 3, x_5) \wedge P(x_5, x_6) \rightarrow C(x_4, x_6) \}.$$

Example 2 We revisit Example 1, and consider a new mapping M in the reverse direction so to speak: $M = \langle S, T, \Sigma \rangle$ with $S = \{C(\cdot, \cdot)\}$ and $T = \{L(\cdot, \cdot, \cdot), P(\cdot, \cdot)\}$ where L, P , and C are as before. Moreover, let Σ be defined as follows:

$$\begin{aligned} \Sigma = \{ C(x_1, x_2) \wedge (y_1, y_2, y_3, y_4) \rightarrow L(y_1, y_2, y_3) \\ L(x_1, 3, y_4) \wedge P(y_4, x_2), \\ C(x_1, db) \wedge (y_1) \rightarrow L(x_1, 3, y_1) \wedge P(y_1, db) \} \end{aligned}$$

Clearly, Σ is equivalent to the singleton $\Sigma = \{ C(x_1, x_2) \wedge (y_4) \rightarrow L(x_1, 3, y_4) \wedge P(y_4, x_2) \}.$

The above schema mapping corresponds to the special case of LAV (local-as-view) [19] with s-t tgds of the form $A(x) \rightarrow \psi(x, y)$, where the antecedent is a single atom $A(x)$ and all variables in $A(x)$ actually do occur in the conclusion. In the most general case (referred to as GLAV mappings), no restrictions are imposed on the CQs in the antecedent and conclusion nor on the variable occurrences. In order to formulate an optimality criterion for schema mappings with s-t tgds of this general form, the analogy with UCQs does not suffice. Indeed, the following example illustrates that we may get a highly unsatisfactory result if we just aim at the minimization of the number of s-t tgds and of the number of atoms inside each s-t tgd.

Example 3 Let $M = \langle S, T, \Sigma \rangle$ with $S = \{L(\cdot, \cdot, \cdot)\}$ and $T = \{C(\cdot, \cdot), E(\cdot, \cdot)\}$ where L and C are as before and E denotes the schema E (Equal-Year(course1, course2), i.e., E contains pairs of courses designed for students in the same year. Moreover, let Σ be defined as follows:

$$\Sigma = \{ L(x_1, x_2, x_3) \wedge (y) \rightarrow C(x_1, y), \\ L(x_1, x_2, x_3) \wedge L(x_4, x_2, x_5) \rightarrow E(x_1, x_4) \}$$

Then, Σ is equivalent to the singleton Σ with the tgd $L(x_1, x_2, x_3) \wedge L(x_4, x_2, x_5) \wedge (y) \rightarrow C(x_1, y) \wedge E(x_1, x_4)$

Now suppose that the title attribute is a key in the source instance. Let l_i denote the title of some lecture in a source instance and suppose that I contains m lectures for students in the same year as l_i . Then, the computation of the canonical universal solution (for details, see Sect. 2) yields two results of significantly different quality depending on whether we take Σ : In case of Σ , we get one tuple $C(l_i, y)$ with this course title l_i . In contrast, for Σ' , we get m tuples $C(l_i, y_1), \dots, C(l_i, y_m)$ with the same course title. The reason for this is that the s-t tgd “res” for every possible combination of key values x_1 and x_4 , although for the conjunct $C(x_1, y)$ in the conclusion, only the value of x_1 is relevant.

We shall refer to the two s-t tgds σ of the above example as the split form of the s-t tgd τ . We shall formally define splitting of s-t tgds in Sect. 3. Intuitively, splitting aims at breaking up the conclusion of an s-t tgd in smaller parts such that the variables in the antecedent are indeed related to the atoms in the conclusion. Without this measure, any target instance would be artificially inflated with marked nulls as we have seen with τ in the above example. Splitting helps to avoid such anomalies. Indeed, it can be seen as an analogous operation to the decomposition of relational schemas into normal form where we also want to exclude that some attributes are fully determined by parts of a key. Carrying over this idea to s-t tgds, we want to exclude that some atoms in the conclusion are fully determined by parts of the atoms in the antecedent. Our optimization goal for schema mappings will therefore be to minimize the number of s-t tgds only to the extent that splitting should be applied whenever possible. Minimizing the size of each s-t tgd and the number of existentially quantified variables in the conclusion will, of course, be pursued as another optimization goal. We thus have the following optimality criteria for sets of s-t tgds:

- cardinality-minimality i.e., the number of s-t tgds in Σ shall be minimal;
- antecedent-minimality i.e., the total size of the antecedents of the s-t tgds in Σ shall be minimal;
- conclusion-minimality i.e., the total size of the conclusions of the s-t tgds in Σ shall be minimal;
- variable-minimality i.e., the total number of existentially quantified variables in the conclusions shall be minimal.

Then, a set of s-t tgds is optimal, if it is minimal w.r.t. each of these four criteria. Following the above discussion, we only take s-t tgds into consideration for which no further splitting is possible. (We shall give a formal definition of this property and of the four optimality criteria in Sect. 3). Cardinality-minimality together with antecedent-minimality means that the cost of the join-operations is minimized when computing a canonical universal solution for some given source instance. Conclusion-minimality and variable-minimality mean that

no unnecessary incomplete facts are introduced in the canonical universal solution. For the transformation of arbitrary sets of s-t tgds into optimal ones, we shall present a normal system of rewrite rules. Moreover, we shall show that the optimal form of a set of s-t tgds is unique up to variable renaming.

In other words, our optimization of schema mappings is also a normalization of schema mappings. As an immediate benefit of a normalization, we get a purely syntactical criterion for testing the equivalence of two schema mappings. Another, even more important application of such a normalization is in the area of defining the semantics of query answering in data exchange. Several definitions in this area depend on the concrete syntactic representation of the s-t tgds. This is, in particular, the case for queries with negated atoms (see e.g., [2, 20]) and for aggregate queries (see [1]). This semantic dependence on the syntax of a mapping clearly is undesirable. Since the minimal set of s-t tgds produced by our rewrite rules is unique up to variable renaming, we can use it as the desired normal form, which eliminates the effect of the concrete representation of the s-t tgds from the semantics of query answering.

Example 4 Consider a schema mapping $\gamma = \langle S, T, \Sigma \rangle$ with $S = \{ S(\cdot, \cdot, \cdot) \}$, $T = \{ L(\cdot, \cdot, \cdot), P(\cdot, \cdot) \}$, where L and P are as in Example 2. S denotes the relational schema $\text{Student}(\text{name}, \text{year}, \text{area})$. Moreover, let Σ express the following constraints: If there exists a student in any year, then there should exist at least one lecture for this year. Moreover, if a student specializes in a particular area, then there should be a professor in this area teaching at least one lecture for this year. We thus have the following set with a single s-t tgd:

$$S(x_1, x_2, x_3) \rightarrow (\exists y_1, y_2, y_3, y_4, y_5) L(y_1, x_2, y_3) \wedge L(y_4, x_2, y_5) \wedge P(y_5, x_3)$$

Clearly, the first atom in the conclusion may be deleted. Now, consider the source instance $\sigma = \{ S(\text{bob}, 3, \text{db}) \}$, and suppose that we want to evaluate the query

$$\text{ans}(x_2) : \exists x_1 L(x_1, x_2, x_3) \wedge \neg P(x_3, x_4)$$

over the target instance, i.e., we want to check if, in some year, there exists a lecture which has not been assigned to a professor. In [2, 20], query answering via the “canonical universal solution” (for details, see Sect. 2) is proposed. Depending on whether the s-t tgd σ has been simplified or not, we either get $J = \{ L(u_1, 3, u_2), L(u_3, 3, u_4), P(u_4, \text{db}) \}$ or the core thereof, $J = \{ L(u_1, 3, u_2), P(u_2, \text{db}) \}$ as the canonical universal solution. In the first case, the query yields the result $\{ 3 \}$ whereas, in the second case, we get

Similarly, a unique normal form of the s-t tgds is crucial for the semantics of aggregate queries in data exchange, whose investigation has been initiated recently by Afrati and

Kolaitis [1]. Aggregate queries are of the form $\text{SELECT } f$ $\text{FROM } R$, where f is an aggregate operator in $\{ \text{min}(R.A), \text{max}(R.A), \text{count}(R.A), \text{count}(\), \text{sum}(R.A), \text{avg}(R.A) \}$, and where R is a target relation symbol or, more generally, a conjunctive query over the target schema Σ and A is an attribute of R . On the one hand, [1] defines an interesting and non-trivial semantics of aggregate queries in data exchange. On the other hand, it is shown that the most important aggregate queries can be evaluated in polynomial time (data complexity). In this paper, we shall show how aggregate queries can benefit from our normalization of schema mappings.

So far, we have only mentioned mappings $\Sigma = \langle S, T, \Sigma \rangle$, where Σ is a set of s-t tgds. In addition, Σ may contain constraints on the target instance alone. One of the most important forms of target constraints are equality-generating target dependencies (egds, for short), which can be considered as a generalization of functional dependencies. Egds are formulas of the form $\exists (\varphi(x) \rightarrow x_i = x_j)$ where φ is a CQ over T and x_i, x_j are variables in x .

Example 5 We modify the setting from Examples 1 and 2. Let $M = \langle S, T, \Sigma \rangle$ with $S = \{ C(\cdot, \cdot, \cdot) \}$ and $T = \{ P(\cdot, \cdot, \cdot) \}$ where C and P denote the relational schemes $\text{course}(\text{title}, \text{course-area}, \text{prof-area})$ and $\text{Prof}(\text{name}, \text{prof-area}, \text{course-area})$. The P relation thus contains information on the area of the professor as well as on the area(s) of the courses taught by him/her. The set Σ of s-t tgds expresses the following constraints: For every course, there exists a professor who teaches courses in his/her own area of expertise and who teaches courses with this combination of course and prof area. Moreover, there exists a professor whose expertise matches the area of the course and vice versa. We thus define a mapping with the following two s-t tgds:

$$C(x_1, x_2, x_3) \rightarrow (\exists y_1, y_2) P(y_1, y_2, y_2) \wedge P(y_1, x_2, x_3)$$

$$C(x_1, x_2, x_3) \rightarrow (\exists y_1) P(y_1, x_3, x_2)$$

This set of dependencies is minimal. However, suppose that we add the egd $P(x_1, x_2, x_3) \rightarrow x_2 = x_3$, expressing that a professor only teaches courses in his/her own area of expertise. Then $P(y_1, y_2, y_2)$ can be eliminated from the conclusion of the first s-t tgd. Moreover, the first and the second s-t tgd imply each other. Hence Σ can be replaced by either Σ' or Σ'' with

$$\Sigma' = \{ C(x_1, x_2, x_3) \rightarrow (\exists y_1) P(y_1, x_2, x_3) \}$$

$$\Sigma'' = \{ C(x_1, x_2, x_3) \rightarrow (\exists y_1) P(y_1, x_3, x_2) \}.$$

Example 5 illustrates that, in the presence of target egds, our rewrite rules for the s-t tgds-only case are not powerful enough. To deal with target egds, we will introduce further rewrite rules. In particular, one of these new rewrite rules will

result in the introduction of source egds to prevent situations where two sets of s-t tgds only differ on source instances which admit no target instance anyway. Indeed, in Example 5, Σ' and Σ'' only differ if $x_2 = x_3$ holds. But this is forbidden by the egd. Hence Σ should be replaced by Σ' with

$$\Sigma' = \{ C(x_1, x_2, x_3) \rightarrow x_2 = x_3, C(x_1, x_2, x_2) \rightarrow (\exists y_1) P(y_1, x_2, x_2) \}.$$

In summary, we shall be able to prove that our extended set of rewrite rules again leads to a normal form which is unique up to variable renaming. The main ingredients of our normalization and optimization are the splitting and simplification of tgds. In the presence of target egds, several pitfalls will have to be avoided when defining appropriate splitting and simplification rules so as not to destroy the uniqueness of the normal form.

Organization of the paper and summary of results In Sect. 2, we recall some basic notions. A conclusion and an outlook to future work are given in Sect. 6. The main results of the paper are detailed in the Sections 3, 4, and 5, namely:

- Optimization and normalization of sets of s-t tgds. In Sect. 3, we give a formal definition of the above-mentioned optimality criteria for sets of s-t tgds and we present rewrite rules to transform any set of s-t tgds into an optimal one (i.e., minimal w.r.t. to these criteria). We shall also show that the normal form obtained by our rewrite rules is unique up to variable renaming. Moreover, we show that if the length of each s-t tgd is bounded by a constant, then this normal form can be computed in polynomial time.
- Extension to target egds. In Sect. 4, the rewrite rule system for s-t tgds is then extended to schema mappings comprising target egds. Several non-trivial extensions (like the introduction of source egds) are required to arrive at a unique normal form again. The extended splitting and simplification rules will have to be defined very carefully so as not to destroy this uniqueness.
- Semantics of aggregate operators. In Sect. 5, we discuss in detail the application of our normalization of schema mappings to the definition of a unique semantics of aggregate operators in data exchange. Some proof details will be omitted due to space limitations. They can be found in the technical report [6].

of a relation for each relation symbol R , s.t. both have the same arity. We only consider finite instances here.

Tuples of the relations may contain two types of constants and variables. The latter are often also called marked nulls or labeled nulls. Two labeled nulls are equal if they have the same label. For every instance I , we write $\text{dom}(I)$, $\text{var}(I)$, and $\text{Const}(I)$ to denote the set of terms, variables, and constants, respectively. Clearly, $\text{dom}(I) = \text{var}(I) \cup \text{Const}(I)$. If we have no particular instance in mind, we write Const to denote the set of all possible constants. We write for a tuple (x_1, x_2, \dots, x_n) . However, by slight abuse of notation, we also refer to the set $\{x_1, \dots, x_n\}$ as x . Hence, we may use expressions like x_i , x or $x \in X$, etc.

Let $S = \{S_1, \dots, S_n\}$ and $T = \{T_1, \dots, T_m\}$ be schemas with no relation symbols in common. We call S the source schema and T the target schema. We write S, T to denote the schemas $\{S_1, \dots, S_n, T_1, \dots, T_m\}$. Instances over S and T are called source and target instances, respectively. If I is a source instance and J a target instance, then their combination I, J is an instance of the schema S, T .

Homomorphisms and substitutions Let I, I' be instances. A homomorphism $h: I \rightarrow I'$ is a mapping $\text{dom}(I) \rightarrow \text{dom}(I')$, s.t. (1) whenever a fact $R(x) \in I$, then $R(h(x)) \in I'$, and (2) for every constant $b(c) = c$. If such h exists, we write $I \rightarrow I'$. Moreover, if $I \rightarrow I'$ then we say that I and I' are homomorphically equivalent. In contrast, if $I \rightarrow I'$ but not vice versa, we say that I is more general than I' , and I' is more specific than I .

If $h: I \rightarrow I'$ is invertible, s.t. h^{-1} is a homomorphism from I' to I , then h is called an isomorphism, denoted $I \cong I'$. An endomorphism is a homomorphism $I \rightarrow I$. An endomorphism is proper if it is not surjective (for finite instances, this is equivalent to being not injective), i.e., if it reduces the domain of I .

If I is an instance, and I' is such that $I \rightarrow I'$ holds but for no other $I'' : I \rightarrow I''$ (that is, I cannot be further “shrunk” by a proper endomorphism), then I' is called a core of I . The core is unique up to isomorphism. Hence, we may speak about the core of I . Cores have the following important property: for arbitrary instances I, J , $J \rightarrow I$ if and only if $\text{core}(J) \rightarrow \text{core}(I)$.

A substitution σ is a mapping which sends variables to other domain elements (i.e., variables or constants). We write $\sigma = \{x_1 \mapsto a_1, \dots, x_n \mapsto a_n\}$ if σ maps each x_i to a_i and σ is the identity outside $\{x_1, \dots, x_n\}$. The application of a substitution is usually denoted in post- x notation, e.g., $\sigma(x)$ denotes the image of x under σ . For an expression $\varphi(x)$, which in the following will normally refer to a conjunctive query with variables in x , we write $\varphi(x\sigma)$ to denote the result of replacing every occurrence of every variable x by $x\sigma$.

Schema mappings and data exchange A schema mapping is given by a triple $M = (S, T, \Sigma)$ where S is the source schema, T is the target schema, and Σ is a set of dependencies expressing the relationship between S and T , and possibly also local constraints on S and T . The data exchange problem associated with M is the following: Given a (ground) source instance I , and a target instance J , s.t. $I, J \models \Sigma$. Such a J is called a solution for I or, simply, a solution if I is clear from the context. The set of all solutions for I under M is denoted by $\text{Sol}^M(I)$. If $J \in \text{Sol}^M(I)$ is such that $J \rightarrow J'$ holds for any other solution $J' \in \text{Sol}^M(I)$, then J is called a universal solution. Since the universal solutions for a source instance I are homomorphically equivalent, the core of the universal solutions for I is unique up to isomorphism. It is the smallest universal solution [1].

In the following, we will often identify a schema mapping $M = (S, T, \Sigma)$ with the set of dependencies Σ , without explicitly mentioning the schemas, for the sake of brevity.

Equivalence of schema mappings Different notions of equivalence of schema mappings have been recently proposed by Fagin et al. [10]. In this paper, we will only consider the strongest one, namely logical equivalence.

Definition 1 [10] Two schema mappings $M = (S, T, \Sigma)$ and $M' = (S, T, \Sigma')$ over the schema S, T are logically equivalent (denoted $M \equiv M'$) if, for every source instance I and target instance J , the equivalence $I, J \models \Sigma \iff I, J \models \Sigma'$ holds. In this case, the equality $\text{Sol}^M(I) = \text{Sol}^{M'}(I)$ holds for every source instance I .

Dependencies Embedded dependencies [6] over a relational schema R are first-order formulas of the form

$$\varphi(x) \rightarrow \exists y \psi(x, y)$$

In case of tuple-generating dependencies (tgds), both the antecedent φ and conclusion ψ are conjunctive queries (CQs) over the relation symbols from R such that all variables in ψ actually do occur in $\varphi(x)$. Equality-generating dependencies (egds) are of the form

$$\varphi(x) \rightarrow x_i = x_j$$

with $x_i, x_j \in x$. Throughout this paper, we shall omit the universal quantifiers: By convention, all variables occurring in the antecedent are universally quantified (over the entire formula). In the context of data exchange, we are mainly dealing with source-to-target dependencies consisting of tuple-generating dependencies (or s-t tgds) over the schema S, T (the antecedent is a CQ over S , the conclusion over T) and target dependencies over T . In the scope of this paper, target dependencies are restricted to equality-generating dependencies (referred to as “target egds”). Moreover, in Sect. 4 we

shall also consider source dependencies consisting of egds over S (referred to as “source egds”).

Database of a conjunctive query Given a conjunctive query χ , we write $At(\chi)$ to denote the database comprising exactly the set of atoms of χ . If the variables of χ are instantiated with distinct fresh constants in $At(\chi)$, this database is called frozen. However, unless otherwise specified, we assume that $At(\chi)$ is not frozen and that the variables of χ are instantiated with distinct labeled nulls in $At(\chi)$. If χ represents an antecedent or conclusion of some dependency τ , $At(\chi)$ is called the antecedent or, respectively, conclusion database of τ .

Chase The data exchange problem can be solved by the chase [4], a sequence of steps, each enforcing a single constraint within some limited set of tuples. More precisely, let I contain a tgd: $\varphi(x) \wedge (\exists y)\psi(x, y)$, s.t. $I \models \varphi(a)$ for some assignment a on x . Then, we extend I with facts corresponding to $\psi(a, z)$, where the elements z are fresh labeled nulls. Note that this definition of the chase differs from the definition in [9], where no new facts are added if $\exists y\psi(a, y)$ is already fulfilled. Omitting this check is referred to as oblivious [17] chase. It is the preferred version of chase if the result of the chase should not depend on the order in which the tgd's are applied (see e.g. [2, 20]).

Now suppose that I contains an egd: $\varphi(x) \wedge x_i = x_j$, s.t. $I \models \varphi(a)$ for some assignment a on x . This egd enforces the equality $a_i = a_j$. We thus choose a null among $\{a_i, a_j\}$ and replace every occurrence of a in I by the other term; if $a_i, a_j \in \text{Const}(I)$ and $a_i = a_j$, the chase halts with failure. We write I^Σ to denote the result of chasing I with the dependencies Σ .

Consider an arbitrary schema mapping $\mathbb{M} = \langle \Sigma_{st}, \Sigma_t \rangle$ where Σ_{st} is a set of source-to-target tgd's and Σ_t is a set of target egds. Then, the solution to a source instance I can be computed as follows: We start off with the instance I , i.e., the source instance I and the target instance is initially empty. Chasing I with Σ_{st} yields the instance J , where J is called the preuniversal instance. This chase always succeeds since Σ_{st} contains no egds. Then, J is chased with Σ_t . This chase may fail on an attempt to unify distinct constants. If the chase succeeds, we end up with J^Σ_t , which is referred to as the canonical universal solution $\text{CanSo}(I)$ or, simply $\text{CanSo}(I)$. Both J and J^Σ_t can be computed in polynomial time w.r.t. the size of the source instance I .

3 Normalization of s-t tgds

In this section, we investigate ways of optimizing sets of s-t tgds. In the first place, we thus formulate some natural optimality criteria. The following parameters of a set of s-t tgds will be needed in the definition of such criteria:

Definition 2 Let Υ be a set of s-t tgds. Then we define:

- $|\Upsilon|$ denotes the number of s-t tgds in Υ .
- $\text{AntSize}(\Upsilon) = \sum \{ |At(\varphi(x))| : \varphi(x) \wedge (\exists y)\psi(x, y) \text{ is an s-t tgd in } \Upsilon \}$, i.e., $\text{AntSize}(\Upsilon)$ is the total number of atoms in all antecedents of tgds in Υ .
- $\text{ConSize}(\Upsilon) = \sum \{ |At(\psi(x, y))| : \varphi(x) \wedge (\exists y)\psi(x, y) \text{ is an s-t tgd in } \Upsilon \}$, i.e., $\text{ConSize}(\Upsilon)$ is the total number of atoms in all conclusions of tgds in Υ .
- $\text{VarSize}(\Upsilon) = \sum \{ |y| : \varphi(x) \wedge (\exists y)\psi(x, y) \text{ is in } \Upsilon \}$, i.e., $\text{VarSize}(\Upsilon)$ is the total number of existentially quantified variables in all conclusions of tgds in Υ .

We would naturally like to transform any set of s-t tgds into an equivalent one where all the above parameters are minimal. Recall however our discussion on the splitting of s-t tgds from Example 3. As we pointed out there, the splitting of s-t tgds is comparable to normal form decomposition of relational schemas. It should clearly be applied in order to avoid anomalies like the introduction of obviously irrelevant atoms in the canonical universal solution as we saw in Example 4 where the set Σ (with two split s-t tgds) was certainly preferable to Σ' even though $|\Sigma| < |\Sigma'|$ and $\text{AntSize}(\Sigma) < \text{AntSize}(\Sigma')$. Note that in Example 3, the equality $\text{ConSize}(\Sigma) = \text{ConSize}(\Sigma')$ holds. Intuitively, the effect of splitting is that the atoms in the conclusion of some s-t tgd are distributed over several strictly smaller s-t tgds. Thus, our goal should be to find an optimal set of s-t tgds (that is, a set where the above-mentioned parameters are minimal) among those sets of s-t tgds for which no further splitting is possible. We now make precise what it means that “no further splitting” is possible and formally define optimality of a set of s-t tgds.

Definition 3 Let Σ be a set of s-t tgds. We say that Σ is split-reduced if there exists no Σ' equivalent to Σ , s.t. $|\Sigma'| > |\Sigma|$ but $\text{ConSize}(\Sigma') = \text{ConSize}(\Sigma)$.

Definition 4 Let Σ be a set of s-t tgds. We say that Σ is optimal if it is split-reduced and if each of the parameters $|\Sigma|$, $\text{AntSize}(\Sigma)$, $\text{ConSize}(\Sigma)$, and $\text{VarSize}(\Sigma)$ is minimal among all split-reduced sets equivalent to Σ .

Of course, given an arbitrary set of s-t tgds, it is a priori not clear that an optimal set equivalent to Σ exists, since it might well be the case that some set minimizes some of the parameters, while another set minimizes the other parameters. The goal of this section is to show that optimality in the above sense can always be achieved and to construct an algorithm, which transforms any set of s-t tgds into an equivalent optimal one. To this end, we introduce a rewrite system that consists of two kinds of rewrite rules: rules that simplify each s-t tgd individually and rules that are applied to the entire set of s-t tgds. The following example illustrates several kinds of redundancy that a single s-t tgd may contain (and which may be eliminated with our rewrite rules).

Example 6 Consider the following dependency:

$$\tau : \begin{array}{l} S(x_1, x_3) \quad S(x_1, x_2) \quad (y_1, y_2, y_3, y_4, y_5) \\ P(x_1, y_2, y_1) \quad R(y_1, y_3, x_2) \quad R(2, y_3, x_2) \\ P(x_1, y_4, 2) \quad P(x_1, y_4, y_5) \quad Q(y_4, x_3) \end{array}$$

Clearly, τ is equivalent to the set $\{\tau_1, \tau_2\}$ of s-t tgds:

$$\begin{array}{l} \tau_1 : \begin{array}{l} S(x_1, x_3) \quad S(x_1, x_2) \quad (y_1, y_2, y_3) \\ P(x_1, y_2, y_1) \quad R(y_1, y_3, x_2) \quad R(2, y_3, x_2) \end{array} \\ \tau_2 : \begin{array}{l} S(x_1, x_3) \quad S(x_1, x_2) \quad (y_4, y_5) \\ P(x_1, y_4, y_5) \quad P(x_1, y_4, 2) \quad Q(y_4, x_3) \end{array} \end{array}$$

Now the antecedents of τ_1 and τ_2 can be simplified:

$$\begin{array}{l} \tau_1 : \begin{array}{l} S(x_1, x_2) \quad (y_1, y_2, y_3) \\ P(x_1, y_2, y_1) \quad R(y_1, y_3, x_2) \quad R(2, y_3, x_2) \end{array} \\ \tau_2 : \begin{array}{l} S(x_1, x_3) \quad (y_4, y_5) \\ P(x_1, y_4, y_5) \quad P(x_1, y_4, 2) \quad Q(y_4, x_3) \end{array} \end{array}$$

Finally, we may also simplify the conclusion of:

$$\tau_2 : S(x_1, x_3) \quad (y_4)P(x_1, y_4, 2) \quad Q(y_4, x_3)$$

In total, τ is equivalent to $\{\tau_1, \tau_2\}$.

For the simplifications illustrated in Example 6, we define the rewrite rules 1–3 in Fig. 1. Rules 1 and 2 replace an s-t tgd τ by a simpler one (i.e., with fewer atoms), while Rule 3 replaces τ by a set $\{\tau_1, \dots, \tau_n\}$ of simpler s-t tgds. These rules make use of the following definitions of the core and the components of CQs.

Rewrite rules to simplify a set of s-t tgds

Rule 1 (Core of the conclusion, see Definition 5).

$$\begin{array}{l} \tau: \varphi(\mathbf{x}) \rightarrow (\exists \mathbf{y})\psi(\mathbf{x}, \mathbf{y}) \implies \\ \tau': \varphi(\mathbf{x}) \rightarrow (\exists \mathbf{y})\psi(\mathbf{x}, \mathbf{y}\sigma), \\ \text{s.t. } \psi(\mathbf{x}, \mathbf{y}\sigma) \text{ is the core of } \psi(\mathbf{x}, \mathbf{y}). \end{array}$$

Rule 2 (Core of the antecedent, see Definition 5).

$$\begin{array}{l} \tau: \varphi(\mathbf{x}_1, \mathbf{x}_2) \rightarrow (\exists \mathbf{y})\psi(\mathbf{x}_1, \mathbf{y}) \implies \\ \tau': \varphi(\mathbf{x}_1, \mathbf{x}_2\sigma) \rightarrow (\exists \mathbf{y})\psi(\mathbf{x}_1, \mathbf{y}), \\ \text{s.t. } \varphi(\mathbf{x}_1, \mathbf{x}_2\sigma) \text{ is the core of } \varphi(\mathbf{x}_1, \mathbf{x}_2). \end{array}$$

Rule 3 (Splitting, see Definition 6).

$$\begin{array}{l} \tau: \varphi(\mathbf{x}) \rightarrow (\exists \mathbf{y})\psi(\mathbf{x}, \mathbf{y}) \implies \{\tau_1, \dots, \tau_n\}, \text{ s.t.} \\ \{\psi_1(\mathbf{x}, \mathbf{y}_1), \dots, \psi_n(\mathbf{x}, \mathbf{y}_n)\} \text{ are the components of } \psi(\mathbf{x}, \mathbf{y}) \\ \text{and } \tau_i: \varphi(\mathbf{x}) \rightarrow (\exists \mathbf{y}_i)\psi_i(\mathbf{x}, \mathbf{y}_i) \text{ for } i \in \{1, \dots, n\}. \end{array}$$

Rule 4 (Implication of an s-t tgd).

$$\begin{array}{l} \Sigma \implies \Sigma \setminus \{\tau\} \\ \text{if } \Sigma \setminus \{\tau\} \models \tau. \end{array}$$

Rule 5 (Implication of atoms in the conclusion).

$$\begin{array}{l} \Sigma \implies (\Sigma \setminus \{\tau\}) \cup \{\tau'\} \\ \text{if } \tau: \varphi(\mathbf{x}) \rightarrow (\exists \mathbf{y})\psi(\mathbf{x}, \mathbf{y}) \\ \text{and } \tau': \varphi(\mathbf{x}) \rightarrow (\exists \mathbf{y}')\psi'(\mathbf{x}, \mathbf{y}'), \\ \text{s.t. } At(\psi'(\mathbf{x}, \mathbf{y}')) \subset At(\psi(\mathbf{x}, \mathbf{y})) \\ \text{and } (\Sigma \setminus \{\tau\}) \cup \{\tau'\} \models \tau. \end{array}$$

Fig. 1 Redundancy elimination from a set of s-t tgds

Definition 5 Let $\chi(u, v)$ be a CQ with variables in $u \cup v$ and let A denote the structure consisting of the atoms $At(\chi(u, v))$, s.t. the variables in u are considered as constants and the variables in v as labeled nulls. Let A denote the core of A with $A \subseteq A$, i.e., there exists a substitution $\sigma: v \rightarrow Const \cup v$ s.t. $At(\chi(u, v\sigma)) = A \subseteq At(\chi(u, v))$. Then, we define the core of $\chi(u, v)$ as the CQ $\chi(u, v\sigma)$.

Definition 6 Let $\chi(u, v)$ be a CQ with variables in $u \cup v$. We set up the dual graph $G(\tau)$ as follows: The atoms of $\chi(u, v)$ are the vertices of $G(\tau)$. Two vertices are connected if the corresponding atoms have at least one variable in common. Let $\{C_1, \dots, C_n\}$ denote the connected components of $G(\tau)$. Moreover, for every $i \in \{1, \dots, n\}$, let v_i with $v_i \subseteq v$ denote those variables from v which actually occur in C_i and let $\chi_i(u, v_i)$ denote the CQ consisting of the atoms in C_i . Then we define the components of $\chi(u, v)$ as the set $\{\chi_1(u, v_1), \dots, \chi_n(u, v_n)\}$.

The splitting rule (i.e., Rule 3 in Fig. 1) was already applied in Example 3. Rule 2 involving core computation of the antecedent was applied in Example 4 when we reduced $L(x_1, x_2, x_3) \rightarrow L(x_4, 3, x_5) \rightarrow P(x_5, x_6)$ to its core $L(x_4, 3, x_5) \rightarrow P(x_5, x_6)$. Likewise, in Example 6, the simplification of τ_1 and τ_2 to τ_1 and τ_2 is due to Rule 2. In a similar way, Rule 1 involving core computation of the conclusion allowed us to reduce $(y_1, y_2, y_3) \rightarrow L(x_1, 3, y_4) \rightarrow P(y_4, x_2)$ in Example 2 to $L(x_1, 3, y_4) \rightarrow P(y_4, x_2)$. In Example 6, Rule 1 was applied when we replaced τ_2 by τ_2 .

The following example illustrates that additional rules are required in order to remove an s-t tgd or a part of an s-t tgd whose redundancy is due to the presence of other s-t tgds.

Example 7 Consider the set $\Sigma = \{\tau_1, \tau_2, \tau_3\}$, where τ_1 and τ_2 are the s-t tgds resulting from the simplification steps in Example 6 and τ_3 is given below:

$$\begin{array}{l} \tau_1 : \begin{array}{l} S(x_1, x_2) \quad (y_1, y_2, y_3) \\ P(x_1, y_2, y_1) \quad R(y_1, y_3, x_2) \quad R(2, y_3, x_2) \end{array} \\ \tau_2 : \begin{array}{l} S(x_1, x_3) \quad (y_4)P(x_1, y_4, 2) \quad Q(y_4, x_3) \end{array} \\ \tau_3 : S(2, x) \quad (y)R(2, y, x) \end{array}$$

The tgd τ_3 generates only a part of the atoms that does and resides in strictly fewer cases than τ_1 . Hence, τ_3 may be deleted. Moreover, considering the combined effect of the rules τ_1 and τ_2 , which reside on exactly the same tuples, and a substitution $\{y_1 \rightarrow 2, y_2 \rightarrow y_4\}$, we notice that the first two atoms in the conclusion of τ_1 are in fact redundant, and it is possible to reduce τ_1 to $\tau_1 : S(x_1, x_2) \quad (y_3)R(2, y_3, x_2)$. In total, Σ may be replaced by $\Sigma' = \{\tau_1, \tau_2\}$.

Rules 4 and 5 in Fig. 1 allow us to eliminate such redundancies from a set Σ of s-t tgds: By Rule 4, we may delete a s-t tgd τ from Σ , if τ is implied by the others, like τ_3 in Example 7. Rule 5 allows us to replace a rule by a

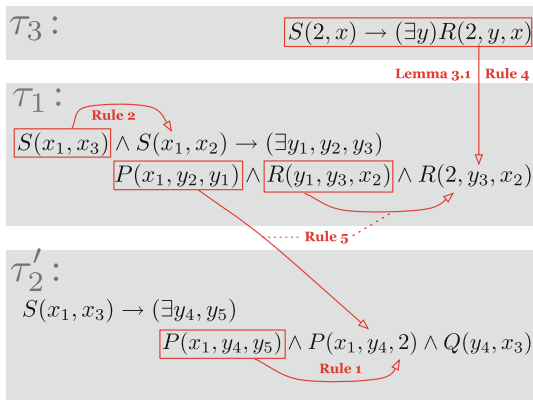


Fig. 2 Tgd optimization. Rectangles mark eliminated atoms, arrows show justifications for elimination

strictly smaller rule (with fewer atoms in the conclusion) if it is implied by τ together with the remaining s-t tgds in Σ (cf. the replacement of τ_1 with τ_1 in Example 7 above). Figure 2 illustrates the elimination of redundant atoms via Rules 1, 2, 4, and 5 in a set $\{\tau_1, \tau_2, \tau_3\}$ of tgds from Examples 6 and 7.

In principle, the implication of a tgd by a set of dependencies can be tested by a procedural criterion based on the chase [4]. For our purposes, the following declarative criterion is more convenient.

Lemma 1 Consider a s-t tgd $\tau : \varphi(x) \rightarrow (\exists y)\psi(x, y)$ and a set Σ of s-t tgds. Then $\Sigma \models \tau$ holds if there exist (not necessarily distinct) s-t tgds τ_1, \dots, τ_k in Σ , such that all s-t tgds $\tau, \tau_1, \dots, \tau_k$ are pairwise variable disjoint and the following conditions hold:

- (a) For every $i \in \{1, \dots, k\}$, there exists a substitution $\lambda_i : x_i \rightarrow \text{Const } x$, s.t. $\text{At}(\varphi_i(x_i \lambda_i)) \subseteq \text{At}(\varphi(x))$.
- (b) A substitution $\mu : y \rightarrow \text{Const } x \cup \bigcup_{i=1}^k y_i$ exists, s.t. $\text{At}(\psi(x, y \mu)) \subseteq \bigcup_{i=1}^k \text{At}(\psi_i(x_i \lambda_i, y_i))$.

Proof For the “ \Leftarrow ”-direction, consider an arbitrary pair S, T of source and target instance, s.t. $S, T \models \Sigma$. It is easy to show that by conditions (a) and (b), then also $S, T \models \tau$ holds. For the “ \Rightarrow ”-direction, we take the source instance $S = \text{At}(\varphi(x))$, where we consider the variables as constants. Moreover, let T denote the target instance which results from the oblivious chase of S with Σ . Let τ_1, \dots, τ_k denote the (not necessarily distinct) s-t tgds whose antecedent can be mapped in S via substitutions $\lambda_1, \dots, \lambda_k$. These substitutions satisfy the condition (a). By $S, T \models \Sigma$ and $\Sigma \models \tau$, we get the desired substitution for condition (b).

Note that Rule 5 generalizes Rule 1 and, in principle, also Rule 4. Indeed, if we restrict Σ in Rule 5 to the singleton $\Sigma = \{\tau\}$, then the replacement of τ by τ means that we

reduce $\psi(x, y)$ to its core. Moreover, Rule 5 allows us to eliminate all atoms from the conclusion of τ if τ may be deleted via Rule 4. Clearly, the deletion of the conclusion of τ essentially comes down to the deletion of τ itself. The correctness of Rules 1–5 is easily established.

Lemma 2 The Rules 1–5 in Fig. 4 are correct, i.e.: Let Σ be a set of s-t tgds and $\tau \in \Sigma$. Suppose that Σ is transformed into Σ' by applying one of the Rules 1–5, that is:

- ⊖ τ is replaced by a single s-t tgd (via Rule 1, 2, 5),
- ⊖ τ is replaced by s-t tgds τ_1, \dots, τ_n (via Rule 3),
- ⊖ τ is deleted (via Rule 4).

Then Σ and Σ' are equivalent.

The following notion of a “proper instance” of an s-t tgd plays an important role for proving that our Rules 1–5 lead to a unique normal form. A proper instance of an s-t tgd is obtained from τ by eliminating at least one existentially quantified variable in the conclusion of τ , while keeping the antecedent unchanged.

Definition 7 Let $\tau : \varphi(x) \rightarrow (\exists y)\psi(x, y)$ be an s-t tgd. We call an s-t tgd τ' a proper instance of τ , if there exists a strict subset y' and a substitution $\sigma : y \rightarrow \text{Const } x \cup y'$, such that τ' is of the form $\tau' : \varphi(x) \rightarrow (\exists y')\psi(x, y \sigma)$.

Example 8 In the following three tgds, each next tgd is a proper instance of the previous ones:

- $\tau_1 : S(x_1, x_2) \rightarrow (\exists y_1, y_2)Q(x_1, y_1, y_2)$
- $\tau_2 : S(x_1, x_2) \rightarrow (\exists y_1)Q(x_1, y_1, y_1)$
- $\tau_3 : S(x_1, x_2) \rightarrow Q(x_1, x_2, x_2)$

Moreover, observe that $\tau_2 \models \tau_1$ and $\tau_3 \models \tau_2$ holds.

The importance of “proper instances” to our investigations comes from the following properties:

Lemma 3 Let τ and τ' be s-t tgds, s.t. τ' is a proper instance of τ . Then the following properties hold:

- (1) $\tau \models \tau'$.
- (2) Suppose that τ is reduced with respect to Rule 1. Then $\tau' \models \tau$.

Proof (Sketch) The proof of the first claim is easy. For the second one, suppose that $\tau \models \tau'$ holds. We have to show that then Rule 1 is applicable to τ . Let S, T denote a pair of source and target instance with $S = \text{At}(\varphi(x))$ and $T = \text{At}(\psi(x, y))$. The variables in x are thus considered as constants, while y are labeled nulls. Clearly, $S, T \models \tau$ and $S \models \varphi(x)$. Thus, by the assumption $\tau \models \tau'$, also $T \models \psi(x, y \sigma)$ holds, i.e., there exists a substitution σ on $y \setminus \text{dom}(T)$ such

that $At(\psi(x, y\sigma\mu)) \vdash \tau$. Hence, also the following inclusion holds: $At(\psi(x, y\sigma\mu)) \subseteq At(\psi(x, y))$. Note that $y\sigma = y$. Hence, also $At(\psi(x, y\sigma\mu)) \subseteq At(\psi(x, y))$. But then $At(\psi(x, y))$ is not a core and, therefore, Rule 1 is applicable to τ .

Lemma 4 Let τ be an s-t tgd reduced w.r.t. the Rules 1 and 3 and let Σ be a set of s-t tgds. $\Sigma \mid \tau$, then one of the following two conditions is fulfilled: Either

- Ⓓ there exists a single s-t tgd $\tau_0 \in \Sigma$, s.t. $\tau_0 \mid \tau$, or
- Ⓔ there exists a proper instance of τ , s.t. $\Sigma \mid \tau$.

Proof Let $\{\tau_1, \dots, \tau_k\} \subseteq \Sigma \setminus \{\tau\}$ with $\{\tau_1, \dots, \tau_k\} \mid \tau$. Suppose that k is minimal with this property and that $k \geq 2$. We show that then $\{\tau_1, \dots, \tau_k\} \mid \tau$ holds for some proper instance of τ . For $i \in \{1, \dots, k\}$, let τ_i s be pairwise variable disjoint and have the form: $\varphi_i(x_i) \wedge (\exists y_i)\psi_i(x_i, y_i)$. By Lemma 1 and the definition of Rule 4, the τ_i s fulfill the following properties:

- (a) For every $i \in \{1, \dots, k\}$, there exists a substitution $\lambda_i : x_i \rightarrow \text{Const } x$, s.t. $At(\varphi_i(x_i\lambda_i)) \subseteq At(\varphi(x))$.
- (b) A substitution $\mu : y \rightarrow \text{Const } x \cup_{i=1}^k y_i$ exists s.t. $At(\psi(x, y\mu)) \subseteq \bigcup_{i=1}^k At(\psi_i(x_i\lambda_i, y_i))$.

Let $At(\psi(x, y)) = \{A_1, \dots, A_n\}$. Clearly, $n \geq k \geq 2$. Suppose that $\{A_1\mu, \dots, A_\alpha\mu\} \subseteq At(\psi_1(x_1\lambda_1, y_1))$ while $\{A_{\alpha+1}\mu, \dots, A_n\mu\} \subseteq \bigcup_{i=2}^k At(\psi_i(x_i\lambda_i, y_i))$. By assumption, τ is reduced w.r.t. Rule 3, i.e., the conclusion of either consists of a single atom without variables from or of atoms forming a single connected component of the dual graph $G(\tau)$. By $n \geq k \geq 2$, the former case can be excluded. Hence, the atoms $\{A_1, \dots, A_\alpha\}$ and $\{A_{\alpha+1}, \dots, A_n\}$ share at least one variable y , i.e., y occurs in some atom A_i with $i \in \{1, \dots, \alpha\}$ and in some atom A_j with $j \in \{\alpha+1, \dots, n\}$. Let $\ell = 1$ denote the index, s.t. $A_j\mu = At(\psi_\ell(x_\ell\lambda_\ell, y_\ell))$. In total, we thus have $A_i\mu = At(\psi_1(x_1\lambda_1, y_1))$ and, therefore, $y\mu = \text{Const } x \cup y_1$. On the other hand $A_j\mu = At(\psi_\ell(x_\ell\lambda_\ell, y_\ell))$ and, therefore $y\mu = \text{Const } x \cup y_\ell$. By assumption y_1 and y_ℓ are disjoint. Thus $y\mu = \text{Const } x$.

We construct the desired proper instance of τ as follows: Let $y := y \setminus \{y\}$ and define the substitution $\sigma : y \rightarrow \text{Const } x \cup y$, s.t. $y\sigma = y\mu$ and σ maps all other variables in y onto themselves. Then, we have $\sigma = \mu$, i.e., for every $y_i \in y$, $y_i\sigma\mu = y_i\mu$. Clearly, $\{\tau_1, \dots, \tau_k\} \mid \tau\mu$ and, therefore, also $\{\tau_1, \dots, \tau_k\} \mid \tau\sigma$ holds. But then τ is the desired proper instance of τ .

Lemma 5 Suppose that an s-t tgd $\tau \in \Sigma$ is reduced w.r.t. Rules 1 and 3 and that cannot be deleted via Rule 4. If there exists a proper instance of τ , s.t. $\Sigma \mid \tau$ holds, then there exists an s-t tgd τ_0 , s.t. τ_0 may be replaced by τ via Rule 5.

Proof Let $\tau : \varphi(x) \wedge (\exists y)\psi(x, y)$ and suppose that $\Sigma \mid \tau$ holds. Then there exist s-t tgds $\tau_1, \dots, \tau_k \in \Sigma$ of the form $\tau_i : \varphi_i(x_i) \wedge (\exists y_i)\psi_i(x_i, y_i)$, s.t. the conditions (a) and (b) of Lemma 4 are fulfilled, i.e.:

- (a) For every $i \in \{1, \dots, k\}$, there exists a substitution $\lambda_i : x_i \rightarrow \text{Const } x$, s.t. $At(\varphi_i(x_i\lambda_i)) \subseteq At(\varphi(x))$.
- (b) There exists a substitution $\mu : y \rightarrow \text{Const } x \cup_{i=1}^k y_i$, s.t. $At(\psi(x, y\mu)) \subseteq \bigcup_{i=1}^k At(\psi_i(x_i\lambda_i, y_i))$.

Let $\tau : \varphi(x) \wedge (\exists y)\psi(x, y)$. We claim that at least one of the τ_i coincides with τ (up to variable renaming). Suppose to the contrary that $\{\tau_1, \dots, \tau_k\} \mid \tau$ holds for every $i \in \{1, \dots, k\}$. Then, the above conditions (a) and (b) imply that $\tau \mid \tau$ holds by Lemma 4. Moreover, $\tau \mid \tau$ holds by Lemma 3, part (1). Thus, $\Sigma \setminus \{\tau\} \mid \tau$ and τ could be deleted by Rule 4, which is a contradiction.

Let $I = \{i \mid 1 \leq i \leq k, \text{ s.t. } \tau_i \text{ is obtained from } \tau \text{ via variable renaming}\}$. We define the CQ $\psi(x, y)$ of the s-t tgd $\tau : \varphi(x) \wedge (\exists y)\psi(x, y)$ as follows: $\Theta = \{A(x, y) \mid A(x, y) = At(\psi(x, y)) \text{ and } i \in I \text{ and } A(x\lambda_i, y) = At(\psi(x, y\mu)) = At(\psi_i(x_i\lambda_i, y_i))\}$. Moreover, we set $\psi(x, y) = \bigwedge_{A(x, y) \in \Theta} A(x, y)$. Clearly, $(\Sigma \setminus \{\tau\}) \cup \{\tau\} \mid \tau$ by Lemma 1. Thus, also $(\Sigma \setminus \{\tau\}) \cup \{\tau\} \mid \tau$, by Lemma 3, part (1). We claim that $At(\psi(x, y)) \subseteq At(\psi(x, y))$ holds. Suppose to the contrary that $At(\psi(x, y)) = At(\psi(x, y))$. Then, by the above definition of Θ and by Lemma 4, $\tau \mid \tau$ would hold. By Lemma 3, part (2), this implies that τ is not reduced w.r.t. Rule 1, which is a contradiction. Hence, τ is indeed the desired s-t tgd, s.t. τ may be replaced by τ via Rule 5.

We now define a normal form of s-t tgds via the rewrite rules of this section. We will then show that this normal form is unique up to isomorphism in the sense defined below.

Definition 8 Let Σ be a set of s-t tgds and let τ be the result of applying the Rules 1–5 of Fig. 2 exhaustively to Σ . Then Σ is the normal form of Σ .

Definition 9 Let $\tau_1 : \varphi_1(x_1) \wedge (\exists y_1)\psi_1(x_1, y_1)$ and $\tau_2 : \varphi_2(x_2) \wedge (\exists y_2)\psi_2(x_2, y_2)$ be two tgds. We say that τ_1 and τ_2 are isomorphic if τ_2 is obtained from τ_1 via variable renamings $\eta : x_1 \rightarrow x_2$ and $\vartheta : y_1 \rightarrow y_2$.

Let Σ_1 and Σ_2 be two sets of tgds. We say that Σ_1 and Σ_2 are isomorphic if $|\Sigma_1| = |\Sigma_2|$, every $\tau_1 \in \Sigma_1$ is isomorphic to precisely one $\tau_2 \in \Sigma_2$ and every $\tau_2 \in \Sigma_2$ is isomorphic to precisely one $\tau_1 \in \Sigma_1$.

We start by showing for two single s-t tgds τ_1 and τ_2 that logical equivalence and isomorphism coincide, provided that the s-t tgds are reduced via our rewrite rules. This result will then be extended to sets Σ_1 and Σ_2 of s-t tgds.

Lemma 6 Let τ_1 and τ_2 be two s-t tgds and suppose that τ_1 and τ_2 are reduced w.r.t. Rules 1–3. Then τ_1 and τ_2 are isomorphic, iff τ_1 and τ_2 are equivalent.

Proof (Sketch) The “ \Rightarrow ”-direction follows immediately from Lemma 1. For the “ \Leftarrow ”-direction, let τ_1 and τ_2 be equivalent s-t tgds with $\tau_1: \varphi_1(x_1, x_2) \wedge (\forall y)\psi_1(x_1, y)$ and $\tau_2: \varphi_2(u_1, u_2) \wedge (\forall v)\psi_2(u_1, v)$.

By Lemma 1, there exist substitutions λ and ρ , s.t.

$$\begin{aligned} \lambda: x_1 & x_2 \quad \text{Const } u_1 \quad u_2, \text{ and} \\ \rho: u_1 & u_2 \quad \text{Const } x_1 \quad x_2, \text{ such that} \\ \text{At}(\varphi_1(x_1\lambda, x_2\lambda)) & \text{At}(\varphi_2(u_1, u_2)) \text{ and} \\ \text{At}(\varphi_2(u_1\rho, u_2\rho)) & \text{At}(\varphi_1(x_1, x_2)). \end{aligned}$$

By exploiting the equivalence of τ_1 and τ_2 and the fact that these s-t tgds are reduced w.r.t. Rule 2, we can show that the antecedents of τ_1 and τ_2 are isomorphic (i.e., the above inclusions are in fact equalities). Moreover, by exploiting that the s-t tgds are reduced w.r.t. Rule 1, we may conclude that τ_1 and τ_2 are isomorphic.

Theorem 1 Let Σ_1 and Σ_2 be equivalent sets of s-t tgds, i.e., $\Sigma_1 \mid \Sigma_2$ and $\Sigma_2 \mid \Sigma_1$. Let Σ_1 and Σ_2 denote the normal form of Σ_1 and Σ_2 , respectively. Then Σ_1 and Σ_2 are isomorphic.

Proof Let Σ_1 and Σ_2 be equivalent. Moreover, let Σ_1 and Σ_2 denote the normal form of Σ_1 and Σ_2 , respectively. By the correctness of our rewrite rules 1–5, of course, Σ_1 and Σ_2 are equivalent.

We first show that every s-t tgd in Σ_1 is isomorphic to some s-t tgd in Σ_2 and vice versa. Suppose to the contrary that this is not the case. W.l.o.g., we assume that there exists a $\tau \in \Sigma_1$, which is not isomorphic to any s-t tgd in Σ_2 . By the equivalence of Σ_1 and Σ_2 , the implication $\Sigma_2 \mid \tau$ clearly holds. By Lemma 4, either $\tau_0 \mid \tau$ for a single s-t tgd $\tau_0 \in \Sigma_2$ or there exists a proper instance of τ , s.t. $\Sigma_2 \mid \tau$.

We start by considering the case that $\tau_0 \mid \tau$ for a single s-t tgd $\tau_0 \in \Sigma_2$. By the equivalence of Σ_1 and Σ_2 , the implication $\Sigma_1 \mid \tau_0$ holds and we can again apply Lemma 4, i.e., either $\tau_1 \mid \tau_0$ for a single s-t tgd $\tau_1 \in \Sigma_1$ or there exists a proper instance τ_0 of τ_0 , s.t. $\Sigma_1 \mid \tau_0$. Again, we consider first the case that a single s-t tgd is responsible for the implication. Actually, if τ_1 was identical to τ then we had the equivalence $\tau_1 \mid \tau$ and $\tau \mid \tau_1$. Since both τ and τ_1 are reduced w.r.t. Rules 1–3, this would mean (by Lemma 6) that τ_1 and τ are isomorphic. This contradicts our original assumption that τ is not isomorphic to any s-t tgd in Σ_2 . Hence, the case that $\tau_1 \mid \tau_0$ for a single s-t tgd $\tau_1 \in \Sigma_1$ means that τ_1 is different from τ . In total, we thus have $\tau_1 \mid \tau_0$ and $\tau_0 \mid \tau$ and, therefore, $\tau_1 \mid \tau$ for a s-t tgd $\tau_1 \in \Sigma_1 \setminus \{\tau\}$. Hence, τ can be deleted from Σ_1 via Rule 4, which contradicts the normal form of Σ_1 .

It thus remains to consider the cases that there exists a proper instance τ_0 of τ , s.t. $\Sigma_2 \mid \tau_0$ or there exists a proper instance τ_0 of τ_0 , s.t. $\Sigma_1 \mid \tau_0$. We only show that the first one leads to a contradiction. The second case is symmetric. So suppose that there exists a proper instance τ_0 of τ , s.t. $\Sigma_2 \mid \tau_0$. By the equivalence of Σ_1 and Σ_2 , we have $\Sigma_1 \mid \Sigma_2$ and, therefore, also $\Sigma_1 \mid \tau_0$. But then τ can be replaced in Σ_1 by τ_0 via Rule 5. Hence, by Lemma 6, τ can be replaced in Σ_1 by some s-t tgd τ_1 via Rule 5. But this contradicts the assumption that τ is in normal form.

Hence, it is indeed the case that every s-t tgd in Σ_1 is isomorphic to some s-t tgd in Σ_2 and vice versa. We claim that every s-t tgd in Σ_1 is isomorphic to precisely one s-t tgd in Σ_2 and vice versa. Suppose to the contrary that there exists a s-t tgd τ which is isomorphic to two s-t tgds τ_1 and τ_2 in the other set. W.l.o.g., $\tau_1 \in \Sigma_1$ and $\tau_2 \in \Sigma_2$. Clearly, τ_1 and τ_2 are isomorphic since they are both isomorphic to τ . Hence, $\tau_1 \mid \tau_2$ and, therefore, $\Sigma_2 \setminus \{\tau_2\} \mid \tau_2$, i.e., Rule 4 is applicable to Σ_2 , which contradicts the assumption that Σ_2 is in normal form.

We now consider the complexity of computing the normal form of a set of s-t tgds. Of course, the application of any of the Rules 1, 2, 4, and 5 is NP-hard, since they involve CQ answering. However, below we show that if the length of each s-t tgd (i.e., the number of atoms) is bounded by a constant, then the normal form according to Definition 4 can be obtained in polynomial time.

Note that a constant upper bound on the length of the s-t tgds is a common restriction in data exchange since, otherwise, even the most basic tasks like computing a target instance fulfilling all s-t tgds would be intractable.

Theorem 2 Suppose that the length (i.e., the number of atoms) of the s-t tgds under consideration is bounded by some constant b . Then, there exists an algorithm which reduces an arbitrary set Σ of s-t tgds to normal form in polynomial time w.r.t. the total size $|\Sigma|$ of (an appropriate representation of) Σ .

Proof (Sketch) First, the total number of applications of each rule is bounded by the total number of atoms in all s-t tgds in Σ . Indeed, Rule 4 deletes an s-t tgd. The Rules 1, 2, and 5 delete at least one atom from an s-t tgd. Rule 3 splits the conclusion of an s-t tgd in 2 or more parts. Hence, also the total number of applications of Rule 3 is bounded by the total number of atoms in Σ . Finally, the application of each rule is feasible in polynomial time since the most expensive part of these rules is the CQ answering where the length of the CQs is bounded by the number of atoms in a single s-t tgd.

The restriction on the number of atoms in each s-t tgd is used in the above proof only in order to show that each rule application is feasible in polynomial time. The argument that

the total number of rule applications is bounded by the total number of atoms in all s-t tgds Σ applies to any set Σ' of s-t tgds. We thus get:

Corollary 1 The rewrite rule system consisting of Rules 1–5 is terminating, i.e., Given an arbitrary set of s-t tgds, the non-deterministic, exhaustive application of the Rules 1–5 terminates.

It can be shown that the unique normal form produced by our rewrite rules is indeed optimal.

Theorem 3 Let Σ be a set of s-t tgds in normal form. Then Σ is optimal according to Definition 4.

Proof The proof proceeds in three stages:

- (1) Σ is split-reduced. Suppose to the contrary that it is not. Then, there exists a set Σ' with $\Sigma \equiv \Sigma'$, $|\Sigma'| < |\Sigma|$ and $\text{ConSize}(\Sigma) = \text{ConSize}(\Sigma')$. Let Σ'' denote the result of exhaustively applying our rewrite rules to Σ . By Theorem 1, Σ and Σ'' are isomorphic. Hence, we have $\text{ConSize}(\Sigma) = \text{ConSize}(\Sigma'')$ and $|\Sigma''| = |\Sigma|$. An inspection of the Rules 1–5 reveals that they may possibly decrement the value of $\text{ConSize}(\Sigma)$ (by either deleting an atom in the conclusion or deleting an entire s-t tgd) but they never increment the value of $\text{ConSize}(\Sigma)$. By $\text{ConSize}(\Sigma) = \text{ConSize}(\Sigma'')$ together with $\text{ConSize}(\Sigma) = \text{ConSize}(\Sigma')$, we immediately have $\text{ConSize}(\Sigma'') = \text{ConSize}(\Sigma')$. Hence, when transforming Σ into Σ'' , we never decrement $\text{ConSize}(\Sigma)$ and, thus, we never delete an s-t tgd. But then $|\Sigma''| = |\Sigma|$ and, therefore, $|\Sigma''| = |\Sigma'|$, which contradicts the assumption that $|\Sigma'| < |\Sigma|$ holds.
- (2) $\text{ConSize}(\Sigma)$ and $\text{VarSize}(\Sigma)$ are minimal. It is easy to verify that by no application of any of the Rules 1–5, the parameters $\text{ConSize}(\Sigma)$ or $\text{VarSize}(\Sigma)$ can increase, i.e., if a set Σ' of s-t tgds is obtained from some set Σ by an application of one of the Rules 1–5, then $\text{ConSize}(\Sigma') \leq \text{ConSize}(\Sigma)$ and $\text{VarSize}(\Sigma') \leq \text{VarSize}(\Sigma)$. Now let Σ be a set of s-t tgds equivalent to Σ' , and let Σ'' denote the result of exhaustively applying our rewrite rules to Σ . By Theorem 1, Σ and Σ'' are isomorphic. Hence, we have the following relations: $\text{ConSize}(\Sigma) = \text{ConSize}(\Sigma'') \leq \text{ConSize}(\Sigma')$ and also $\text{VarSize}(\Sigma) = \text{VarSize}(\Sigma'') \leq \text{VarSize}(\Sigma')$.
- (3) $|\Sigma|$ and $\text{AntSize}(\Sigma)$ are minimal. Let Σ' be an arbitrary split-reduced set of s-t tgds equivalent to Σ . We first show that $|\Sigma| \leq |\Sigma'|$. Suppose to the contrary that $|\Sigma| > |\Sigma'|$. We derive a contradiction by showing that then Σ is not split-reduced. By (2), we know that $\text{ConSize}(\Sigma) = \text{ConSize}(\Sigma')$ holds. Analogously to the proof of Lemma 7, we can transform Σ into $\bar{\Sigma}$ with $\text{ConSize}(\bar{\Sigma}) = \text{ConSize}(\Sigma)$ simply by choosing an s-t

tgd τ in Σ and inating its conclusion by sufficiently many atoms of the form $P(u_1, \dots, u_k)$. In total, we then have $\bar{\Sigma} \equiv \Sigma$, $\text{ConSize}(\bar{\Sigma}) = \text{ConSize}(\Sigma)$, and $|\bar{\Sigma}| > |\Sigma|$. Hence, Σ is not split-reduced. It remains to prove $\text{AntSize}(\Sigma) = \text{AntSize}(\Sigma')$ as the final inequality. It is easy to verify that the parameter $\text{AntSize}(\Sigma)$ can never increase when one of the Rules 1, 2, 4, or 5 is applied. Moreover, by Lemma 7 we know that Rule 3 is never applicable when we transform a split-reduced set of s-t tgds into normal form. Now let Σ'' denote the normal form of Σ . By Theorem 1, Σ and Σ'' are isomorphic. Hence, we have $\text{AntSize}(\Sigma) = \text{AntSize}(\Sigma'')$ and $\text{AntSize}(\Sigma'') = \text{AntSize}(\Sigma')$.

An important motivation for seeking a conclusion-minimal mapping Σ is to keep the redundancies in the target instance small when using Σ in data exchange. The following theorem establishes that our normal form indeed serves this purpose.

Theorem 4 Let $M = \langle S, T, \Sigma \rangle$ be a schema mapping where Σ is a set of s-t tgds and Σ' is in normal form. Moreover, let Σ'' be another set of s-t tgds, Σ' and Σ'' are equivalent, and let I be an arbitrary source instance. Then, there exists a universal solution $\text{CanSol}(I)$, s.t. $\text{CanSol}(I) \equiv \text{CanSol}(I)$ holds, i.e., the canonical instance produced by Σ is subset minimal up to variable renaming.

Proof It is easy to verify that every application of any of the Rules 1–5 either leaves the corresponding canonical universal solution unchanged or prevents the introduction of some atoms in the canonical universal solution, i.e., let the set of s-t tgds be obtained from some set Σ by an application of one of the Rules 1–5, then there exists a substitution μ s.t. $\text{CanSol}(I)\mu \equiv \text{CanSol}(I)$ holds. The theorem follows by an easy induction argument.

We conclude this section by two remarks on the splitting rule:

- (1) The purpose of the splitting rule is to enable a further simplification of the antecedents of the resulting s-t tgds. Of course, it may happen that no further simplification is possible. As an example, consider a schema mapping $\Sigma = \{ R(x, y) \rightarrow R(y, z) \rightarrow S(x, z) \rightarrow T(z, x) \}$. Splitting yields $\Sigma' = \{ R(x, y) \rightarrow R(y, z) \rightarrow S(x, z); R(x, y) \rightarrow R(y, z) \rightarrow T(z, x) \}$, which cannot be further simplified. In cases like this, one may either “undo” the splitting or simply keep track of s-t tgds with identical (possibly up to variable renaming) antecedents in order to avoid multiple evaluation of the same antecedent by the chase.

(2) Definition 3 gives a “semantical” definition of “split reduced” while the splitting rule is a “syntactical” criterion. The following lemma establishes the close connection between them.

Lemma 7 Let Σ be a split-reduced set of s-t tgds, and let Σ' denote the normal form of Σ . Then, for every possible sequence of rewrite rule applications, this normal form is obtained from Σ without ever applying Rule 3 (i.e., splitting).

Proof Suppose to the contrary that there exists a sequence of rewrite rule applications including the splitting rule on the way from Σ to Σ' . An inspection of the rewrite rules shows that an application of Rule 4 (i.e., deletion of an s-t tgd) is never required as a precondition in order to be able to apply another rule. Hence, w.l.o.g., we may assume that Rule 4 is applied at the very end of the transformation into Σ' , so that Rule 4 does not precede the application of any other rule. Let $\Sigma_0, \dots, \Sigma_n$ with $\Sigma_0 = \Sigma$ and $\Sigma_n = \Sigma'$ denote the sequence of intermediate results along this transformation of Σ into Σ' . Then, there exists an $i \in \{1, \dots, n\}$, s.t. Σ_i is obtained from Σ_{i-1} by an application of Rule 3. Moreover, suppose that this is the first application of Rule 3 along this transformation of Σ into Σ' . Since we are assuming that all applications of Rule 4 occur at the very end of this transformation from Σ to Σ' , we have $|\Sigma_{i-1}| = |\Sigma|$ and, therefore, $|\Sigma_i| > |\Sigma|$. An inspection of the Rules 1, 2, 3, and 5 reveals that they may possibly decrement the value $\text{ConSize}(\Sigma)$ (by deleting an atom in the conclusion via Rule 1 or 5) but they never increment the value $\text{ConSize}(\Sigma)$. Hence, we have $\text{ConSize}(\Sigma_i) = \text{ConSize}(\Sigma)$. We derive a contradiction by constructing a set Σ'' equivalent to Σ_i (and, hence, to Σ), with $\text{ConSize}(\Sigma'') = \text{ConSize}(\Sigma)$ and $|\Sigma''| > |\Sigma|$. In other words, we show that Σ is not split-reduced.

Let τ with $\tau: \varphi(x) \rightarrow \psi(x, y)$ be an arbitrary s-t tgd in Σ_i and let $P(z_1, \dots, z_k)$ with $\{z_1, \dots, z_k\} \cap x \cap y = \emptyset$ be an atom in the conclusion of τ . Clearly, we may add atoms of the form $P(u_1, \dots, u_k)$ for fresh, existentially quantified variables u_1, \dots, u_k to the conclusion without changing the semantics of Σ . Indeed, any such atom could be removed again by our Rule 1 (Core of the conclusion). Then, we transform Σ_i into Σ'' with $\text{ConSize}(\Sigma'') = \text{ConSize}(\Sigma_i)$ simply by inserting the conclusion of τ in Σ_i by sufficiently many atoms of the form $P(u_1, \dots, u_k)$. In total, we then have $\Sigma'' \equiv \Sigma_i \equiv \Sigma$, $\text{ConSize}(\Sigma'') = \text{ConSize}(\Sigma)$, and $|\Sigma''| = |\Sigma_i| > |\Sigma|$. Hence, Σ is not split-reduced, which contradicts the assumption of this lemma.

If a mapping Σ contains redundancies in the sense that one of the Rules 1, 4, 5 is applicable, then the notion of “split-reduced” according to Definition 3 and the non-applicability do not necessarily coincide as the following example illustrates. However, if Rules 1, 4, 5 are not applicable, then Definition 3 is exactly captured by the splitting rule (Rule 3).

Example 9 Consider the set $\Sigma = \{\tau\}$ of s-t tgds with $\tau: P(x_1, x_2) \wedge (y_1, y_2) R(x_1, x_2, y_1) \wedge R(x_1, y_1, y_2)$. On the one hand, Rule 3 is not applicable because the conclusion of τ consists of a single connected component.

On the other hand, Σ may be also simplified (via Rule 1 or Rule 5) to $\tau: P(x_1, x_2) \wedge (y) R(x_1, x_2, y)$. Now let Σ' consist of two “copies” of τ , i.e., $\Sigma' = \{\tau, \tau\}$ with $\tau: P(z_1, z_2) \wedge (y) R(z_1, z_2, y)$. Then we have the equivalence $\Sigma \equiv \Sigma'$. Moreover, $|\Sigma'| > |\Sigma|$ and $\text{ConSize}(\Sigma') = \text{ConSize}(\Sigma)$. Hence, Σ is not split-reduced in the sense of Definition 3.

Lemma 8 Let Σ be a set of s-t tgds, s.t. Σ is reduced w.r.t. Rules 1, 4, 5. Then the following equivalence holds: Σ is split-reduced (according to Definition 3) iff Rule 3 (i.e., splitting) is not applicable.

Proof If Rule 3 is applicable to Σ , then Σ can obviously be transformed into an equivalent set with $|\Sigma'| > |\Sigma|$ and $\text{ConSize}(\Sigma') = \text{ConSize}(\Sigma)$, i.e., Σ is not split-reduced according to Definition 3.

Now suppose that the splitting rule is not applicable to Σ . We have to show that then Σ is split-reduced. Suppose to the contrary that it is not split-reduced, i.e., there exists an equivalent set Σ' with $|\Sigma'| > |\Sigma|$ and $\text{ConSize}(\Sigma') = \text{ConSize}(\Sigma)$. We derive a contradiction by exploiting Theorem 1 (i.e., the uniqueness of the normal form according to Definition 3).

First, we observe that the normal form Σ'' of Σ can be obtained via Rule 2 only. Indeed, by assumption, none of Rules 1, 3, 4, 5 is applicable to Σ . Hence, either Σ already is in normal form (i.e., Rule 2 is not applicable either) or Σ can be simplified via Rule 2. Clearly, an application of Rule 2 does not enable the application of any of the other rules. Hence, Σ'' is obtained by iterated applications of Rule 2 only. Note that Rule 2 has no influence on the cardinality and on the conclusion size of a mapping. Hence, we have $|\Sigma''| = |\Sigma|$ and $\text{ConSize}(\Sigma'') = \text{ConSize}(\Sigma)$.

Second, let us transform Σ' into normal form. By Theorem 1, this normal form is unique up to isomorphism. Hence, w.l.o.g., this normal form of Σ' is Σ'' . As far as the cardinality of the involved mappings is concerned, we have $|\Sigma'| > |\Sigma|$ and $|\Sigma''| = |\Sigma|$. Hence, during the transformation of Σ' into Σ'' , eventually Rule 4 or 5 must be applied thus reducing the conclusion size. Hence, we have $\text{ConSize}(\Sigma') > \text{ConSize}(\Sigma'')$. But this is a contradiction to the above equalities $\text{ConSize}(\Sigma') = \text{ConSize}(\Sigma)$ and $\text{ConSize}(\Sigma'') = \text{ConSize}(\Sigma)$.

4 Extension to target Egds
We now extend our rewrite rule system to schema mappings with both s-t tgds and target egds. Several additional considerations and measures are required to arrive at a unique

normal form and a basis for the s-t tgds optimization also in this case. The outline of this section is as follows:

- (1) We have already seen in Example 5 that the presence of egds may have an effect on the equivalence between two sets of s-t tgds. We shall, therefore, first present a method of “propagating” the effects of the egds into the s-t tgds.
- (2) Splitting has played an important role in all our considerations so far. It will turn out that splitting via Rule 3 as in the tgd-only case is not powerful enough if egds are present. We shall, therefore, present a generalization of the notion of “split-reduced” and of the splitting rule to the case when also egds are present. This will lead to the notion of “egd-split-reduced” mappings.
- (3) The intuition of “egd-split-reduced” mappings is that it is not possible to generate the atoms in the conclusion of some tgd by means of several tgds. The antecedent may thus possibly be left unchanged. It can easily be shown that, in general, there does not exist a unique “egd-split-reduced” normal form. We, therefore, restrict this notion to “antecedent-split-reduced” mappings, i.e., a tgd is replaced by new tgds only if the new tgds have strictly smaller antecedent than the original one. With this concept, we shall manage to prove that there always exists a unique (up to isomorphism) normal form also in the presence of egds.
- (4) Finally, we leave aside the considerations on splitting and concentrate on the optimization of the set of s-t tgds according to the criteria of Sect. 3. We shall show that grouping the s-t tgds by homomorphically equivalent antecedents is the key to any optimization tasks in this area.
- (5) We also look at the operation opposite to splitting: namely, merge of s-t tgds with homomorphically equivalent antecedents. As we will show, unlike the s-t tgds only case, in presence of target dependencies, the splitting of s-t tgds can cause an increase in the total number of conclusion atoms. Hence, for some cases, the merge operation can be a reasonable alternative. We will show, however, that with respect to the unique normal form, the “merged” form of the mappings is no more useful than the “egd-split-reduced” form.

$C(x_1, x_2, x_3) \quad (y_1) P(y_1, x_3, x_2),$
 $C(x_1, x_2, x_2) \quad Q(x_1)$
 $\Sigma_{st} = \{ C(x_1, x_2, x_3) \quad (y_1) P(y_1, x_2, x_3),$
 $C(x_1, x_2, x_3) \quad Q(x_1) \}$
 $\Sigma_t = \{ P(x_1, x_2, x_3) \quad x_2 = x_3 \}$
 We have $\Sigma_{st} \equiv \Sigma_t \equiv \Sigma_{st} \equiv \Sigma_t$. Moreover, both Σ_{st} and Σ_{st} are in normal form w.r.t. the Rules 1–5 from Sect. 4. However, $\Sigma_{st} \equiv \Sigma_{st}$ holds.
 In contrast, the equivalence of two sets of target egds is not influenced by the presence of s-t tgds, as the following lemma shows.
Lemma 9 Suppose that $\Sigma = \Sigma_{st} \cup \Sigma_t$ and $\Upsilon = \Upsilon_{st} \cup \Upsilon_t$ are two logically equivalent sets of s-t tgds and target egds. Then, Σ_t and Υ_t are equivalent.
Proof W.l.o.g. assume that there exists an: $\varphi(x) \in \Sigma_t \cup \Upsilon_t$ s.t. $\Sigma_t \models \varepsilon$. That is, the set $L = At(\varphi(x)) \cup \Sigma_t$ of atoms of the antecedent of φ chased with Σ_t does not satisfy ε . However, it does satisfy Σ_t . Now, consider the pair of instances $\langle L, L \rangle$. Since $L \models \Sigma_t$, $\langle L, L \rangle \models \Sigma$ and $\langle L, L \rangle \models \Upsilon$, which is a contradiction.
 Recall that we are only considering logical equivalence of dependencies here. The study of weaker notions of equivalence [10] which only take attainable target instances into account (which is not the case for the above proof) has been initiated in [24].
 In order to work with logical equivalence, we need a way to test logical implication of mappings. However, since we are now dealing with s-t tgds and egds, the declarative implication criterion from Lemma 4 no longer works. Instead, we take the chase-based procedure by Beeri and Vardi [4] applicable to any embedded dependencies that cannot cause an infinite chase (which is clearly the case when all tgds are s-t tgds).
Lemma 10 [4] Let Σ be a set of acyclic tgds and egds and let φ be either a tgd or an egd. Let $\varphi(x)$ denote the antecedent of φ and let T denote the database obtained by chasing φ with Σ . The variables in x are considered as labeled nulls. Then, $\Sigma \models \delta$ if $T \models \delta$ holds.

Propagating the effect of egds into s-t tgds. An important complication introduced by the egds has already been hinted at in Sect. 1, namely the equivalence of two sets of s-t tgds may be affected by the presence of egds:

Example 10 (Example 5 slightly extended).

$$\Sigma_{st} = \{ C(x_1, x_2, x_3) \quad (y_1, y_2) P(y_1, y_2, y_2) \quad P(y_1, x_2, x_3),$$

Analogously to Rule 4 in Fig. 1, we also need a rule for deleting redundant tgds in the presence of target egds. We shall refer to this rule as the Rule E1 in the rewrite rule system to be constructed in this section, which is specified in Fig. 4. As in the tgd-only case, the primary goal of such a rewrite rule system is the definition of a unique normal form of the s-t tgds—but now taking also the target egds into account. The first step toward this goal is to incorporate the effects of egds into s-t tgds. As we have already pointed out in Sect. 4,

Procedure PROPAGATE
Input: A set of s-t tgds and target egds $\Sigma = \Sigma_{st} \cup \Sigma_t$
Output: Sets of source egds Σ_s and rewritten s-t tgds Σ_{st}^*

1. Set $\Sigma_s = \Sigma_{st}^* = \emptyset$;
2. **for each** s-t tgd $\tau: \varphi(\mathbf{x}) \rightarrow (\exists \mathbf{y})\psi(\mathbf{x}, \mathbf{y})$ in Σ_{st} **do**
 /* (a) non-frozen antecedent database */
 $I := At(\varphi(\mathbf{x}))$;
 /* (b) chase with $\Sigma = \Sigma_{st} \cup \Sigma_t$ */
 $\langle J_S, J_T \rangle := \langle I, \emptyset \rangle^\Sigma$;
 /* (c) transform s-t tgd τ into τ' */
 Let $J^* = core(J_T)$, whereby the terms occurring in J_S
 are considered as constants.
 Let \mathbf{y} be a tuple of all variables from $var(J_T) \setminus var(J_S)$;
 $\tau' := (\bigwedge_{A \in J_S} A) \rightarrow (\exists \mathbf{y}) \bigwedge_{B \in J^*} B$;
 $\Sigma_{st}^* := \Sigma_{st}^* \cup \{\tau'\}$;
 /* (d) generate source egds */
 Compute a substitution λ s.t. $At(\varphi(\mathbf{x}\lambda)) = J_S$;
for each pair of variables $x_j, x_k \in \mathbf{x}$ **do**
 if $x_j\lambda = x_k\lambda$ **then**
 $\Sigma_s := \Sigma_s \cup \{\varphi(\mathbf{x}) \rightarrow x_j = x_k\}$;
end for;

Fig. 3 Procedure propagate

may require the introduction of source egds. Since we only consider source instances containing no variables (and not the recent semantics of [2]), there will be no source chase. The source egds are only meant to capture the failure conditions, which cannot be detected otherwise after the rewriting of the s-t tgds.

In Fig. 3, we present the procedure PROPAGATE, which incorporates, to some extent, the effect of the target egds into the s-t tgds and thereby possibly generates source egds. The idea of this procedure is that, for every s-t tgd we identify all egds that will be applicable whenever it is. Moreover, we want that all equalities enforced by these egds should already be enforced in the s-t tgd. Note that the chase in step 2.(b) is not the usual chase in data exchange. Here, in order to propagate backwards the effect of the target egds, in step 2.(b), we chase the database $At(\varphi(\mathbf{x}))$ with labeled nulls, which instantiate the variables from \mathbf{x} . We assume that this chase always succeeds: the only reason for failure could be constants occurring in s-t tgds. If this is the case, however, the chase is certain to fail on any source instance satisfying the antecedent of τ . Thus, such a τ can be simply replaced by a source egd of the form $\varphi(\mathbf{x}) \rightarrow x_j = x_k$ to rule out instances on which τ would re.

Example 11 We now apply the PROPAGATE procedure to $\Sigma = \Sigma_{st} \cup \Sigma_t$ from Example 10. We start the loop in step 2 of the procedure with the first tgd of Σ_{st} $\tau: C(x_1, x_2, x_3) \rightarrow (y_1, y_2) P(y_1, y_2, y_2) \wedge P(y_1, x_2, x_3)$.

(a) $I := \{C(x_1, x_2, x_3)\}$. (We now consider every x_i as a labeled null).

Rewrite Rules in the Presence of Egds

Rule E1 (General implication).
 $\Sigma \implies \Sigma \setminus \{\tau\}$
 if $\Sigma \setminus \{\tau\} \models \tau$.

Rule E2 (Restriction of an antecedent to endomorphic images).
 $\Sigma \implies (\Sigma \setminus \{\tau\}) \cup \{\tau_1, \dots, \tau_n\}$
 if $\tau: \varphi(\mathbf{x}) \rightarrow (\exists \mathbf{y})\psi(\mathbf{x}, \mathbf{y})$
 and $(\Sigma \setminus \{\tau\}) \cup \{\tau_1, \dots, \tau_n\} \models \tau$
 and for each $i \in \{1, \dots, n\}$
 $\tau_i: \varphi_i(\mathbf{x}_i) \rightarrow (\exists \mathbf{y}_i)\psi_i(\mathbf{x}_i, \mathbf{y}_i)$,
 s.t. $\exists \lambda, At(\varphi(\mathbf{x}\lambda)) = At(\varphi_i(\mathbf{x}_i)) \subseteq At(\varphi(\mathbf{x}))$
 and $\psi_i(\mathbf{x}_i, \mathbf{y}_i) = core(At(\varphi_i(\mathbf{x}_i)))^\Sigma$.

Rule E3 (Implication of atoms in the conclusion).
 $\Sigma \implies (\Sigma \setminus \{\tau\}) \cup \{\tau'\}$
 if $\tau: \varphi(\mathbf{x}) \rightarrow (\exists \mathbf{y})\psi(\mathbf{x}, \mathbf{y})$
 and $\tau': \varphi(\mathbf{x}) \rightarrow (\exists \mathbf{y}')\psi'(\mathbf{x}, \mathbf{y}')$,
 s.t. $At(\psi'(\mathbf{x}, \mathbf{y}')) \subseteq At(\psi(\mathbf{x}, \mathbf{y}))$
 and $(\Sigma \setminus \{\tau\}) \cup \{\tau'\} \models \tau$.

Rule ES (generalized splitting in the presence of egds).
 $\Sigma \implies (\Sigma \setminus \{\tau\}) \cup \{\tau_1, \dots, \tau_n\}$
 if $\tau: \varphi(\mathbf{x}) \rightarrow (\exists \mathbf{y})\psi(\mathbf{x}, \mathbf{y})$
 and $(\Sigma \setminus \{\tau\}) \cup \{\tau_1, \dots, \tau_n\} \models \tau$
 and for each $i \in \{1, \dots, n\}$
 $\tau_i: \varphi_i(\mathbf{x}_i) \rightarrow (\exists \mathbf{y}_i)\psi_i(\mathbf{x}_i, \mathbf{y}_i)$,
 s.t. $\emptyset \subset At(\varphi_i(\mathbf{x}_i)) \subseteq At(\varphi(\mathbf{x}))$
 and $At(\psi_i(\mathbf{x}_i, \mathbf{y}_i\mu_i)) \subseteq At(\psi(\mathbf{x}, \mathbf{y}))$ for a substitution μ_i .

Rule E2-eager (Restriction of an antecedent to subsets).
 $\Sigma \implies (\Sigma \setminus \{\tau\}) \cup \{\tau_1, \dots, \tau_n\}$
 if $\tau: \varphi(\mathbf{x}) \rightarrow (\exists \mathbf{y})\psi(\mathbf{x}, \mathbf{y})$
 and $(\Sigma \setminus \{\tau\}) \cup \{\tau_1, \dots, \tau_n\} \models \tau$
 and for each $i \in \{1, \dots, n\}$
 $\tau_i: \varphi_i(\mathbf{x}_i) \rightarrow (\exists \mathbf{y}_i)\psi_i(\mathbf{x}_i, \mathbf{y}_i)$,
 s.t. $\emptyset \subset At(\varphi_i(\mathbf{x}_i)) \subseteq At(\varphi(\mathbf{x}))$
 and $\psi_i(\mathbf{x}_i, \mathbf{y}_i) = core(At(\varphi_i(\mathbf{x}_i)))^\Sigma$.

Fig. 4 Rewrite rules in the presence of egds

(b) Chasing I , with Σ_{st} yields $I = \{C(x_1, x_2, x_3), P(y_1, y_2, y_2), P(y_1, x_2, x_3), P(y_1, x_3, x_2)\}$. The egd of Σ_t is then applied, resulting in $I = \{C(x_1, x_2, x_2), P(y_1, y_2, y_2), P(y_1, x_2, x_2), P(y_1, x_2, x_2)\}$. Note, that the egd application affected the “source” atom C . Now, the third tgd in Σ_{st} becomes applicable, producing the ultimate instance $J_T = I$, $\Sigma = \{C(x_1, x_2, x_2), P(y_1, y_2, y_2), P(y_1, x_2, x_2), P(y_1, x_2, x_2), Q(x_1)\}$. We got instances $J_S = \{C(x_1, x_2, x_2)\}$ and $J_T = \{P(y_1, y_2, y_2), P(y_1, x_2, x_2), P(y_1, x_2, x_2), Q(x_1)\}$. Core computation of J_T yields $J = \{P(y_1, x_2, x_2), Q(x_1)\}$. The s-t tgd τ is thus transformed into the following $\tau: C(x_1, x_2, x_2) \rightarrow y_1 P(y_1, x_2, x_2) \wedge Q(x_1)$.
 (d) We compute the substitution $\lambda = \{x_3 = x_2\}$, which maps the only atom $C(x_1, x_2, x_3)$ in $\varphi(\mathbf{x})$ onto instance $J_S = \{C(x_1, x_2, x_2)\}$. Hence, we get one source egd $C(x_1, x_2, x_3) \rightarrow x_2 = x_3$.

Finally, after the first iteration of the loop, we have $\Sigma_{st} = \{C(x_1, x_2, x_2) \rightarrow (y_1) P(y_1, x_2, x_2) \wedge Q(x_1)\}$ and $\Sigma_s = \{C(x_1, x_2, x_3) \rightarrow x_2 = x_3\}$. The remaining two iterations

do not change the tgds σ (and thus also introduce no further source egds).

The PROPAGATE procedure never increases the size of the antecedents of s-t tgds. Hence, the cost of the join operations when computing the canonical universal solution is not affected. On the other hand, the size of the conclusions is normally increased by this procedure. Note however that all atoms thus accumulated in the conclusion of some s-t tgd would be generated in a target instance anyway, whenever τ res. We will ultimately discuss the deletion of redundant atoms from the conclusion of the s-t tgds (via a rule similar to Rule 5 from Fig.1). However, for the time being, it is convenient to have all these atoms present. This ensures that dependencies resulting from the PROPAGATE procedure possess the following essential properties.

Lemma 11 Consider a set $\Sigma = \Sigma_{st} \cup \Sigma_t$ of s-t tgds Σ_{st} and target egds Σ_t . Moreover, let (Σ_s, Σ_{st}) denote the result of PROPAGATE(Σ_{st}, Σ_t). Then, the following conditions hold:

- (1) For every s-t tgd $\tau \in \Sigma_{st}$, let I_τ be a database obtained from the antecedent $At(\tau)$ of τ by instantiating the variables of x with fresh distinct constants. Then, the chase of I_τ with $\Sigma_{st} \cup \Sigma_t$ is successful.
- (2) For every source instance I : if $I \models \Sigma_s$ then the chase of I with $\Sigma_{st} \cup \Sigma_t$ fails.

Proof (1) After the successful completion of the chase in step 2.(b) of the PROPAGATE procedure, all necessary unifications in the antecedent relations have been performed. Hence, the instance $At(\varphi(x)), At(\psi(x, y))$ for an s-t tgd $\tau: \varphi(x) \wedge \psi(x, y)$ in Σ_{st} satisfies both Σ_{st} and Σ_t . Freezing the variables in $At(\varphi(x))$ (i.e., considering them as constants) makes no difference.

(2) An inspection of steps 2.(b) and (d) of the PROPAGATE procedure reveals that Σ_s enforces only those equalities which are implied by $\Sigma_{st} \cup \Sigma_t$. Therefore, a violation of Σ_s means that also $\Sigma_{st} \cup \Sigma_t$ is violated.

Lemma 12 The PROPAGATE procedure is correct, i.e.: let $\Sigma = \Sigma_{st} \cup \Sigma_t$ and let (Σ_s, Σ_{st}) result from a call of PROPAGATE(Σ_{st}, Σ_t). Then $\Sigma \models \Sigma$ for $\Sigma = \Sigma_s \cup \Sigma_{st} \cup \Sigma_t$.

Proof The PROPAGATE procedure leaves the set Σ_t unchanged. Moreover, Lemma 11, part (2), implies $\Sigma \models \Sigma_s$. It thus remains to show $\Sigma \models \tau$ for every $\tau \in \Sigma_{st}$ and $\Sigma \models \tau$ for every $\tau \in \Sigma_{st}$. These relationships are proved by inspecting the loop in PROPAGATE (in particular, step 2.b) and checking that the implication criterion of [4] recalled in Lemma 10 is fulfilled.

[$\Sigma \models \tau$] Let $\tau: \varphi(x) \wedge \psi(x, y)$ be obtained by applying the loop of PROPAGATE to some s-t tgd $\sigma \in \Sigma_{st}$ with $\tau: \varphi(x) \wedge \psi(x, y)$. The unifications applied to

$\varphi(x)$ in order to get $\varphi(x)$ are precisely the ones enforced by the chase of $At(\varphi(x))$ with Σ in step 2.(b) of PROPAGATE.

Therefore, chasing $At(\varphi(x))$ with Σ yields the same result as the chase of $\sigma = At(\varphi(x))$ with Σ , namely I^Σ . Hence, the conjunction of the atoms in the set in step 2.(b) is satisfied by I^Σ . Now consider the conclusion $\psi(x, y)$ of τ , which is obtained via core computation from σ . $\psi(x, y)$ is the conjunction of a subset of σ , which is clearly also satisfied by I^Σ .

[$\Sigma \models \tau$] Let $\tau: \varphi(x) \wedge \psi(x, y)$ be an s-t tgd in Σ_{st} and let $\tau: \varphi(x) \wedge \psi(x, y)$ denote the result of applying the loop of PROPAGATE to τ . Consider the (non-frozen) antecedent database $At(\varphi(x))$ of τ . Chasing I with Σ_s comes down to enforcing Σ_s (which transforms $At(\varphi(x))$ into $At(\varphi(x))$) followed by chasing $At(\varphi(x))$ with τ . The result of this chase is $I \models At(\varphi(x)) \wedge At(\psi(x, y))$. Note that $At(\psi(x, y))$ is the core of the chase $At(\varphi(x))$ with Σ . Hence, I must satisfy τ , from which the claim follows, by Lemma 10.

The following property is easy to see and will be helpful subsequently.

Lemma 13 Let $\Sigma = \Sigma_{st} \cup \Sigma_t$, and let (Σ_s, Σ_{st}) denote the result of PROPAGATE(Σ_{st}, Σ_t). Moreover, let $\tau \in \Sigma_{st}$ with $\tau: \varphi(x) \wedge \psi(x, y)$, and let I be a source instance with $I \models At(\varphi(x))$, such that elements of I are instantiated with distinct fresh constants in I .

Then, the chase of I both with $\Sigma = \Sigma_{st} \cup \Sigma_t$ and with $\Sigma = \Sigma_s \cup \Sigma_{st} \cup \Sigma_t$ succeeds. Moreover, $core(I^\Sigma) = core(I^{\Sigma'})$.

Proof By condition (1) of Lemma 11, the chase with Σ and with Σ' succeeds on the frozen antecedent database I . Hence, the chase must also succeed for any subset of the equality $core(J^\Sigma) = core(J^{\Sigma'})$ immediately follows from the correctness of the PROPAGATE procedure, see Lemma 12.

Splitting in the presence of egds The following example illustrates that the splitting rule (i.e., Rule 3 in Fig.1) does not suffice to detect the possibility of splitting a “bigger” tgd into smaller ones in the presence of target egds:

Example 12 Consider the following mapping $\sigma = \{ \tau, \epsilon \}$

$$\begin{aligned} \tau : & S(x, z_1) \wedge S(x, z_2) \wedge R(z_1, y) \wedge Q(z_2, y) \\ \epsilon : & R(x_1, y_1) \wedge Q(x_2, y_2) \wedge y_1 = y_2 \end{aligned}$$

It is easy to check that σ is equivalent to the mapping $\sigma' = \{ \tau_1, \tau_2, \epsilon \}$ with the same target egd and two s-t tgds, each containing only a subset of the antecedent and conclusion atoms of τ :

$$\begin{aligned} \tau_1 : & S(x, z_1) \wedge R(z_1, y) \\ \tau_2 : & S(x, z_2) \wedge Q(z_2, y) \end{aligned}$$

The Rule 3 from Sect 3 does not allow such a splitting, however.

In some sense, the splitting in the above example still has significant similarities with splitting in the absence of egds, namely: The basic idea of distributing the conclusion atoms over several dependencies is still present when target egds have to be taken into account. However, we have to deal with a significant extension here: Without egds, it would never be possible to split the connected component (w.r.t. the existential variables) of the conclusion of a tgdc . As we have seen in the above example, egds may allow us to tear a connected component apart. Moreover, splitting in the presence of egds is not merely distributing atoms of the conclusion of some dependency over several ones. The following example illustrates that we may also have to copy atoms in order to further split the conclusion of a tgdc .

Example 13 Consider the following mapping σ :

$$\begin{aligned} \tau : S(x_1, x_2) \quad S(x_1, x_3) \\ R(x_2, y) \quad P(y, x_2) \quad Q(y, x_3) \\ \epsilon : R(x, y_1) \quad R(x, y_2) \quad y_1 = y_2 \end{aligned}$$

The s-t tgdc can be rewritten in the following way:

$$\begin{aligned} \tau_1 : S(x_1, x_2) \quad S(x_1, x_3) \quad R(x_2, y) \quad Q(y, x_3) \\ \tau_2 : S(x_1, x_2) \quad R(x_2, y) \quad P(y, x_2) \end{aligned}$$

Both τ_1 and τ_2 must contain an R -atom.

We observe that the total number of atoms in all conclusions in the resulting mapping in Example 13 has increased. But compared with the original mapping each conclusion is strictly smaller than the original one. (i.e., is obtained by deletion of at least one atom and possibly the renaming of some variable occurrences). We thus generalize the notion of split-reduced mappings from the s-t tgdc -only case:

Definition 10 Let $\Sigma = \Sigma_{st} \cup \Sigma_t$ be a mapping. The tgdc $\tau : \varphi(\bar{x}, \bar{z}) \quad \psi(\bar{x}, \bar{y}) \quad \Sigma_{st}$ is *egd-split-reduced* if it is not possible to replace it by a set of new s-t tgdc s with antecedent φ_i and conclusion ψ_i , s.t. $\text{At}(\varphi_i) \subseteq \text{At}(\varphi)$ and $|\text{At}(\psi_i)| < |\text{At}(\psi)|$. Σ_{st} is said to be *egd-split-reduced* if each dependency in it is.

The above notion of *egd-split-reduced* mappings generalizes the notion of *split-reduced* mappings from Definition 10 to mappings with target egds. The connection between the two notions of splitting is formalized by the following lemma:

Lemma 14 Let $\Sigma = \Sigma_{st} \cup \Sigma_t$ with $\Sigma_t = \emptyset$ and suppose that Σ_{st} cannot be simplified by any of the Rules 1, 4, and 5 from Fig. 1 (i.e., the rules which would reduce $\text{ConSize}(\Sigma)$ are not applicable). Then the following equivalences hold: Σ_{st} is *egd-split-reduced* iff Σ_{st} is *split-reduced* if Rule 3 (i.e., splitting) cannot be applied.

Proof The second equivalence was already shown in Lemma 8. Below we show that Σ_{st} is *egd-split-reduced* if Rule 3 (i.e., splitting) cannot be applied.

First, suppose that Rule 3 (i.e., splitting) actually can be applied to Σ_{st} . Then, some $\tau \in \Sigma_{st}$ can be replaced by tgdc s τ_1, \dots, τ_n , s.t. the antecedent of each τ_i coincides with the antecedent of τ and the conclusion of each τ_i is a proper subset of the conclusion of τ . Hence, Σ_{st} is not *egd-split-reduced*.

Now suppose that Σ_{st} is not *egd-split-reduced*. We have to show that then Rule 3 can be applied. Suppose to the contrary that Rule 3 cannot be applied. We derive a contradiction by showing that then one of the Rules 1, 4, 5 is applicable to Σ : Since Σ_{st} is not *egd-split-reduced*, there exists a $\tau \in \Sigma_{st}$ with antecedent φ which can be replaced by a set of new tgdc s $\{\tau_1, \dots, \tau_n\}$, s.t. for every i , $\text{At}(\varphi_i) = \text{At}(\varphi)$ and $|\text{At}(\psi_i)| < |\text{At}(\psi)|$ hold, where φ_i and ψ_i , respectively, denote the antecedent and conclusion of τ_i . Moreover, $\Sigma \cup \Sigma'$ holds with $\Sigma' = (\Sigma \setminus \{\tau\}) \cup \{\tau_1, \dots, \tau_n\}$. In particular, $\Sigma \cup \Sigma' \models \tau$. By Lemma 4, then either (1) $\Sigma \cup \Sigma' \models \tau$ holds for some proper instance σ of τ (see Definition 7) or (2) τ is already implied by a single tgdc $\sigma \in \Sigma'$.

In case (1), we clearly also have $\Sigma \cup \Sigma' \models \tau$ for the proper instance σ of τ . But then, by Lemma 5, Rule 5 is applicable to τ , which is a contradiction. It remains to consider case (2). Clearly, $\sigma \in \Sigma \setminus \{\tau\}$ since otherwise σ could be deleted from Σ via Rule 4. So let $\sigma = \tau_j$ for some j . We thus have $\Sigma \cup \Sigma' \models \tau$ with $\Sigma' = (\Sigma \setminus \{\tau\}) \cup \tau_j$ and, therefore, also $\Sigma \cup \Sigma' \models \tau$. Moreover, $|\text{At}(\psi_j)| < |\text{At}(\psi)|$ holds, which implies $\text{ConSize}(\Sigma \cup \Sigma') < \text{ConSize}(\Sigma)$. Now suppose that we transform both Σ and Σ' into the unique (up to isomorphism) normal form $\bar{\Sigma}$ according to Definition 8. By assumption, none of the Rules 1, 3, 4, and 5 is applicable to Σ . Hence, by the same considerations as in the proof of Lemma 8, $\bar{\Sigma}$ is obtained by successive applications of Rule 2, which leaves the conclusions of the tgdc s unchanged. Hence, we have $\text{ConSize}(\bar{\Sigma}) = \text{ConSize}(\Sigma)$. On the other hand, if we transform Σ' into the normal form $\bar{\Sigma}'$, then we never increase the conclusion size. Hence, from the inequality $\text{ConSize}(\bar{\Sigma}') < \text{ConSize}(\bar{\Sigma})$, we may infer $\text{ConSize}(\bar{\Sigma}') < \text{ConSize}(\bar{\Sigma})$, which contradicts the equality that we have just derived.

In Fig. 4 we present the Rule ES, whose exhaustive application obviously transforms any mapping into an *egd-split-reduced* one. Alas, the following example shows that we cannot hope to get a unique *egd-split-reduced* mapping.

Example 14 Consider the schema mapping σ consisting of a single s-t tgdc and a number of target egds:

$$\begin{aligned} S(1, x) \quad S(1, 2) \quad S(y, 2) \quad T(x, y, z) \quad P(x, z) \\ R(y, z) \quad Q(z, v, w) \end{aligned}$$

$T(x, y, z) \quad P(x, z) \quad Q(z, v, w) \quad v = w$
 $T(x, y, z) \quad R(y, z) \quad Q(z, v, w) \quad v = w$
 $T(x, y, z_1) \quad P(x, z_2) \quad Q(w, v, v) \quad z_1 = z_2$
 $T(x, y, z_1) \quad R(y, z_2) \quad Q(w, v, v) \quad z_1 = z_2$

This mapping is not in the egd-split-reduced form, since the antecedent of the s-t tgd can be shrunk by extracting either the P or the Q atom from the conclusion. Σ_{st} and Σ_{st} are two possible transformations of Σ into egd-split-reduced form via the Rule ES.

$\Sigma_{st} = \{ S(1, x) \quad S(1, 2) \quad S(y, 2) \quad T(x, y, z) \quad R(y, z) \quad Q(z, v, w), \quad S(1, x) \quad S(1, 2) \quad P(x, z) \}$ and
 $\Sigma_{st} = \{ S(1, x) \quad S(1, 2) \quad S(y, 2) \quad T(x, y, z) \quad P(x, z) \quad Q(z, v, w), \quad S(1, 2) \quad S(y, 2) \quad R(y, z) \}$

Clearly, the problem in Example 14 is not just due to the definition of the Rule ES. Instead, it is an intrinsic problem of the notion of egd-split-reduced mappings. Apparently, this extent of splitting is too strong. We shall therefore relax the notion of egd-split-reduced. This will be the topic of the next paragraph.

Antecedent-split-reduced mappings In Example 14, we observed that certain antecedent atoms may be freely distributed between several tgds, if the idea of splitting from Sect. 3 is directly adopted in the setting with target constraints. Therefore, in order to arrive at an intuitive definition of a unique normal form, we shift our focus to the antecedents:

Definition 11 Let $\Sigma = \Sigma_{st} \cup \Sigma_t$ be a mapping. The s-t tgd $\tau: \varphi(\bar{x}, \bar{z}) \quad \psi(\bar{x}, \bar{y})$ is antecedent-split-reduced if it is not possible to replace it with a set of new s-t tgds τ_i each having strictly smaller antecedent, i.e., for the antecedents φ_i , we get $|\text{At}(\varphi_i)| < |\text{At}(\varphi)|$. Σ_{st} is said to be antecedent-split-reduced if each dependency in it is.

In order to transform a mapping into an antecedent-split-reduced one, we define the rule E2-eager in Fig. 4. It can be shown that any normal form under a rule rewrite system containing Rule E2-eager is antecedent-split-reduced and vice versa. In this rule, we have to inspect all subsets of the antecedent database of each tgd. Actually, we will show that it suffices to check all subsets of an antecedent $\varphi(x)$, such that φ_i is a proper endomorphic image of $\varphi(x)$. This is what the Rule E2 in Fig. 4 does. Clearly, the number of endomorphic images is, in general, far smaller than the number of all subsets. In particular, we never have to check antecedents smaller than the core of the already present antecedents (here, we mean the core of the conjunctive

query—without distinguishing two groups of variables as we did in the definition of the Rules 1 and 2). The following theorem shows that both the rule E2-eager and the rule E2 exactly capture the notion of antecedent-split-reduced mappings.

Theorem 5 Let $\Sigma = \Sigma_{st} \cup \Sigma_t$ be a schema mapping in which no tgd can be deleted via the Rule E1. Then, the following properties are equivalent:

1. Σ is antecedent-split-reduced.
2. Σ is reduced w.r.t. Rule E2-eager.
3. Σ is reduced w.r.t. Rule E2.

Proof (Sketch) (1) \rightarrow (2) follows directly from the definition of antecedent-split-reduced form.

(2) \rightarrow (3). is trivial, as every proper endomorphic image of a set of atoms $\text{At}(\varphi)$ is a proper subset of $\text{At}(\varphi)$.

(3) \rightarrow (2). Consider an application of the rule E2-eager, in which it substitutes some s-t tgd in Σ with a set of s-t tgds T , s.t. $\Sigma \setminus \{ \tau \} \cup T$. Let now φ_i be the antecedent of τ . By definition of E2-eager, $\text{At}(\varphi_i) \subseteq \text{At}(\varphi)$. We have to show that $\text{At}(\varphi_i)$ is an endomorphic image of $\text{At}(\varphi)$. Suppose to the contrary that it is not. We show that then τ is “super uous” in T : Namely, the property $(\Sigma \setminus \{ \tau \} \cup T) \models \tau$ holds. Of course, if T only contains such tgds which are super uous, then the tgd τ itself can be deleted by the E1 rule, which is a contradiction to the assumption of this theorem. On the other hand, if T is non-empty and contains only tgds whose antecedent is an endomorphic image of φ , then also the E2 Rule is applicable.

To show that τ is super uous in T , consider the following two cases:

- (a) There exists no homomorphism φ_i . Then, τ_i is super uous in T , in a sense that the property $(\Sigma \setminus \{ \tau \} \cup T) \models \tau_i$ holds. Indeed, according to E2-eager, the conclusion of τ_i was created by chasing $\text{At}(\varphi_i)$ with Σ . Since φ_i played no role in that chase, and thus $(\Sigma \setminus \{ \tau \} \cup T) \models \tau_i$ holds, and hence τ_i is indeed super uous.
- (b) There exists a homomorphism φ_i . Let Λ denote the set of all homomorphisms φ_i . Then, we define $T_\varphi = \{ \tau_j \in T \mid \text{At}(\tau_j) = \varphi_i \lambda \text{ for some } \lambda \in \Lambda \}$, i.e., the antecedents of the tgds in T_φ are subsets of φ_i and, at the same time, endomorphic images of φ . Similarly to the previous case, one can show that $(\Sigma \setminus \{ \tau \} \cup T_\varphi) \models \tau_i$ holds. Thus τ_i is super uous.

There is a close connection between antecedent-split-reduced mappings and the split-reduced form from Definition 3:

Lemma 15 Let $\Sigma = \Sigma_{st} \cup \Sigma_t$ with $\Sigma_t = \dots$, and suppose that Σ_{st} cannot be simplified by any of the Rules 1, 4, and 5 from Fig. 1 (i.e., the rules which would reduce Con(S₂) are not applicable). Then, the following equivalence holds: Σ is antecedent-split-reduced if Rule 3 (i.e., splitting) cannot be applied in such a way that the antecedents of all resulting dependencies can be further simplified (by Rule 2).

Proof First, suppose that Rule 3 (i.e., splitting) followed by a simplification of the antecedent of each new tgds is applicable. Then, in the first place, some Σ_{st} can be replaced by tgds τ_1, \dots, τ_n , s.t. the antecedent of each τ_i coincides with the antecedent of ψ and the conclusion of each τ_i is a proper subset of the conclusion of ψ . Moreover, each τ_i can then be transformed via Rule 2 into τ_i , s.t. the antecedent of τ_i is a proper subset of the antecedent of ψ and, therefore, also of τ . Hence, Σ_{st} is not antecedent-split-reduced.

For the opposite direction, suppose that Σ is not antecedent-split-reduced. We have to show that then Rule 3 can be applied followed by applying Rule 2 to each of the new tgds. Since Σ_{st} is not antecedent-split-reduced, there exists a $\tau \in \Sigma_{st}$ with antecedent ψ which can be replaced by a set of new tgds $\{\tau_1, \dots, \tau_n\}$, s.t. for every i , $At(\varphi_i) \subseteq At(\psi)$ and $|At(\psi_i)| < |At(\psi)|$ hold, where φ_i and ψ_i , respectively, denote the antecedent and conclusion of τ_i . Moreover, $\Sigma \equiv \Sigma'$ holds with $\Sigma' = (\Sigma \setminus \{\tau\}) \cup \{\tau_1, \dots, \tau_n\}$.

Analogously to the proof of Theorem 5 we may assume w.l.o.g., that each of the new antecedents is an endomorphic image of ψ . Moreover, we may assume w.l.o.g., that each ψ_i contains only one connected component since otherwise we simply split τ_i further via Rule 3. We claim that for every connected component ϕ of ψ , there is one τ_i , s.t. this connected component corresponds to ϕ . Suppose to the contrary that there is a connected component ϕ of ψ which does not have a corresponding τ_i . Then, we derive a contradiction as follows. The tgds τ obtained from τ by reducing the conclusion ψ to χ is clearly implied by Σ' . Hence, by Lemma 4, either (1) $\Sigma' \models \tau$ holds for some proper instance of τ (see Definition 7) or (2) τ is already implied by a single tgds $\sigma \in \Sigma'$. In case (1), we thus have $\Sigma' \models \tau$ for the proper instance α of τ . But then, also $\Sigma \models \tau$, where τ denotes the tgds obtained from τ by replacing the connected component χ by the conclusion of τ (i.e., a proper instance of χ) and leaving all other connected components unchanged. By Lemma 5, Rule 5 is applicable to $\alpha \in \Sigma$, which is a contradiction. Now consider case (2), i.e., τ is implied by a single tgds $\sigma \in \Sigma'$. Clearly, σ cannot be contained in $\Sigma \setminus \{\tau\}$ since this would mean that the connected component of the conclusion of τ could be deleted from Σ via Rule 5. So suppose that $\sigma = \tau_j$ for some j , i.e., we have $\tau_j \models \tau$. By Lemma 10, this means that the conclusion of τ can be obtained by chasing the antecedent of τ with χ . Note however that χ is a single connected component. Hence, all

of χ are obtained in a single chase step, since otherwise we conclude that also a proper instance of τ is implied by τ_j and we proceed as in case (1). Since τ_j is obtained in a single chase step, the conclusion of τ indeed comprises all of χ . To conclude the proof, recall the above observation that each of the antecedents ψ_i is an endomorphic image of ψ . But then we can indeed apply the Rule 2 in the reverse direction to extend each τ_i to τ_i . Let the resulting tgds be called $\{\bar{\tau}_1, \dots, \bar{\tau}_n\}$. Then, we indeed have that Σ may be replaced by $\{\bar{\tau}_1, \dots, \bar{\tau}_n\}$ via Rule 3, and each $\bar{\tau}_i$ may be further simplified via Rule 2 to τ_i with strictly smaller antecedent.

Most importantly, the notion of antecedent-split-reduced mappings allows us to define a unique (up to isomorphism) normal form of the set of s-t tgds. To this end, we consider the transformation of an arbitrary mapping consisting of s-t tgds and target egds by the PROPAGATE procedure from Fig. 3 followed by exhaustive application of the rules E1 and E2 from Fig. 4. Below we show that the resulting normal form is indeed unique up to isomorphism:

Lemma 16 The rewrite rules E1 and E2 in Fig. 4 are correct, i.e.: Let Σ be a set of dependencies and Σ' be the result of applying one of the rules E1 or E2 to Σ . Then $\Sigma \equiv \Sigma'$.

Proof The correctness follows directly from the fact that a logical implication test is built into the rules E1 and E2.

Theorem 6 Let $\Sigma = \Sigma_{st} \cup \Sigma_t$ and $\Upsilon = \Upsilon_{st} \cup \Upsilon_t$ be two logically equivalent sets consisting of s-t tgds and target egds and let $\Sigma_s, \Sigma_{st}, \Sigma_t$ and $\Upsilon_s, \Upsilon_{st}, \Upsilon_t$ be obtained from Σ respectively Υ by first applying the PROPAGATE procedure and then exhaustively applying the rules E1 and E2 to these mappings. Then $\Sigma \equiv \Upsilon$ holds and Σ_{st} and Υ_{st} are isomorphic.

Proof The equivalence $\Sigma_t \equiv \Upsilon_t$ was shown in Lemma 8. It remains to show that Σ_{st} and Υ_{st} are isomorphic.

Let $\sigma \in \Sigma_{st}$ be an arbitrary s-t tgds in Σ_{st} . We have to show that it has an isomorphic analogue in Υ_{st} (and vice versa). Let $\bar{\Sigma}_\sigma$ denote the set of s-t tgds whose antecedents are the proper subsets of ψ and whose conclusions are obtained by chasing the corresponding antecedent database \mathcal{D} with Σ_t (i.e., we get s-t tgds analogous to those in Rule E2). By Lemma 3, this is the same as chasing these particular source instances with Σ .

By $\Upsilon \models \sigma$, there exists a subset $\Upsilon_s \cup \Upsilon_{st} \cup \Upsilon_t$, s.t. $\Upsilon \models \Sigma_s \cup \bar{\Sigma}_\sigma \cup \Sigma_t \cup \Sigma_{st} \setminus \{\sigma\} \cup \sigma$. We claim that there even exists a set $\Upsilon_\sigma \cup \Upsilon_{st}$ with $\Upsilon_\sigma \cup \Sigma_s \cup \bar{\Sigma}_\sigma \cup \Sigma_t \cup \Sigma_{st} \setminus \{\sigma\} \cup \sigma$, s.t. every $\tau \in \Upsilon_\sigma$ fulfills the following properties:

1. The antecedents $\psi_\tau(x_\tau)$ of τ and $\psi_\sigma(x_\sigma)$ of σ are homomorphically equivalent;

2. there exists a substitution $\varphi_\tau(x_k \lambda) = \varphi_\sigma(x_\sigma)$. That is, the antecedent of τ can be mapped onto the entire antecedent of σ ;
3. τ is not equivalent to any dependency in $\Sigma_{st} \setminus \{\sigma\}$

In order to prove this claim, we start with a set Υ_s $\Upsilon_{st} \Upsilon_t$, s.t. $\Upsilon_s \Sigma_s \bar{\Sigma}_\sigma \Sigma_t \Sigma_{st} \setminus \{\sigma\} \mid \sigma$ and remove all parts from Υ until a subset $\Upsilon_\sigma \Upsilon$ with the desired properties is obtained. It is convenient to write Υ_σ as a short-hand for $\Sigma_s \bar{\Sigma}_\sigma \Sigma_t \Sigma_{st} \setminus \{\sigma\}$.

(a) Eliminate Υ_s from Υ . This is justified by the fact that $\Upsilon_{st} \Upsilon_t \mid \sigma$ holds. Suppose to the contrary that this fact does not hold: that is, let τ be an instance over the schema $S \setminus T$, in which the only non-empty relations are those of the antecedent database $At(\varphi_\sigma(x_\sigma))$ of σ . Then, chasing with $\Upsilon_{st} \Upsilon_t$ leads to an instance $\Upsilon_{st} \Upsilon_t \mid \sigma$, whereas $\Upsilon_s \Upsilon_{st} \Upsilon_t \mid \sigma$. Since source dependencies are only applicable to relations of the source schema, it must hold that Υ_s modifies $At(\varphi_\sigma(x_\sigma))$; otherwise there would be no difference between the two chase results. That is, $\Upsilon_s \mid At(\varphi_\sigma(x_\sigma))$. By Lemma 11, part (2), this means that the chase of $At(\varphi_\sigma(x_\sigma))$ with Υ fails. Thus, also the chase with Υ_σ fails, which contradicts Lemma 11, part (1).

(b) Eliminate Υ_t from Υ . The correctness of this step follows immediately from the equivalence $\Upsilon_s \Sigma_t$ that we showed in Lemma 9.

(c) Eliminate every τ from Υ which is equivalent to some $\sigma \in \Sigma_{st} \setminus \{\sigma\}$. Clearly, after such a reduction, we still have $\Upsilon_\sigma \Sigma \mid \sigma$ with $\Sigma = \Sigma_s \bar{\Sigma}_\sigma \Sigma_t \Sigma_{st} \setminus \{\sigma\}$.

(d) Eliminate from Υ all dependencies with the antecedent $\varphi_i(x_i)$ which is not homomorphically equivalent to the antecedent $\varphi_\sigma(x_\sigma)$ of σ . Indeed, for every s-t τ with the antecedent “more specific” than $\varphi_\sigma(x_\sigma)$, we may conclude that for arbitrary Σ , such that $\Sigma \mid \sigma$, it holds that $\Sigma \setminus \{\tau\} \mid \sigma$. For every τ_j with the antecedent “more general” than $\varphi_\sigma(x_\sigma)$, we have that $\Sigma \setminus \{\sigma\} \mid \tau_j$, and therefore, τ_j is redundant in $\Upsilon \setminus \{\sigma\}$.

(e) Eliminate from Υ all s-t tgds with the antecedent $\varphi_k(x_k)$ such that there exists no variable substitution $\varphi_k(x_k \lambda) = \varphi_\sigma(x_\sigma)$, where $\varphi_\sigma(x_\sigma)$ again denotes the antecedent of σ . First, observe that there are no dependencies whose antecedents under any variable substitution are supersets of $\varphi_\sigma(x_\sigma)$, since they are “more specific” than $\varphi_\sigma(x_\sigma)$ and have therefore been removed in the previous step.

Now consider the substitutions $\varphi_k(x_k \lambda_{ki}) = \varphi_\sigma(x_\sigma)$ and the corresponding s-t tgds τ_k . We claim that the following property holds: For any set of dependencies K such that $\tau_k \in K, K \mid \sigma$ iff $(K \setminus \{\tau_k\}) \mid \sigma$, where K_{τ_k} is

the set of all instantiations of τ_k with $\lambda_{ki} : \tau_{ki} = \varphi_k(x_k \lambda_{ki})$ $y_k \psi(x_k \lambda_{ki}, y_k)$.

The claim follows from the consideration of the implication test by Beeri and Vardi [4]. To chase the antecedent database $At(\varphi_\sigma(x_\sigma))$ of σ , τ_k is instantiated by every λ_{ki} and thus has the same effect in the chase as φ_k . Hence, every τ_k in Υ whose antecedent cannot be projected onto the entire $At(\varphi_\sigma(x_\sigma))$ may be replaced by the respective instantiations.

We now recall that the antecedents of the s-t tgds $\bar{\Sigma}_\sigma \Sigma$ range over all possible subsets of $\varphi_\sigma(x_\sigma)$. That is, for each τ_{ki} with the antecedent $\varphi_k(x_k \lambda_{ki})$, there exists ρ_{ki} with the identical antecedent and with the conclusion $\varphi_k(x_k \lambda_{ki})$ with Σ . Since Σ and Υ are equivalent, we conclude that $\rho_{ki} \mid \tau_{ki}$ and thus $\Sigma \mid K_{\tau_k}$ for every τ_k . Hence, it is indeed allowed to eliminate from all s-t tgds with the antecedent $\varphi_k(x_k)$ such that there exists no variable substitution $\varphi_k(x_k \lambda) = \varphi_\sigma(x_\sigma)$.

After the above elimination steps, Υ is indeed reduced to a set Υ_σ of the desired form. Note that Υ_σ is non-empty. This can be seen as follows: The s-t tgds reduced w.r.t. rules E1 and E2. Hence $\Sigma_s \bar{\Sigma}_\sigma \Sigma_t \Sigma_{st} \setminus \{\sigma\} \mid \sigma$ and, therefore, Υ_σ must be non-empty.

By obvious symmetry reasons, the same holds for any s-t tgd $\tau \in \Upsilon_{st}$ as well: each τ must also have such a corresponding non-empty set $S_\tau \Sigma_{st}$, with elements satisfying the conditions 1–3.

We now construct a directed bipartite graph $G = (V_1, V_2, E)$ as follows: We associate the s-t tgds Σ_{st} and Υ_{st} with the vertices, $s.M_1 = \Sigma_{st}$ and $V_2 = \Upsilon_{st}$. Moreover, whenever $\tau \in T_\sigma$ (resp. $\sigma \in S_\tau$), then there is an edge from τ to σ (resp. from σ to τ).

The conditions 1–3 of T_σ and S_τ translate into the following properties of the graph G :

- a. Every vertex has an incoming edge, since the sets T_σ and S_τ are non-empty.
- b. Cycles in G have length at most 2. Indeed, by property 2, an edge from τ to σ implies that the size of the antecedent of τ is no less than the size of σ . But then all s-t tgds associated with the vertices in a cycle must have antecedents of equal size. By properties 1 and 2, all such antecedents are isomorphic. This means that the conclusions are isomorphic as well, since they are obtained as cores of the chase of isomorphic source instances with equivalent sets of dependencies (procedure PROPAGATE). Vertices that participate in such a two-edge cycle are disconnected from the rest of the graph. This follows from the fact that the corresponding s-t tgds are equivalent, and thus, any other edge would contradict the property 3.

We obtained a graph, of which each vertex should be connected by an incoming path to a cycle (there is only a finite number of vertices, and from each vertex, an infinite

incoming path can be traced, by the property “a”). Consider “c”, this is only possible if each vertex itself belongs to a cycle, and, by “b”, G must consist of connected components of size 2. In total, this means that every s - t egd $\sigma \in \Sigma_{st}$ has an isomorphic counterpart γ_{st} and vice versa.

The question now is how to further simplify the set of s - t tgds. Due to the egds, we could strengthen Rule 5 from Sect. 3 (i.e., deletion of redundant atoms from some conclusion) to the Rule E3 in Fig. 4. Unfortunately, this would again lead to a non-unique normal form as the following example illustrates.

Example 15 Consider the mapping consisting of two s - t tgds and one egd:

$$\begin{aligned} S(x, y) & \quad P(x, z) \quad Q(x, z) \\ S(x, y) & \quad R(x, z) \quad Q(x, z) \\ P(x, z_1) & \quad R(x, z_2) \quad z_1 = z_2 \end{aligned}$$

It is easy to verify that the atom Q can be eliminated by the rule E3 from the conclusion of any of the two tgds, but not from both.

However, if we content ourselves with the simplifications from the s - t tgds only case (i.e., Rules 1–5 from Sect. 3), then we get an intuitive normal form which is simplified to a large extent and which is guaranteed to be unique up to isomorphism. As was mentioned earlier, it is sometimes important in data exchange to arrive at a unique canonical universal solution (this is in particular the case for defining the semantics of queries in a way that the semantics does not depend on the syntax of the dependencies). In these situations, the normal form defined below should be chosen.

Definition 12 Consider a set $\Sigma = \Sigma_{st} \cup \Sigma_t$ of s - t tgds and target egds Σ_{st} and let the result of $\text{PROPAGATE}(\Sigma_{st}, \Sigma_t)$ be denoted by Σ_s, Σ_{st} . Moreover, let Σ_{st} denote the set of s - t tgds resulting from Σ_{st} by exhaustive application of the rules E1, E2 as well as the rules 1–5 from Sect. 3, and let Σ_s denote the result of exhaustive reduction of Σ_s via rule E1. Then, we call $\Sigma_s, \Sigma_{st}, \Sigma_t$ the normal form of Σ .

Theorem 7 Let $\Sigma = \Sigma_{st} \cup \Sigma_t$ and $\Upsilon = \Upsilon_{st} \cup \Upsilon_t$ be equivalent sets consisting of s - t tgds and target egds and let $\Sigma_s, \Sigma_{st}, \Sigma_t$ and $\Upsilon_s, \Upsilon_{st}, \Upsilon_t$ be the corresponding normal forms. Then Σ_{st} and Υ_{st} are isomorphic. Moreover, $\Sigma_t \cup \Upsilon_t$ holds.

Proof The fact that Σ_{st} and Υ_{st} are isomorphic follows immediately from Theorems 5 and 6. The equivalence $\Sigma_t \cup \Upsilon_t$ was proved in Lemma 8.

Homomorphically equivalent components The normal form obtained by the PROPAGATE procedure followed by the

Rules E1 and E2 is not optimal in all respects yet. In particular, both the PROPAGATE procedure and the Rule E2 may have introduced more atoms than needed in the conclusion of s - t tgds. Moreover, by the E2 rule, we may have split the antecedent of tgds into several smaller ones, such that the total number of atoms in the antecedents is increased. Of course, we may now simply apply the rules from Fig. 4 to further simplify the set of s - t tgds. However, in the final part of this section, we want to look in a principled way at further optimizations of the normal form of s - t tgds in the presence of egds. The following concept is crucial.

Definition 13 Let $\Sigma = \Sigma_{st} \cup \Sigma_t$. We say that two tgds τ_1 and τ_2 in Σ_{st} are homomorphically equivalent if their antecedents are. Moreover, we say two sets S of tgds are homomorphically equivalent if the tgds in one set and the tgds in the other set are homomorphically equivalent.

Obviously, homomorphical equivalence is indeed an equivalence relation on Σ_{st} . We refer to the equivalence classes of this relation as the HE-components of Σ_{st} .

We now define a partial order on the HE-components of a set of s - t tgds by considering a “more general” component as greater than a “more specific” one (i.e., there are homomorphisms from the more general one into the “more specific” one but not vice versa). Moreover, we also consider the closure under the greater than relation. Below, we show that the closure of each HE-component is unique up to logical equivalence.

Definition 14 Let $\Sigma = \Sigma_{st} \cup \Sigma_t$ be a mapping, and let $S = \{S_1, \dots, S_m\}$ denote the HE-components of Σ_{st} . We define a partial order as follows: For any pair of indices j, k , we define $S_i \succ S_j$ if for every antecedent $\varphi(x)$ of the tgds in S_i and every antecedent $\psi(z)$ of the tgds in S_j , $\text{At}(\varphi(x)) \wedge \text{At}(\psi(z))$ holds (i.e., there is a homomorphism from $\varphi(x)$ to $\psi(z)$). If $S_i \succ S_j$, S_i is said to be strictly greater than S_j .

For $i \in \{1, \dots, n\}$, we define the closure of S_i above as $\text{Cl}(S_i, \Sigma) = \{ \tau \mid \tau \succ S_j \text{ for some } j \text{ with } S_j \succ S_i \}$.

Example 16 Consider a source schema consisting of a single relation symbol $P(\cdot, \cdot)$ and a schema mapping $\Sigma = \{ \tau_1, \tau_2, \tau_3, \tau_4 \}$, where the τ_i s are defined as follows:

$$\begin{aligned} \tau_1 : & P(x_1, x_2) \quad P(x_2, x_3) \\ & P(y_1, y_2) \quad P(y_2, y_3) \quad P(y_2, y_3) \quad Q(x_1, y_3) \\ \tau_2 : & P(u_1, u_2) \quad P(u_2, u_3) \quad P(u_2, u_3) \quad Q(u_3, u_3) \\ \tau_3 : & P(v_1, v_2) \quad T(v_1, v_2) \\ \tau_4 : & P(v_1, v_1) \quad Q(v_1, v_1) \end{aligned}$$

Intuitively, the binary relation symbol $P(\cdot, \cdot)$ can be thought of as defining edges of a directed graph. Then, the antecedent of the tgd τ_1 consists of two connected components: two

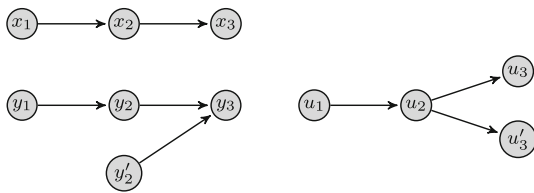


Fig. 5 Antecedents of τ_1 (left) and τ_2 (right), Example 16

paths of length two, one having an additional edge pointing to the peak. The antecedent of the top corresponds to a Y-shaped graph (see Fig. 5). The antecedent of τ_2 consists of a single edge, and the antecedent of τ_3 consists of a single self-loop.

The antecedents of τ_1 and τ_2 have the same cores (a path of length 2) and thus are homomorphically equivalent. Hence, τ_1 and τ_2 are part of the same HE-component. The tgd τ_3 belongs to a different HE-component with $S_2 > S_1$. Indeed, there is a homomorphism sending (v_1, v_2) either to the antecedent of τ_1 or the antecedent of τ_2 , but not vice versa. For the same reason, τ_3 gives rise to yet another HE-component S_3 with $S_1 > S_3$. In total, Σ has three HE-components. As far as the “closure above” is concerned, we thus have $Cl(S_1, \Sigma) = \{\tau_1, \tau_2, \tau_3\}$, $Cl(S_2, \Sigma) = \{\tau_3\}$, and $Cl(S_3, \Sigma) = \Sigma$.

Lemma 17 Let $\Sigma = \Sigma_{st} \cup \Sigma_t$ and $\Upsilon = \Upsilon_{st} \cup \Upsilon_t$ be two logically equivalent mappings. Moreover, let S be an HE-component in Σ_{st} , and let T be an HE-component in Σ_t , s.t. S and T are homomorphically equivalent. Then, $Cl(S, \Sigma) \cup \Sigma_t = Cl(T, \Upsilon) \cup \Upsilon_t$ holds.

Proof By Lemma 9, we have $\Sigma_t \cup \Upsilon_t$. It remains to show that, for every $\tau \in Cl(T, \Upsilon)$, the implication $Cl(S, \Sigma) \cup \Sigma_t \mid \tau$ holds. The implication $Cl(T, \Upsilon_{st}) \cup \Upsilon_t \mid \sigma$ for every $\sigma \in S$ follows by symmetry.

By $\Sigma \cup \Upsilon$, we clearly have $\Sigma \mid \tau$. Let $\varphi(x)$ denote the frozen antecedent of τ and let $I = At(\varphi(x))$. Now consider the result Σ_{st} of chasing I with Σ_{st} : Clearly, only those tgds $\sigma \in \Sigma_{st}$ re, s.t. there is a homomorphism from the antecedent of σ to $\varphi(x)$. These are precisely the tgds in $Cl(S, \Sigma_{st})$. Hence, we have $\Sigma_{st} = I \mid Cl(S, \Sigma)$. But then, by the implication criterion of [4] recalled in Lemma 0, $\Sigma \mid \tau$ holds iff $Cl(S, \Sigma) \cup \Sigma_t \mid \tau$ holds.

The following lemma shows that, unless a mapping contains redundant dependencies, the HE-components of a mapping are in a sense invariant under logical equivalence. Moreover, HE-components may be exchanged between logically equivalent mappings.

Lemma 18 Let $\Sigma = \Sigma_{st} \cup \Sigma_t$ and $\Upsilon = \Upsilon_{st} \cup \Upsilon_t$ be two logically equivalent mappings, s.t. Rule E1 is not applicable to them. Let $S = \{S_1, \dots, S_n\}$ denote the HE-components of Σ_{st} and $T = \{T_1, \dots, T_n\}$ the HE-components of Υ_{st} .

Then, the following properties hold: $\# m$ and for every $S_j \in S$, there exists exactly one T_j , s.t. the tgds in S_j are homomorphically equivalent to the tgds in T_j .

Proof W.l.o.g. suppose that there exists an HE-component S_j of Σ_{st} , which is not homomorphically equivalent to any HE-component T_j of Υ_{st} . By assumption $\Sigma \cup \Upsilon$. Hence, $\Upsilon \mid S_j$. Let $T = \cup\{T_j \mid T_j \in S_j\}$. By the same considerations as in the proof of Lemma 17, only the HE-components in T are used to test the implication $\Upsilon \mid S_j$ via Lemma 10. Hence, we have $\Upsilon \mid S_j$.

On the other hand, also $\Sigma \mid T$. Now define $S = \cup\{S_k \mid S_k \in T_j \text{ for some } T_j \in T\}$. Again, we may conclude $S \mid T$ and, therefore, also $\Sigma \mid S_j$. By assumption, Υ does not contain an HE-component whose tgds are homomorphically equivalent to S_j . Therefore, all HE-components in T are strictly greater than S_j . But then, all HE-components S_k in S are also strictly greater than S_j . Thus, $\Sigma \setminus S_j \mid S_j$. In other words, every dependency σ can be removed from Σ_{st} by the Rule E1, which contradicts the assumption that the E1 Rule is not applicable.

Lemma 19 Let $\Sigma = \Sigma_{st} \cup \Sigma_t$ and $\Upsilon = \Upsilon_{st} \cup \Upsilon_t$ be two logically equivalent mappings, s.t. the E1 Rule is not applicable to them. Moreover, let S be an HE-component in Σ_{st} and let T be an HE-component in Υ_{st} , s.t. S and T are homomorphically equivalent. Then, the logical equivalence $Cl(S, \Sigma) \cup \Sigma_t = Cl(T, \Upsilon) \cup \Upsilon_t$ holds (i.e., we may replace the HE-component S from Σ by the corresponding HE-component T from Υ).

Proof Let $S = \{S_1, \dots, S_n\}$ denote the HE-components of Σ_{st} and $T = \{T_1, \dots, T_n\}$ the HE-components of Υ_{st} . By Lemma 18, we may assume w.l.o.g., that every S_j is homomorphically equivalent to T_j . Now let S and T of this lemma correspond to S_j and T_j for some $j \in \{1 \dots n\}$.

We apply Lemma 17 to all HE-components that are strictly greater than S_j resp. T_j : Let $I = \{i \mid S_i > S_j\}$. Clearly, $I = \{i \mid T_i > T_j\}$. For every $i \in I$, we have $Cl(S_i, \Sigma) \cup \Sigma_t = Cl(T_i, \Upsilon) \cup \Upsilon_t$ by Lemma 17. Then, also $(\cup_{i \in I} Cl(S_i, \Sigma)) \cup \Sigma_t = (\cup_{i \in I} Cl(T_i, \Sigma)) \cup \Upsilon_t$ holds, i.e.: $(Cl(S_j, \Sigma) \setminus S_j) \cup \Sigma_t = (Cl(T_j, \Sigma) \setminus T_j) \cup \Upsilon_t$, i.e., the HE-components strictly greater than S_j and T_j lead to logical equivalence.

Now if we apply Lemma 17 to S_j and T_j , we may conclude $Cl(S_j, \Sigma) \cup \Sigma_t = Cl(T_j, \Upsilon) \cup \Upsilon_t$. By the above considerations, we may exchange $Cl(T_j, \Upsilon)$ in all HE-components that are strictly greater than T_j by the corresponding HE-components from Σ . That is, $Cl(S_j, \Sigma) \cup \Sigma_t = (Cl(S_j, \Sigma) \setminus S_j) \cup T_j \cup \Upsilon_t$. By adding all remaining HE-components of Σ to both sides of the equivalence, we get the desired equivalence $Cl(S, \Sigma) \cup \Sigma_t = Cl(T, \Upsilon) \cup \Upsilon_t$.

HE-components will turn out to be crucial for optimizing the s-t tgds. Indeed, we show that for all optimization criteria

considered here, local optimization inside every HE-component yields a global optimum.

Definition 15 An optimization problem on sets of dependencies is called **sum-minimization problem** if the goal of the optimization is to minimize a function F with the following property: (1) $F(\Sigma) \geq 0$ holds for every set of dependencies Σ and (2) for any two sets of dependencies Σ, Σ' with $\Sigma \sqsubseteq \Sigma'$, we have $F(\Sigma) \leq F(\Sigma')$.

Clearly, all optimization criteria studied here (like cardinality-minimality, antecedent-minimality, conclusion-minimality, and variable-minimality, see Definition 4) are sum-minimization problems.

Definition 16 Let $\Sigma = \Sigma_{st} \sqcup \Sigma_t$ be a mapping, s.t. the E1 Rule is not applicable to it, i.e. Σ contains no s-t tgds that may be deleted. Now consider a sum-minimization problem whose goal is to minimize some function F over sets of s-t tgds.

We say that Σ is **globally optimal** (or simply **optimal**) if, for every mapping $\Upsilon = \Upsilon_{st} \sqcup \Upsilon_t$ with $\Sigma \sqsubseteq \Upsilon$, we have $F(\Sigma) \leq F(\Upsilon)$.

We say that Σ is **locally optimal** if the following conditions are fulfilled: let $\Upsilon = \Upsilon_{st} \sqcup \Upsilon_t$ be an arbitrary mapping with $\Sigma \sqsubseteq \Upsilon$. Moreover, let S be an arbitrary HE-component of Σ and let T be the corresponding HE-component of Υ , s.t. S and T are homomorphically equivalent. Then $F(S) \leq F(T)$ holds.

Theorem 8 Let $\Sigma = \Sigma_{st} \sqcup \Sigma_t$ be a mapping, s.t. the E1 Rule is not applicable to it. Now consider a sum-minimization problem whose goal is to minimize some function F over sets of s-t tgds. Then Σ is globally optimal if it is locally optimal.

Proof Let $\Upsilon = \Upsilon_{st} \sqcup \Upsilon_t$ be an arbitrary mapping with $\Sigma \sqsubseteq \Upsilon$. By Lemma 18, there exist sets of s-t tgds $S = \{S_1, \dots, S_n\}$ and $T = \{T_1, \dots, T_n\}$, s.t. S denotes the set of HE-components of Σ_{st} , T denotes the set of HE-components of Υ_{st} , and for every i , the tgds in S_i are homomorphically equivalent to the tgds in T_i .

First, suppose that Σ is globally optimal. We have to show that then Σ is also locally optimal. Assume to the contrary that $F(S_i) > F(T_i)$ holds for some i . We define $\Sigma' = (\Sigma_{st} \setminus S_i) \sqcup T_i$.

By Lemma 19, $\Sigma \sqsubseteq \Sigma'$. Moreover, since we are considering a sum-minimization problem, we clearly have $F(\Sigma_{st}) = F(\Sigma_{st} \setminus S_i) + F(S_i) > F(\Sigma_{st} \setminus S_i) + F(T_i) = F((\Sigma_{st} \setminus S_i) \sqcup T_i) = F(\Sigma')$. This contradicts the assumption that Σ is globally optimal.

Now suppose that Σ is locally optimal. We have to show that then Σ is also globally optimal. The local optimality implies that $F(S_i) \leq F(T_i)$ holds for every i . Since F defines a sum-minimization problem, we have $F(\Sigma_{st}) = \sum_{i=1}^n F(S_i)$ and $F(\Upsilon_{st}) = \sum_{i=1}^n F(T_i)$. But then also $F(\Sigma_{st}) \leq F(\Upsilon_{st})$ holds, i.e., Σ is globally optimal.

Theorem 8 says, that for the optimization of an HE-component, it does not matter how and if other HE-components have already been optimized. However, this does not mean that one can optimize a single HE-component in isolation. In particular, the closure above must be considered.

As demonstrated by the Examples 14 and 15, aggressive splitting and conclusion optimization lead to a non-unique normal form. In the rest of the section, we consider an operation opposite to splitting: Namely, merging of multiple s-t tgds, to enforce cardinality-minimality. As we will see, also this approach leads to non-unique normal forms. The following theorem contains a property that any merge operation must fulfill:

Theorem 9 Consider a mapping $\Sigma = \Sigma_s \sqcup \Sigma_{st} \sqcup \Sigma_t$ created by the PROPAGATE procedure and additionally, reduced by the rules E1 and E2. Assume that dependency Σ_{st} with the antecedent τ can be substituted by the dependency τ' with the antecedent τ' , such that $\Sigma \setminus \{\tau\} \cup \{\tau'\} \sqsubseteq \Sigma$ holds, and $At(\varphi)$ does not cause a chase failure under φ . Then, φ must coincide (up to isomorphism) with some endomorphic image of φ .

Proof By Theorem 6, exhaustive application of the rules E1 and E2 allows us to obtain a unique normal form of s-t dependencies. Hence, if the mapping $\Sigma \setminus \{\tau\} \cup \{\tau'\}$ is logically equivalent to Σ , it is possible to bring it back in the form isomorphic to Σ by applying the procedure PROPAGATE, followed by the rules E1 and E2.

Since $At(\varphi)$ does not cause a chase failure, we know that PROPAGATE does not affect φ in any way. Moreover, all the remaining rules in $\Sigma \setminus \{\tau\}$ remain unchanged after Σ is transformed by E1 and E2. Hence, it must be the case that one can obtain τ' from τ by (possibly successive) applications of E2, and hence φ' has to be among the endomorphic images of φ .

To achieve cardinality-minimality, we will replace each HE-component with a single tgd. As Theorem 9 suggests, the antecedent of this tgd must contain every antecedent from the original HE-component as an endomorphic image. The following example illustrates that there is no unique minimal “merged” antecedent.

Example 17 Recall the mapping Σ from Example 16, with the HE-components S_1 containing tgds τ_1, τ_2 :

$$\begin{aligned} \tau_1 : & P(x_1, x_2) \wedge P(x_2, x_3) \\ & P(y_1, y_2) \wedge P(y_2, y_3) \wedge P(y_2, y_3) \wedge Q(x_1, y_3) \\ \tau_2 : & P(u_1, u_2) \wedge P(u_2, u_3) \wedge P(u_2, u_3) \wedge Q(u_3, u_3) \end{aligned}$$

Let φ_1, φ_2 denote the antecedents of τ_1 and τ_2 , respectively. Recall the graphical representation of φ_1 and φ_2 that was given in Fig. 5. Obviously, φ_1 and φ_2 are not isomorphic.

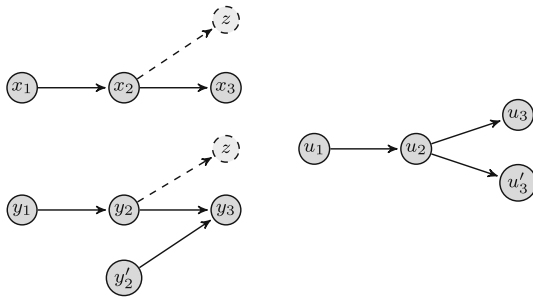


Fig. 6 Possible merges of antecedents φ_2 , Example 17

```

Procedure MERGE
Input. A set  $\Phi$  of homomorphically equivalent CQs;
Output. CQ  $\varphi$  having endomorphism onto each  $\varphi_i \in \Phi$ .
  while  $|\Phi| > 1$  do
    Choose distinct  $\varphi_i, \varphi_j \in \Phi$ .
    Find an endomorphism  $e_i$  for  $\varphi_i$  and  $e_j$  for  $\varphi_j$ , such that
       $e_i(\varphi_i) \cong e_j(\varphi_j)$  and  $|e_i(\varphi_i)|$  is maximized.
    Let  $\lambda$  be a variable renaming  $e_i(\varphi_i) \rightarrow e_j(\varphi_j)$ .
    Set  $\varphi_{ij} := \varphi_i \lambda \wedge \varphi_j$ .
    Set  $\Phi := \Phi \cup \{\varphi_{ij}\} \setminus \{\varphi_i, \varphi_j\}$ .
  od;
  return the element remaining in  $\Phi$ ;

Procedure MERGETGDS
Input. A mapping  $\Sigma = \Sigma_s \cup \Sigma_{st} \cup \Sigma_t$ , a CQ  $\chi$ ;
Output. A tgd  $\tau$  and set  $\Sigma'_s$  of source egds, such that
  the equality  $\Sigma \equiv (\Sigma \setminus \Sigma[\chi]) \cup \Sigma'_s \cup \{\tau\}$  holds.
  /*(a) collect and merge the antecedents of  $\Sigma[\chi]$  */
  Set  $\Phi = \{\varphi_i(\varphi_i(\mathbf{x}_i) \rightarrow (\exists \mathbf{y}_i) \psi(\mathbf{x}_i, \mathbf{y}_i) \in \Sigma_{st}) \wedge \varphi_i \leftrightarrow \chi\}$ ;
   $I := At(MERGE(\Phi))^{\Sigma_s}$ 
  /*(b) initialize  $\tau'$  */
   $J := I^{\Sigma_{st}}$ .
  Let  $\mathbf{y}$  be a tuple of all labeled nulls from  $var(J) \setminus var(I)$ 
   $\tau' := \bigwedge_{A \in I} A \rightarrow (\exists \mathbf{y}) \bigwedge_{B \in J} B$ .
  /*(c) propagate  $\Sigma_t$  through  $\tau'$  */
   $(\Sigma'_s, \Sigma'_{st}) := PROPAGATE(\Sigma \cup \{\tau'\})$ .
  /*(d) return the result */
  Let  $\tau$  be the s-t tgd in  $\Sigma'_{st}$  corresponding to  $\tau'$ .
  return  $(\tau, \Sigma'_s)$ ;
  
```

Fig. 7 Procedures MERGE and MERGETGDS

Now, there are two ways of adding a single edge in order to get a minimum conjunctive query containing both antecedents as its endomorphic images namely, $\varphi_1 \wedge \varphi_2$ and $\varphi_1 \wedge \varphi_2$, see Fig. 6. Clearly, the resulting antecedents φ_1 and φ_2 are not isomorphic.

Example 17 shows that there is no unique optimal way transforming it into some $\varphi_1 \in \Phi$. For every such φ_1 , we have of merging s-t tgds from a single HE-component. Notably, $At(\varphi(x\mu)) \mid \Sigma_s$, so whenever two nulls, $w \in \text{dom}(At(\varphi))$ egds play no role here. On the other hand, an obvious unique way to be unified by Σ_s , necessarily $\varphi\mu = w\mu$ holds. In total, (though hardly optimal) way of merging would be to take also the instance \mathcal{I} created in step (a) of MERGETGDS as conjunction of all antecedents in a HE-component of a mapping endomorphisms onto every $At(\varphi_i)$. Then, also the tgd τ , created in the step (b), is as powerful as τ_1 , as for each $\tau_1 \in \Sigma_{st}[\chi]$, we know that the chase with

We conclude this discussion by presenting a procedure that merges several homomorphically equivalent conjunctive queries in one, of reasonable size and satisfying the condition of Theorem 9. At every iteration, the procedure MERGE takes two conjunctive queries φ_i and φ_j , finds a greatest common (up to isomorphism) endomorphic image in them, which in the worst case is the core, and renames the variables in such a way that it is stitched to along this greatest common endomorphic subquery. The resulting query is thus sure to have an endomorphism φ as well as to φ_j .

This operation is then used in the Procedure MERGETGDS, which produces a s-t tgd to substitute a given HE-component in a mapping Σ . To build the conclusion of such a merged tgd, MERGETGDS uses the PROPAGATE procedure, which chases the merged antecedent with Σ and then takes the conjunction of atoms in the core of the resulting target instance as the conclusion.

Definition 17 Let Σ be a set of s-t tgds, and let χ be a conjunctive query. Then, we write $\Sigma[\chi]$ to denote the HE-component of those tgds in Σ , whose antecedents are homomorphically equivalent to χ .

Theorem 10 Let $\Sigma = \Sigma_s \cup \Sigma_{st} \cup \Sigma_t$ be a mapping consisting of source egds, s-t tgds, and target egds, and Σ_{st} being produced by the PROPAGATE procedure, and let χ be a CQ. Moreover, let (τ, Σ'_s) be the output of MERGETGDS (Σ, χ) . Then, the following equivalence holds: $(\Sigma \setminus \Sigma[\chi]) \cup \Sigma'_s \cup \{\tau\} \equiv \Sigma$.

Proof First, the new s-t tgd τ was created in step (b) of MERGETGDS by chasing with Σ , so $\Sigma \mid \tau$ holds and thus $\Sigma \equiv \Sigma \cup \{\tau\}$. But then, also $\Sigma \equiv \Sigma \cup \Sigma'_s \cup \{\tau\}$ follows by Lemma 12, as both Σ_s and τ were produced by applying PROPAGATE Procedure to $\Sigma \cup \{\tau\}$.

In the other direction, the merged antecedent φ is at least as “powerful” as any of the antecedents φ_i in $\Sigma[\chi]$, in the following sense: whenever a substitution μ for the variables \mathbf{x}_i exists, such that $At(\varphi_i(x_i\lambda)) = 1$, then also for some substitution μ for \mathbf{x} , $At(\varphi(x\mu)) = 1$.

Indeed, suppose that at step (a) of MERGETGDS, the property $MERGE(\Phi) \mid \Sigma$ holds. Then, the claim is immediate, since MERGE is designed to deliver a CQ that satisfies Theorem 9. If, however, MERGE (Φ) has to be updated with Σ_s , the unifications performed by Σ_s do not affect the property that every CQ φ_i in Φ is an endomorphic image of φ . Indeed, let μ be an endomorphism of \mathbf{x} such that $At(\varphi_i(x_i\lambda)) = 1$.

Σ_{st} has produced at least all the conclusion atoms of the conclusion of τ .

The step (c) with PROPAGATE procedure does not affect dependencies other than λ , since, by precondition of theorem, Σ_{st} results from PROPAGATE procedure, and thus, no frozen antecedent database Σ_{st} can cause chase failure under Σ . Moreover, an application of PROPAGATE to τ cannot deteriorate any endomorphism, which makes the antecedent $\varphi(x)$ of τ isomorphic to some $\varphi_i \mid \Phi$. Indeed, suppose this happens and the unification of variables $x_k, x_l \mid x$ cancels some endomorphism. That is, λ is such that $\varphi(x \mid \lambda) = \varphi_i(x_i)$, and $x_k \lambda = x_l \lambda$. Then, in step 2.(d) of the PROPAGATE procedure, the source egd $\varphi(x) \mid x_k = x_l$ is produced, and $\text{At}(\varphi_i) \mid \varepsilon$ must be the case. Hence, by $\Sigma \mid \Sigma \{ \tau \}$ and Claim 2 of Lemma 1, the chase of frozen $\text{At}(\varphi_i)$ with Σ fails, which contradicts Claim 1 of Lemma 1 and the fact that Σ_{st} is the output of PROPAGATE.

Example 18 Recall the tgds τ_1 and τ_2 with the antecedents φ_1 and φ_2 from Example 17. As illustrated by that example, there are two possible ways to merge φ_1 and φ_2 , resulting in two possible merged antecedents $\varphi_1 = \varphi_1 \mid P(x_2, z)$ and $\varphi_2 = \varphi_2 \mid P(y_2, z)$. The corresponding outputs of the procedure MERGETGDS are

$$\tau_1 : P(y_1, y_2) \mid P(y_2, y_3) \mid P(y_2, y_3) \mid P(x_1, x_2) \mid P(x_2, x_3) \mid P(x_2, z) \mid Q(x_1, y_3) \mid Q(x_3, z)$$

and

$$\tau_2 : P(y_1, y_2) \mid P(y_2, y_3) \mid P(y_2, y_3) \mid P(y_2, z) \mid P(x_1, x_2) \mid P(x_2, x_3) \mid Q(x_1, y_3) \mid Q(y_3, z),$$

respectively.

Summary To sum up, the following lessons have been learned from our analysis of the normalization and optimization of s-t tgds in the presence of egds: In contrast to the tgds-only case, we have seen that one has to be very careful with the definition of splitting and optimization so as not to produce a non-unique normal form: If we aim at a strict generalization of the splitting rule from Sec. 3 via the Rule ES in Fig. 4, then there does not exist a unique normal form. This also happens if we aim at a strict generalization of the Rule 5 (deletion of redundant atoms from the conclusion of a tgd) via the Rule E3 in Fig. 4. For most purposes, we therefore consider the transformation of an arbitrary mapping (consisting of s-t tgds and target egds) via the PROPAGATE procedure and exhaustive application of the rules E1 and E2 from Fig. 4 followed by the Rules 1–5 from Sec. 3 as the best choice: The resulting normal form is unique up to isomorphism and incorporates a reasonable amount of splitting and simplification. From the splitting point of view, the resulting normal form is referred to as “antecedent-split-reduced”. This corresponds to a restriction of the splitting rule in the tgds-only

case to those situations where subsequent antecedent simplifications of all resulting dependencies are possible. Such a restriction is justifiable by the fact that one of the main motivations for splitting is indeed to further reduce the antecedents. From the optimization point of view, the Rules 1–5 guarantee that we do not perform worse than in the tgds-only case. But of course, this leaves some additional potential of further optimization in the presence of egds (in particular the Rule E3) unexploited.

We have also identified the HE-components (components of tgds with homomorphically equivalent antecedents) as an important handle for the most common optimization tasks on the s-t tgds (in particular, for all optimization criteria according to Definition 2). We have seen that a global optimum according to the optimization criteria studied here is obtained by locally optimizing the s-t tgds inside each HE-component. In particular, this allowed us to define a simple procedure, which transforms a mapping into an equivalent one with the smallest possible number of s-t tgds. Of course, also in this case, uniqueness is not guaranteed.

We have entirely concentrated on the normalization and optimization of the s-t tgds, while a transformation of the egds has not been considered. Indeed, a normal form of the (source or target) egds is not important for our purposes since we will show in Theorem 11 that the unique (up to isomorphism) canonical universal solution in data exchange only depends on the normalization of the s-t tgds—the equivalence of the source egds and the concrete syntax of the egds are irrelevant.

5 Aggregate queries

As an application for the schema mappings normalization, in this chapter, we discuss the semantics and evaluation of aggregate queries in data exchange, i.e., queries of the form $\text{SELECT } f \text{ FROM } R$, where f is an aggregate operator $\text{min}(R.A)$, $\text{max}(R.A)$, $\text{count}(R.A)$, $\text{count}(_)$, $\text{sum}(R.A)$, $\text{oravg}(R.A)$, and where R is a target relation symbol or, more generally, a conjunctive query over the target schema \mathcal{A} and A is an attribute of R . For this purpose, we first recall some basic notions on query answering in data exchange as well as some fundamental results on aggregate queries from [11].

Definition 18 Let Σ be a schema mapping over the schema \mathcal{S} . Then, the certain answer for a query q over \mathcal{T} and for the source instance I is

$\text{certain}(q, I, W(I)) = \bigcap \{q(J) \mid J \subseteq W(I)\}$,
 where $W(I)$ is the set of possible worlds for I and Σ .

Several proposals can be found in the literature [9, 15, 20, 21] as to which solutions should be taken as possible worlds for $W(I)$. Typical examples are the set of all solutions, the set of universal solutions, the core of the universal solutions, or the CWA-solutions. For conjunctive queries, all these proposals lead to identical results.

Aggregate certain answers Afrati and Kolaitis [1] initiated the study of the semantics of aggregate queries in data exchange. They adopted the notion of aggregate certain answers for inconsistent databases by Arenas et al. [3] to data exchange:

Definition 19 [1] Let query q be of the form $\text{SELECT } f \text{ FROM } R$, where R is a target relation symbol or, more generally, a first-order query over the target schema T , and f is one of the following aggregate operators: $\min(R.A)$, $\max(R.A)$, $\text{count}(R.A)$, $\text{count}()$, $\text{sum}(R.A)$, or $\text{avg}(R.A)$ for some attribute A of R . For all aggregate operators but $\text{count}()$, tuples with a null value in attribute $R.A$ are ignored in the computation.

- $\text{Value } r$ is a possible answer of q w.r.t. I and $W(I)$ if there exists an instance $J \subseteq W(I)$ for which $f(q)(J) = r$.
- $\text{pos}_q(f(q), I, W(I))$ denotes the set of all possible answers of the aggregate query $f(q)$ w.r.t. I and $W(I)$.
- For the aggregate query $f(q)$, the aggregate certain answer $\text{agg-certain}(f, I, W(I))$ w.r.t. I and $W(I)$ is the interval

$$[\text{glb}(\text{pos}_q(f(q), I, W(I))), \text{lub}(\text{pos}_q(f(q), I, W(I)))]$$

where glb and lub stand, respectively, for the greatest lower bound and the least upper bound.

Semantics of aggregate queries via endomorphic images
 A key issue in defining the semantics of queries in data exchange is to define which set of possible worlds should be considered. In [1], Afrati and Kolaitis showed that all previously considered sets of possible worlds yield a trivial semantics of aggregate queries. Therefore, they introduced a new approach via the endomorphic images of the canonical universal solution. Let $\text{Endom}(I, M)$ denote the endomorphic images of the canonical universal solution $\text{CanSol}(I)$, i.e.: $J \in \text{Endom}(I, M)$ if there exists an endomorphism $h: J \rightarrow J$, s.t. $J = h(J)$. As shown in [1], taking $W(I) = \text{Endom}(I, M)$ leads to an interesting and non-trivial semantics of aggregate queries. However, in general, the semantics definition depends on the concrete syntactic representation of the s-t tgds.

Example 19 Consider the source schema $S = \{P\}$, target schema $T = \{R\}$ and the pair of schema mappings $\Sigma_1 =$

S, T, Σ_1 and $M_2 = S, T, \Sigma_2$ with the following s-t tgds:

$$\Sigma_1 = \{ P(x) \rightarrow (\exists y)R(1, x, y) \} \text{ and}$$

$$\Sigma_2 = \{ P(x) \rightarrow (\exists y_1 \dots y_n)R(1, x, y_1) \dots R(1, x, y_n) \}$$

Clearly, M_1 and M_2 are logically equivalent. However, for the source instance $I = \{P(a)\}$, they yield different canonical universal solutions $M_1 = \{R(1, a, y)\}$ and $M_2 = \{R(1, a, y_1), \dots, R(1, a, y_n)\}$. Let A denote the name of the first attribute of R . Then, all of the three aggregate queries $\text{count}(R.A)$, $\text{count}()$, and $\text{sum}(R.A)$ have the range semantics $[1, 1]$ in M_1 and $[1, n]$ in M_2 , i.e.: M_1 admits only one possible world and the three aggregate queries evaluate to 1 in this world. In contrast, M_2 gives rise to a number of possible worlds with $\{R(1, a, y_1), \dots, R(1, a, y_n)\}$ being the biggest one and $\{R(1, a, y)\}$ the smallest. Thus, the three aggregate queries may take values between 1 and n .

In order to eliminate the dependence on the concrete syntactic representation of the s-t tgds, we have defined a new normal form of s-t tgds in Definition 12. Below, we show that we thus get a unique canonical universal solution also in the presence of target tgds.

Theorem 11 Let $M = S, T, \Sigma_{st} \rightarrow \Sigma_t$ be a schema mapping, and let $\Sigma_s \rightarrow \Sigma_{st} \rightarrow \Sigma_t$ be the normal form of $\Sigma_{st} \rightarrow \Sigma_t$. Moreover, let I be a source instance and J the canonical universal solution for I under M obtained via an oblivious chase with Σ_{st} followed by a chase with Σ_t in arbitrary order. Then, J is unique up to isomorphism. We denote J as $\text{CanSol}(I)$.

Proof (Sketch) By Theorem 6, the normal form of the s-t tgds is unique up to isomorphism. Hence, also the result of the oblivious chase with the s-t tgds is unique up to isomorphism. Finally, also the chase with equivalent sets of tgds produces isomorphic canonical universal instances. This property is proved by induction on the length of one of the chase sequences.

To obtain a unique range semantics of the aggregate functions \min , \max , count , $\text{count}()$, sum , and avg , we therefore propose to follow the approach of [1] with the only difference that we take the unique target instance $\text{CanSol}(I)$ from Theorem 11.

6 Conclusion

We have initiated the study of a theory of schema mapping optimization. We have thus presented several natural optimality criteria and a rewrite rule system for transforming any set of s-t tgds into an equivalent optimal one. Recently, several other works have also presented rewrite rules for transforming a set of s-t tgds into an equivalent one with better computational properties. In [23] and [26], the authors aim at the transformation of a set of s-t tgds into an equivalent set Σ , s.t. chasing a source instance with directly

yields the core of the universal solutions of the corresponding data exchange problem. [22], this transformation of s-t tgds is extended to mappings, which comprise also functional dependencies as target dependencies. The transformations in [22,23,26] insert negated atoms and/or inequalities in the antecedents of some s-t tgds so as to block certain forms of applying these s-t tgds in the chase. The goal pursued by these transformations is to avoid the expensive core computation by post-processing of the canonical universal solution and to obtain the core directly as the chase result. Normalization and optimization of the mappings are not in the scope of those transformations.

In order to extend our rewrite rule system to schema mappings including target egds, the most important ingredients of our transformation (namely splitting and simplification of tgds) had to be defined very carefully so as not to destroy the uniqueness of the normal form. We have investigated several forms of splitting and of optimization, and we have identified a rewrite rule system which indeed guarantees to produce a normal form that is again unique up to variable renaming. Finally, we have applied the normalization of schema mappings containing target egds to aggregate queries in data exchange. An implementation of the presented algorithms is freely available from <http://www.dbai.tuwien.ac.at/proj/sm>.

In this paper, we have only considered the optimization of mappings with respect to logical equivalence. As pointed out in [10], weaker notions of equivalence such as “data exchange equivalence” and “conjunctive query equivalence” may sometimes be more appropriate. Unfortunately, many optimization tasks in these relaxed settings are undecidable [25].

As future work, our results should be extended to more expressive schema mappings, including second-order s-t tgds or target tgds. Not surprisingly, a unique normal form via redundancy elimination is not feasible for the target tgds: Consider the set of target tgds $\Sigma = \{P(x) \rightarrow R(x), S(x), R(x) \rightarrow S(x), S(x) \rightarrow R(x)\}$. Now the tgd $P(x) \rightarrow R(x) \wedge S(x) \rightarrow R(x)$ can be either reduced to $P(x) \rightarrow R(x)$ or to $P(x) \rightarrow S(x)$. Of course, even if no unique normal form of the tgds exists, it is still conceivable that one may obtain a unique (up to isomorphism) canonical universal solution via redundancy elimination from the tgds.

Acknowledgments R. Pichler and V. Savenkov acknowledge support by the Vienna Science and Technology Fund (WWTF), project ICT08-032. V. Savenkov is a holder of a scholarship from the European program “Erasmus Mundus External Cooperation Window”. G. Gottlob’s work was supported by EPSRC grant EP/E010865/1 “Schema Mappings and Automated Services for Data Integration and Exchange”. G. Gottlob is the holder of a Royal Society Wolfson Research Merit Award.

References

1. Afrati, F.N., Kolaitis, P.G.: Answering aggregate queries in data exchange. In: Proceedings PODS’08, pp. 129–138. ACM (2008)

2. Arenas, M., Barceló, P., Fagin, R., Libkin, L.: Locally consistent transformations and query answering in data exchange. In: Proceedings PODS’04, pp. 229–240. ACM (2004)
3. Arenas, M., Bertossi, L.E., Chomicki, J., He, X., Raghavan, V., Spinrad, J.: Scalar aggregation in inconsistent databases. *Theor. Comput. Sci.*3(296), 405–434 (2003)
4. Beeri, C., Vardi, M.Y.: A proof procedure for data dependencies. *J. ACM*31(4), 718–741 (1984)
5. Bernstein, P.A., Green, T.J., Melnik, S., Nash, A.: Implementing mapping composition. *VLDB J.*17(2), 333–353 (2008)
6. Bernstein, P.A., Melnik, S.: Model management 2.0: manipulating richer mappings. In: Proceedings SIGMOD’07, pp. 1–12. ACM (2007)
7. Chandra, A.K., Merlin, P.M.: Optimal implementation of conjunctive queries in relational data bases. In: Proceedings STOC’77, pp. 77–90. ACM Press (1977)
8. Fagin, R.: Horn clauses and database dependencies. *J. ACM*29(4), 952–985 (1982)
9. Fagin, R., Kolaitis, P.G., Miller, R.J., Popa, L.: Data exchange: semantics and query answering. *Theor. Comput. Sci.*336(1), 89–124 (2005)
10. Fagin, R., Kolaitis, P.G., Nash, A., Popa, L.: Towards a theory of schema-mapping optimization. In: Proceedings PODS’08, pp. 33–42. ACM (2008)
11. Fagin, R., Kolaitis, P.G., Popa, L.: Data exchange: getting to the core. *ACM Trans. Database Syst.*30(1), 174–210 (2005)
12. Fagin, R., Kolaitis, P.G., Popa, L., Tan, W.-C.: Reverse data exchange: coping with nulls. In: Proceedings PODS’09, pp. 23–32. ACM (2009)
13. Gottlob, G., Pichler, R., Savenkov, V.: Optimization and normalization of schema mappings. Technical Report DBAI-TR-2011-69, Vienna University of Technology (2011)
14. Halevy, A.Y., Rajaraman, A., Ordille, J. J.: Data integration: the teenage years. In: Proceedings VLDB’06, pp. 9–16. ACM (2006)
15. Hernich, A., Schweikardt, N.: Cwa-solutions for data exchange settings with target dependencies. In: Proceedings PODS’07, pp. 113–122. ACM (2007)
16. Imieliński, T., Lipski, W. Jr.: Incomplete information in relational databases. *J. ACM*31(4), 761–791 (1984)
17. Johnson, D.S., Klug, A.C.: Testing containment of conjunctive queries under functional and inclusion dependencies. *J. Comput. Syst. Sci.*28(1), 167–189 (1984)
18. Kolaitis, P.G.: Schema mappings, data exchange, and metadata management. In: Proceedings PODS’05, pp. 61–75. ACM (2005)
19. Lenzerini, M.: Data integration: a theoretical perspective. In: Proceedings PODS’02, pp. 233–246. ACM (2002)
20. Libkin, L.: Data exchange and incomplete information. In: Proceedings PODS’06, pp. 60–69. ACM Press (2006)
21. Libkin, L., Sirangelo, C.: Data exchange and schema mappings in open and closed worlds. In: Proceedings PODS’08, pp. 139–148. ACM (2008)
22. Marnette, B., Mecca, G., Papotti, P.: Scalable data exchange with functional dependencies. *PVLDB*1(1), 105–116 (2010)
23. Mecca, G., Papotti, P., Raunich, S.: Core schema mappings. In: Proceedings SIGMOD’09, pp. 655–668 (2009)
24. Pichler, R., Sallinger, E., Savenkov, V.: Relaxed notions of schema mapping equivalence revisited. In: Proceedings ICDT’11, pp. 90–101. ACM (2011)
25. Sagiv, Y., Yannakakis, M.: Equivalences among relational expressions with the union and difference operators. *J. ACM*27(4), 633–655 (1980)
26. ten Cate, B., Chiticariu, L., Kolaitis, P.G., Tan, W.C.: Laconic schema mappings: computing the core with sql queries. *PVLDB*2(1), 1006–1017 (2009)