



Rewriting the Infinite Chase

Michael Benedikt
Oxford University
Oxford, United Kingdom
michael.benedikt@cs.ox.ac.uk

Maxime Buron
LIRMM, Inria, Univ. of Montpellier
Montpellier, France
maxime.buron@inria.fr

Stefano Germano
Oxford University
Oxford, United Kingdom
stefano.germano@cs.ox.ac.uk

Kevin Kappelmann
Technical University of Munich
Munich, Germany
kevin.kappelmann@tum.de

Boris Motik
Oxford University
Oxford, United Kingdom
boris.motik@cs.ox.ac.uk

ABSTRACT

Guarded tuple-generating dependencies (GTGDs) are a natural extension of description logics and referential constraints. It has long been known that queries over GTGDs can be answered by a variant of the *chase*—a quintessential technique for reasoning with dependencies. However, there has been little work on concrete algorithms and even less on implementation. To address this gap, we revisit *Datalog rewriting* approaches to query answering, where GTGDs are transformed to a Datalog program that entails the same base facts on each base instance. We show that the rewriting can be seen as containing “shortcut” rules that circumvent certain chase steps, we present several algorithms that compute the rewriting by simulating specific types of chase steps, and we discuss important implementation issues. Finally, we show empirically that our techniques can process complex GTGDs derived from synthetic and real benchmarks and are thus suitable for practical use.

PVLDB Reference Format:

Michael Benedikt, Maxime Buron, Stefano Germano, Kevin Kappelmann, and Boris Motik. Rewriting the Infinite Chase. PVLDB, 15(11): 3045–3057, 2022.

doi:10.14778/3551793.3551851

PVLDB Artifact Availability:

The source code, data, and/or other artifacts have been made available at <https://github.com/KRR-Oxford/Guarded-saturation>.

1 INTRODUCTION

Tuple-generating dependencies (TGDs) are a natural extension of description logics and referential constraints, and they are extensively used in databases. For example, they are used in data integration to capture semantic restrictions on data sources, mapping rules between data sources and the mediated schema, and constraints on the mediated schema. A fundamental computational problem in such applications is *query answering under TGDs*: given a query Q , a collection of facts I , and a set of TGDs Σ , find all the answers to Q that logically follow from I and Σ . This problem has long been seen as a key component of a declarative data integration systems

[24, 32], and it also arises in answering querying using views and accessing data sources with restrictions [19, 25, 36].

The *chase* is a quintessential technique for reasoning with TGDs. It essentially performs “forward reasoning” by extending a set of given facts I to a set I' of all facts implied by I and a set of TGDs Σ . To answer a query, one can compute I' using the chase and then evaluate the query in I' . Unfortunately, the chase does not necessarily terminate, and in fact query answering for general TGDs is undecidable. Considerable effort was devoted to identifying classes of TGDs for which query answering is decidable. One line of work has focused on TGDs where the chase terminates; *weakly-acyclic TGDs* [20] are perhaps the best-known such class. Another line of work focused on *guarded TGDs* (GTGDs). GTGDs are interesting since they can capture common constraints used in data integration, and ontologies expressed in variants of *description logic* (DL) [6] can be translated directly into GTGDs. Example 1.1 illustrates the use of GTGDs used in a practical data integration scenario.

EXAMPLE 1.1. The IEC Common Information Model (CIM) is an open model for describing power generation and distribution networks. It is frequently used as a semantic layer in applications that integrate data about power systems [21]. CIM is defined in UML, but its formal semantics has been provided by a translation into an OWL ontology. The domain of CIM is described using *classes* and *properties*, which correspond to unary and binary relations, respectively. Moreover, semantic relationships between classes and properties are represented as OWL axioms, many of which can be translated into GTGDs. A significant portion of CIM describes power distribution equipment using GTGDs such as (1)–(4).

$$\text{ACEquipment}(x) \rightarrow \exists y \text{ hasTerminal}(x, y) \wedge \text{ACTerminal}(y) \quad (1)$$

$$\text{ACTerminal}(x) \rightarrow \text{Terminal}(x) \quad (2)$$

$$\text{hasTerminal}(x, z) \wedge \text{Terminal}(z) \rightarrow \text{Equipment}(x) \quad (3)$$

$$\text{ACTerminal}(x) \rightarrow \exists y \text{ partOf}(x, y) \wedge \text{ACEquipment}(y) \quad (4)$$

Data integration is then achieved by populating the vocabulary using *mappings*, which can be seen queries over the data sources that produce a set of facts called a *base instance*. A key issue in data integration is dealing with incompleteness of data sources. For example, it is not uncommon that one data source mentions two switches sw_1 and sw_2 , while another data source provides information about connected terminals only for switch sw_1 .

$$\text{ACEquipment}(\text{sw}_1) \quad \text{ACEquipment}(\text{sw}_2) \quad (5)$$

$$\text{hasTerminal}(\text{sw}_1, \text{trm}_1) \quad \text{ACTerminal}(\text{trm}_1) \quad (6)$$

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit <https://creativecommons.org/licenses/by-nc-nd/4.0/> to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing info@vldb.org. Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.

Proceedings of the VLDB Endowment, Vol. 15, No. 11 ISSN 2150-8097.
doi:10.14778/3551793.3551851

GTGDs can be used to complete the data. For example, if a user asks to list all pieces of equipment known to the system, both sw_1 and sw_2 will be returned, even though the base instance does not explicitly classify either switch as a piece of equipment. \blacktriangleleft

Even though the chase for GTGDs does not necessarily terminate, query answering for GTGDs is decidable [33]. To prove decidability, one can argue that the result of a chase is *tree-like*—that is, the facts derived by the chase can be arranged into a particular kind of tree. Next, one can develop a *finite representation* of potentially infinite trees. One possibility is to describe the trees using a *finite tree automaton*, so query answering can be reduced to checking automaton emptiness. While theoretically elegant, this method is not amenable to practical use: building the automaton and the emptiness test are both complex and expensive, and such algorithms always exhibit worst-case complexity. Alternatively, one can use *blocking* to identify a tree prefix sufficient for query evaluation. Blocking is commonly used in description logic reasoning [6], and it was later lifted to *guarded logic* [26]. However, blocking was shown to be impractical for query answering: the required tree prefix can be much larger than the base instance I so, as I grows in size, the size of the tree prefix becomes unmanageable.

More promising query answering techniques for GTGDs are based on *Datalog rewriting* [34]. The idea was initially proposed by Marnette [34], and it was later extended to broader classes of TGDs [10, 23] and settings [11]. The main idea is to transform an input set of GTGDs Σ into a set $\text{rew}(\Sigma)$ of *Datalog* rules such that Σ and $\text{rew}(\Sigma)$ entail the same base facts on each base instance. Thus, given a base instance I , instead of computing the chase of I and Σ (which may not terminate), we compute the chase I' of I and $\text{rew}(\Sigma)$. Since *Datalog* rules essentially correspond to existential-free TGDs, I' is always finite and it can be computed using optimized *Datalog* engines. Moreover, Σ and $\text{rew}(\Sigma)$ entail the same base facts on I , so we can answer any *existential-free* conjunctive query (i.e., queries where all variables are answer variables) by evaluating in I' . The restriction to existential-free queries is technical: existentially quantified variables in a query can be matched to objects introduced by existential quantification, and these are not preserved in a *Datalog* rewriting. However, practical queries are typically existential-free since all query variables are usually answer variables.

EXAMPLE 1.2. A *Datalog* program consisting of rules (2)–(3) and (7) is a rewriting of GTGDs (1)–(4).

$$\text{ACEquipment}(x) \rightarrow \text{Equipment}(x) \quad (7)$$

Rule (7) is a logical consequence of GTGDs (1)–(3), and it provides a “shortcut” for the inferences of the other GTGDs. \blacktriangleleft

The advantage of rewriting-based approaches is scalability in the size of the base instance I . Such techniques have been implemented and practically validated in the context of description logics [27, 28], but practical algorithms have not yet been proposed for GTGDs. This raises several theoretical and practical questions.

How to compute the *Datalog* rules needed for completeness? Existing *Datalog* rewriting algorithms often prove their correctness indirectly. For example, completeness of a rewriting algorithm for description logics [28] uses a proof-theoretic argument, which does not provide an intuition about why the algorithm actually works.

Our first contribution is to *relate Datalog rewriting approaches to the chase*. Towards this goal, we introduce the *one-pass* variant of the chase, which we use to develop a general completeness criterion for *Datalog* rewriting algorithms. This, in turn, provides us with a better understanding of how rewriting algorithms work, and it allows us to discover new algorithms in a systematic way.

What does the space of rewriting algorithms look like? Computing the rewriting $\text{rew}(\Sigma)$ usually requires extending Σ with certain logical consequences of Σ . We show that we can select the relevant consequences using different criteria. Some methods require deriving TGDs with existential quantifiers in the head, others generate *Datalog* rules directly, and yet other methods derive logical implications with function symbols. We relate all of these methods to the one-pass chase mentioned earlier, and we provide theoretical worst-case guarantees about their performance.

How do we ensure scalability of rewriting algorithms? Implementations of *Datalog* rewriting algorithms have thus far been mainly considered in the setting of description logics [28, 37]. *To the best of our knowledge, we provide the first look at optimization and implementation of Datalog rewriting algorithms for GTGDs.* We achieve scalability by developing and combining various indexing and redundancy elimination techniques.

How do we evaluate rewriting algorithms? We provide a benchmark for GTGD query answering algorithms, and we use it to evaluate our methods. To the best of our knowledge, this is the first attempt to evaluate query answering techniques for GTGDs.

Summary of contributions. We give an extensive account of *Datalog* rewriting for GTGDs. In particular, we develop a theoretical framework that allows us to understand, motivate, and show completeness of rewriting algorithms. Moreover, we present several concrete algorithms, establish worst-case complexity bounds, and discuss their relationships. We complement this theoretical analysis with a discussion of how to adapt techniques from first-order theorem proving to the setting of GTGDs. Finally, we empirically evaluate our techniques using an extensive benchmark. All proofs and the details of one algorithm are given in the extended version [13]. Our implementation and a more detailed account of our experimental results can be found online [13].

2 RELATED WORK

Answering queries via rewriting has been extensively considered in description logics. For example, queries over ontologies in the DL-Lite family of languages can be rewritten into first-order queries [17], and fact entailment for *SHIQ* ontologies can be rewritten to disjunctive *Datalog* [28]. These techniques provide the foundation for the Ontop [16] and KAON2 [37] systems, respectively.

In the context of TGDs, first-order rewritings were considered in data integration systems with inclusion and key dependencies [15]. *Datalog* rewritings have been considered for GTGDs [34] and their extensions such as frontier-guarded TGDs [11], and nearly frontier-guarded and nearly guarded TGDs [23]. The focus in these studies was to identify complexity bounds and characterize expressivity of TGD classes rather than provide practical algorithms. Existing implementations of query answering for TGDs use first-order rewriting for linear TGDs [47], chase variants for TGDs with terminating chase

[14], chase with blocking for warded TGDs [12], chase with the magic sets transformation for shy TGDs [3], and Datalog rewriting for separable and weakly separable TGDs [48]. These TGD classes are all different from GTGDs, and we are unaware of any attempts to implement and evaluate GTGD rewriting algorithms.

Our algorithms are related to resolution-based decision procedures for variants of guarded logics [18, 22, 49]. Moreover, our characterization of Datalog rewritings is related to a chase variant used to answer queries over data sources with access patterns [4]. Finally, a variant of the one-pass chase from Section 4 was generalized to the broader context of disjunctive GTGDs [30].

3 PRELIMINARIES

In this section, we recapitulate the well-known definitions and notation that we use to formalize our technical results.

TGDs. Let consts , vars , and nulls be pairwise disjoint, infinite sets of *constants*, *variables*, and *labeled nulls*, respectively. A *term* is a constant, a variable, or a labeled null; moreover, a term is *ground* if it does not contain a variable. For α a formula or a set thereof, $\text{consts}(\alpha)$, $\text{vars}(\alpha)$, $\text{nulls}(\alpha)$, and $\text{terms}(\alpha)$ are the sets of constants, free variables, labeled nulls, and terms, respectively, in α .

A *schema* is a set of relations, each of which is associated with a nonnegative integer arity. A *fact* is an expression of the form $R(\vec{t})$, where R is an n -ary relation and \vec{t} is a vector of n ground terms; moreover, $R(\vec{t})$ is a *base fact* if \vec{t} contains only constants. An *instance* I is a finite set of facts, and I is a *base instance* if it contains only base facts. An *atom* is an expression of the form $R(\vec{t})$, where R is an n -ary relation and \vec{t} is a vector of n terms not containing labeled nulls. Thus, each base fact is an atom. We often treat conjunctions as sets of conjuncts; for example, for γ a conjunction of facts and I an instance, $\gamma \subseteq I$ means that each conjunct of γ is contained in I .

A *tuple generating dependency* (TGD) is a first-order formula of the form $\forall \vec{x}[\beta \rightarrow \exists \vec{y} \eta]$, where β and η are conjunctions of atoms, η is not empty, the free variables of β are \vec{x} , and the free variables of η are contained in $\vec{x} \cup \vec{y}$. Conjunction β is the *body* and formula $\exists \vec{y} \eta$ is the *head* of the TGD. We often omit $\forall \vec{x}$ when writing a TGD. A TGD is *full* if \vec{y} is empty; otherwise, the TGD is *non-full*. A TGD is in *head-normal form* if it is full and its head contains exactly one atom, or it is non-full and each head atom contains at least one existentially quantified variable. Each TGD can be easily transformed to an equivalent set of TGDs in head-normal form. A full TGD in head-normal form is a *Datalog rule*, and a *Datalog program* is a finite set of Datalog rules. The *head-width* (hwidth) and the *body-width* (bwidth) of a TGD are the numbers of variables in the head and body, respectively; these are extended to sets of TGDs by taking the maxima over all TGDs. The notion of an instance satisfying a TGD is inherited from first-order logic. A base fact F is *entailed* by an instance I and a finite set of TGDs Σ , written $I, \Sigma \models F$, if $F \in I'$ holds for each instance $I' \supseteq I$ that satisfies Σ .

A *substitution* σ is a function that maps finitely many variables to terms. The domain and the range of σ are $\text{dom}(\sigma)$ and $\text{rng}(\sigma)$, respectively. For γ a term, a vector of terms, or a formula, $\sigma(\gamma)$ is obtained by replacing each free occurrence of a variable x in γ such that $x \in \text{dom}(\sigma)$ with $\sigma(x)$.

Fact Entailment for Guarded TGDs. Fact entailment for general TGDs is semidecidable, and many variants of the *chase* can be used

to define a (possibly infinite) set of facts that is homomorphically contained in each model of a base instance and a set of TGDs.

Fact entailment is decidable for *guarded* TGDs (GTGDs): a TGD $\forall \vec{x}[\beta \rightarrow \exists \vec{y} \eta]$ is guarded if β contains an atom (called a *guard*) that contains all variables of \vec{x} . Note that a guard need not be unique in β . Let Σ be a finite set of GTGDs. We say that a set of ground terms G is Σ -*guarded* by a fact $R(\vec{t})$ if $G \subseteq \vec{t} \cup \text{consts}(\Sigma)$. Moreover, G is Σ -*guarded* by a set of facts I if G is Σ -guarded by some fact in I . Finally, a fact $S(\vec{u})$ is Σ -guarded by a fact $R(\vec{t})$ (respectively a set of facts I) if \vec{u} is Σ -guarded by $R(\vec{t})$ (respectively I).

By adapting the reasoning techniques for guarded logics [5, 46] and referential database constraints [29], fact entailment for GTGDs can be decided by a chase variant that works on tree-like structures. A *chase tree* T consists of a directed tree, one tree vertex that is said to be *recently updated*, and a function mapping each vertex v in the tree to a finite set of facts $T(v)$. A chase tree T can be transformed to another chase tree T' in the following two ways.

- One can apply a *chase step* with a GTGD $\tau = \forall \vec{x}[\beta \rightarrow \exists \vec{y} \eta]$ in head-normal form. The precondition is that there exist a vertex v in T and a substitution σ with domain \vec{x} such that $\sigma(\beta) \subseteq T(v)$. The result of the chase step is obtained as follows.
 - If τ is full (and thus η is a single atom), then chase tree T' is obtained from T by making v recently updated in T' and setting $T'(v) = T(v) \cup \{\sigma(\eta)\}$.
 - If τ is not full, then σ is extended to a substitution σ' that maps each variable in \vec{y} to a labeled null not occurring in T , and chase tree T' is obtained from T by introducing a fresh child v' of v , making v' recently updated in T' , and setting $T(v') = \sigma'(\eta) \cup \{F \in T(v) \mid F \text{ is } \Sigma\text{-guarded by } \sigma'(\eta)\}$.
- One can apply a *propagation step* from a vertex v to a vertex v' in T . Chase tree T' is obtained from T by making v' recently updated in T' and setting $T'(v') = T(v') \cup S$ for some nonempty set S satisfying $S \subseteq \{F \in T(v) \mid F \text{ is } \Sigma\text{-guarded by } T(v')\}$.

A *tree-like chase sequence* for a base instance I and a finite set of GTGDs Σ in head-normal form is a finite sequence of chase trees T_0, \dots, T_n such that T_0 contains exactly one *root vertex* r that is recently updated in T_0 and $T_0(r) = I$, and each T_i with $0 < i \leq n$ is obtained from T_{i-1} by a chase step with some $\tau \in \Sigma$ or a propagation step. For each vertex v in T_n and each fact $F \in T_n(v)$, this sequence is a *tree-like chase proof of F from I and Σ* . It is well known that $I, \Sigma \models F$ if and only if there exists a tree-like chase proof of F from I and Σ (e.g., [33]). Example 4.3 in Section 4 illustrates these definitions. One can decide $I, \Sigma \models F$ by imposing an upper bound on the size of chase trees that need to be considered [33].

Rewriting. A *Datalog rewriting* of a finite set of TGDs Σ is a Datalog program $\text{rew}(\Sigma)$ such that $I, \Sigma \models F$ if and only if $I, \text{rew}(\Sigma) \models F$ for each base instance I and each base fact F . If Σ contains GTGDs only, then a Datalog rewriting $\text{rew}(\Sigma)$ is guaranteed to exist (which is not the case for general TGDs). Thus, we can reduce fact entailment for GTGDs to Datalog reasoning, which can be solved using highly optimized Datalog techniques [1, 37]. For example, given a base instance I , we can compute the *materialization* of $\text{rew}(\Sigma)$ on I by applying the rules of $\text{rew}(\Sigma)$ to I up to a fixpoint. This will compute precisely all base facts entailed by $\text{rew}(\Sigma)$ (and thus also by Σ) on I , and it can be done in time polynomial in the size of I .

Encoding Existentials by Function Symbols. It is sometimes convenient to represent existentially quantified values using functional terms. In such cases, we use a slightly modified notions of terms, atoms, and rules. It will be clear from the context which definitions we use in different parts of the paper.

We adjust the notion of a term as either a constant, a variable, or an expression of the form $f(\vec{t})$ where f is an n -ary function symbol and \vec{t} is a vector of n terms. The notions of ground terms, (base) facts, and (base) instances, and atoms are the same as before, but they use the modified notion of terms. A rule is a first-order implication of the form $\forall \vec{x}[\beta \rightarrow H]$ where β is a conjunction of atoms whose free variables are \vec{x} , and H is an atom whose free variables are contained in \vec{x} ; as for TGDs, we often omit $\forall \vec{x}$. A rule thus contains no existential quantifiers, but its head contains exactly one atom that can contain function symbols. Also, a Datalog rule, a function-free rule, and a full TGD in head-normal form are all synonyms. Finally, a base fact still contains only constants.

Skolemization allow us to replace existential quantifiers in TGDs by functional terms. Specifically, let $\tau = \forall \vec{x}[\beta \rightarrow \exists \vec{y} \eta]$, and let σ be a substitution defined on each $y \in \vec{y}$ as $\sigma(y) = f_{\tau,y}(\vec{x})$ where $f_{\tau,y}$ is a fresh $|\vec{x}|$ -ary Skolem symbol uniquely associated with τ and y . Then, the Skolemization of τ produces rules $\forall \vec{x}[\beta \rightarrow \sigma(\eta)]$ for each atom $H \in \eta$. Moreover, the Skolemization Σ' of a finite set of TGDs Σ is the union of the rules obtained by Skolemizing each $\tau \in \Sigma$. It is well known that $I, \Sigma \models F$ if and only if $I, \Sigma' \models F$ for each base instance I and each base fact F .

Unification. A unifier of atoms A_1, \dots, A_n and B_1, \dots, B_n is a substitution θ such that $\theta(A_i) = \theta(B_i)$ for $1 \leq i \leq n$. Such θ is a *most general unifier* (MGU) if, for each unifier σ of A_1, \dots, A_n and B_1, \dots, B_n , there exists a substitution ρ such that $\sigma = \rho \circ \theta$ (where \circ is function composition). An MGU is unique up to variable renaming if it exists, and it can be computed in time $O(\sum_{i=1}^n |A_i| + |B_i|)$ where $|A_i|$ and $|B_i|$ are the encoding sizes of A_i and B_i [40, 41].

4 CHASE-BASED DATALOG REWRITING

Our objective is to develop rewriting algorithms that can handle complex GTGDs. Each algorithm will derive Datalog rules that provide “shortcuts” in tree-like chase proofs: instead of introducing a child vertex v' using a chase step with a non-full GTGD at vertex v , performing some inferences in v' , and then propagating a derived fact F back from v' to v , these “shortcuts” will derive F in one step without having to introduce v' . The main question is how to derive all “shortcuts” necessary for completeness while keeping the number of derivations low. In this section we lay the technical foundations that will allow us to study different strategies for deriving “shortcuts” in Section 5. We show that, instead of considering arbitrary chase proofs, we can restrict our attention to chase proofs that are *one-pass* according to Definition 4.1. Then, we identify the parts of such proofs that we need to be able to circumvent using “shortcuts”. Finally, we present sufficient conditions that guarantee completeness of rewriting algorithms. We start by describing formally the structure of tree-like chase proofs.

Definition 4.1. A tree-like chase sequence T_0, \dots, T_n for a base instance I and a finite set of GTGDs Σ in head-normal form is one-pass if, for each $0 < i \leq n$, chase tree T_i is obtained by applying one of the following two steps to the recently updated vertex v of T_{i-1} :

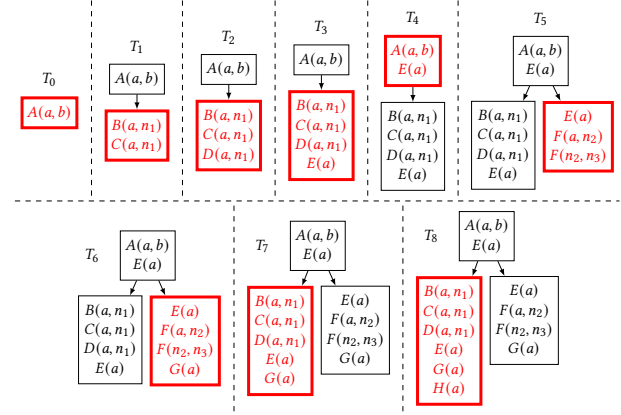


Figure 1: Tree-Like Chase Sequence for Example 4.3

- a propagation step copying exactly one fact from v to its parent, or
- a chase step with a GTGD from Σ provided that no propagation step from v to the parent of v is applicable.

Thus, each step in a tree-like chase sequence is applied to a “focused” vertex; steps with non-full TGDs move the “focus” from a parent to a child, and propagation steps move the “focus” in the opposite direction. Moreover, once a child-to-parent propagation takes place, the child cannot be revisited in further steps. Theorem 4.2 states a key property about chase proofs for GTGDs: whenever a proof exists, there exists a one-pass proof too. Example 4.3 illustrates important aspects of Definition 4.1 and Theorem 4.2.

THEOREM 4.2. For each base instance I , each finite set of GTGDs Σ in head-normal form, and each base fact F such that $I, \Sigma \models F$, there exists a one-pass tree-like chase proof of F from I and Σ .

EXAMPLE 4.3. Let $I = \{A(a, b)\}$ and let Σ contain GTGDs (8)–(13).

$$A(x_1, x_2) \rightarrow \exists y B(x_1, y) \wedge C(x_1, y) \quad (8)$$

$$C(x_1, x_2) \rightarrow D(x_1, x_2) \quad (9)$$

$$B(x_1, x_2) \wedge D(x_1, x_2) \rightarrow E(x_1) \quad (10)$$

$$A(x_1, x_2) \wedge E(x_1) \rightarrow \exists y_1, y_2 F(x_1, y_1) \wedge F(y_1, y_2) \quad (11)$$

$$E(x_1) \wedge F(x_1, x_2) \rightarrow G(x_1) \quad (12)$$

$$B(x_1, x_2) \wedge G(x_1) \rightarrow H(x_1) \quad (13)$$

A tree-like chase sequence for I and Σ is shown in Figure 1, and it provides a proof of the base fact $H(a)$ from I and Σ . The recently updated vertex of each chase tree is shown in red. We denote the root vertex by r , and its left and right children by v_1 and v_2 , respectively. The step producing T_7 from T_6 does not satisfy the requirements of one-pass chase: it propagates the fact $G(a)$ from v_2 to v_1 , where the latter is a “sibling” of the former.

To obtain a one-pass chase sequence, we could try to “slow down” the propagation of $G(a)$: we first propagate $G(a)$ from v_2 to r , and then from r to v_1 . The former step is allowed in one-pass chase, but the latter step is not: once we leave the subtree rooted at v_1 , we are not allowed to revisit it later. Note, however, that $B(a, n_1)$ and $G(a)$ must occur jointly in a vertex of a chase tree in order to derive $H(a)$. Moreover, note that no reordering of chase steps will

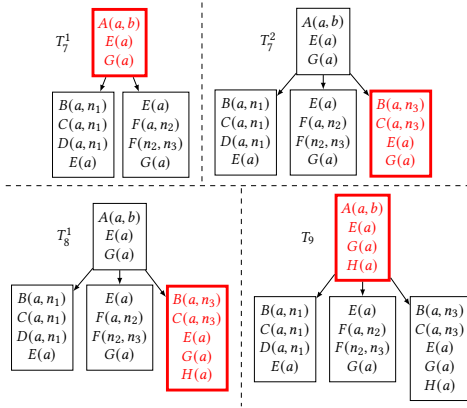


Figure 2: One-Pass Chase Sequence Obtained from Figure 1

derive $H(a)$: we must first produce v_1 to be able to derive v_2 , and we must combine $G(a)$ from v_2 and $B(a, n_1)$ from v_1 .

The solution, which is used in the proof of Theorem 4.2, is to replace propagation to the child by “regrowing” the entire subtree. In our example, we replace the steps producing T_7 and T_8 with the steps shown in Figure 2. Chase tree T_7^1 is obtained from T_6 by propagating $G(a)$ from v_2 to r . Then, instead of propagating $G(a)$ from r to v_1 , a new vertex v_3 is created in T_7^2 by reapplying (8) and fact $G(a)$ is pushed to v_3 as part of the chase step with a non-full GTGD. This allows $H(a)$ to be derived in vertex v_3 of T_8^1 .

Fact $D(n_3)$ can be derived in vertex v_3 , but this is not needed to prove $H(a)$. Moreover, our chase is *oblivious* [33]: a non-full TGD can be applied to the same facts several times, each time introducing a fresh vertex and fresh labeled nulls. The number of children of a vertex is thus not naturally bounded, and our objective is not to apply all chase steps exhaustively to obtain a universal model of Σ . Instead, we are interested only in *chase proofs*, which must only contain steps needed to demonstrate entailment of a specific fact. ◀

One-pass chase proofs are interesting because they can be decomposed into *loops* as described in Definition 4.4.

Definition 4.4. For T_0, \dots, T_n a one-pass tree-like chase sequence for some I and Σ , a loop at vertex v with output fact F is a subsequence T_i, \dots, T_j with $0 \leq i < j \leq n$ such that

- T_{i+1} is obtained by a chase step with a non-full GTGD,
- T_j is obtained by a propagation step that copies F , and
- v is the recently updated vertex of both T_i and T_j .

The length of the loop is defined as $j - i$.

EXAMPLE 4.5. Subsequence T_0, T_1, T_2, T_3, T_4 of the chase trees from Example 4.3 is a loop at the root vertex r with output fact $E(a)$: chase tree T_1 is obtained by applying a non-full GTGD to r , and chase tree T_4 is obtained by propagating $E(a)$ back to r . Analogously, T_4, T_5, T_6, T_7^1 is another loop at r with output fact $G(a)$. Finally, T_7^1, T_7^2, T_8^1, T_9 is a loop at r with output fact $H(a)$. ◀

Thus, a loop is a subsequence of chase steps that move the “focus” from a parent to a child vertex, perform a series of inferences in the child and its descendants, and finally propagate one fact back to

the parent. If non-full TGDs are applied to the child, then the loop can be recursively decomposed into further loops at the child. The properties of the one-pass chase ensure that each loop is finished as soon as a fact is derived in the child that can be propagated to the parent, and that the vertices introduced in the loop are not revisited at any later point in the proof. In this way, each loop at vertex v can be seen as taking the set $T_i(v)$ as input and producing the output fact F that is added to $T_j(v)$. This leads us to the following idea: for each loop with the input set of facts $T_i(v)$, a rewriting should contain a “shortcut” Datalog rule that derives the loop’s output.

EXAMPLE 4.6. One can readily check that rules (14)–(16) provide “shortcuts” for the three loops identified in Example 4.5.

$$A(x_1, x_2) \rightarrow E(x_1) \quad (14)$$

$$A(x_1, x_2) \wedge E(x_1) \rightarrow G(x_1) \quad (15)$$

$$A(x_1, x_2) \wedge G(x_1) \rightarrow H(x_1) \quad (16)$$

Moreover, these are all relevant “shortcuts”: the union of rules (14)–(16) and the Datalog rules from Example 4.3—that is, rules (9), (10), (12), and (13)—is a rewriting of the set Σ from Example 4.1. ◀

These ideas are formalized in Proposition 4.7, which will provide us with a correctness criterion for our algorithms.

PROPOSITION 4.7. A Datalog program Σ' is a rewriting of a finite set of GTGDs Σ in head-normal form if

- Σ' is a logical consequence of Σ ,
- each Datalog rule of Σ is a logical consequence of Σ' , and
- for each base instance I , each one-pass tree-like chase sequence T_0, \dots, T_n for I and Σ , and each loop T_i, \dots, T_j at the root vertex r with output fact F , there exist a Datalog rule $\beta \rightarrow H \in \Sigma'$ and a substitution σ such that $\sigma(\beta) \subseteq T_i(r)$ and $\sigma(H) = F$.

Intuitively, the first condition ensures soundness: rewriting Σ' should not derive more facts than Σ . The second condition ensures that Σ' can mimic direct applications of Datalog rules from Σ at the root vertex r . The third condition ensures that Σ' can reproduce the output of each loop at vertex r using a “shortcut” Datalog rule.

5 REWRITING ALGORITHMS

We now consider ways to produce “shortcut” Datalog rules satisfying Proposition 4.7. In Subsection 5.1 we present the ExbDR algorithm that manipulates GTGDs directly, and in Subsections 5.2 and 5.3 we present the SkDR and HypDR algorithms, respectively, that manipulate rules obtained by Skolemizing the input GTGDs. All of these algorithms can produce intermediate GTGDs/rules that are not necessarily Datalog rules. In the online appendix [13] we present the FullDR algorithm that manipulates GTGDs, but derives only Datalog rules. However, the performance of FullDR proved to not be competitive, so we do not discuss it any further here.

Each algorithm is defined by an inference rule Inf that can be applied to several TGDs/rules to derive additional TGDs/rules. For simplicity, we use the same name for the rule and the resulting algorithm. Given a set of GTGDs Σ , the algorithm applies Inf to (the Skolemization of) Σ as long as possible and then returns all produced Datalog rules. This process, however, can derive a large

number of TGDs/rules, so it is vital to eliminate TGDs/rules whenever possible. We next define notions of *redundancy* that can be used to discard certain TGDs/rules produced by Inf.

Definition 5.1. A TGD $\tau_1 = \forall \vec{x}_1 [\beta_1 \rightarrow \exists \vec{y}_1 \eta_1]$ is a syntactic tautology if it is in head-normal form and $\beta_1 \cap \eta_1 \neq \emptyset$. TGD τ_1 subsumes a TGD $\tau_2 = \forall \vec{x}_2 [\beta_2 \rightarrow \exists \vec{y}_2 \eta_2]$ if there exists a substitution μ such that $\text{dom}(\mu) = \vec{x}_1 \cup \vec{y}_1$, $\mu(\vec{x}_1) \subseteq \vec{x}_2$, $\mu(\vec{y}_1) \subseteq \vec{y}_2 \cup \vec{y}'_2$, $\mu(y) \neq \mu(y')$ for distinct y and y' in \vec{y}_1 , $\mu(\beta_1) \subseteq \beta_2$, and $\mu(\eta_1) \supseteq \eta_2$.

A rule $\tau_1 = \forall \vec{x}_1 [\beta_1 \rightarrow H_1]$ is a syntactic tautology if $H_1 \in \beta_1$. Rule τ_1 subsumes a rule $\tau_2 = \forall \vec{x}_2 [\beta_2 \rightarrow H_2]$ if there exists a substitution μ such that $\mu(\beta_1) \subseteq \beta_2$ and $\mu(H_1) = H_2$.

A TGD/rule τ is contained in a set of TGDs/rules S up to redundancy if τ is a syntactic tautology or some $\tau' \in S$ subsumes τ .

The following example illustrates Definition 5.1.

EXAMPLE 5.2. Rule $A(x) \wedge B(x) \rightarrow A(x)$ is a syntactic tautology: applying a chase step with it cannot produce a new fact. A non-full TGD in head-normal form cannot be a syntactic tautology since each head atom of such a TGD contains an existentially quantified variable that does not occur in the TGD body.

Rule $\tau_1 = A(f(x_1), f(x_1)) \wedge B(x_1) \rightarrow B(f(x_1))$ is subsumed by rule $\tau_2 = A(x_2, x_3) \rightarrow B(x_2)$ using substitution μ_1 that maps both x_2 and x_3 to $f(x_1)$. If τ_1 derives $B(f(t))$ in one step from a set of facts I by a substitution σ where $\sigma(x_1) = t$, then τ_2 also derives $B(f(t))$ from I in one step by substitution $\sigma \circ \mu_1$. Thus, rule τ_1 is not needed when rule τ_2 is present, so τ_1 can be discarded.

While syntactic tautologies and rule subsumption are standard in first-order theorem proving [8], subsumption of TGDs is more involved. TGD $\tau_3 = A(x_1, x_1) \wedge B(x_1) \rightarrow \exists y_1 C(x_1, y_1)$ is subsumed by TGD $\tau_4 = A(x_2, x_3) \rightarrow \exists y_2, y_3 C(x_2, y_2) \wedge D(x_3, y_3)$ by substitution μ_2 where $\mu_2(x_2) = \mu_2(x_3) = x_1$, $\mu_2(y_2) = y_1$, and $\mu_2(y_3) = y_3$. The conditions on substitution μ_2 in Definition 5.1 ensure that y_2 and y_3 are not mapped to each other or to x_1 . Thus, as in the previous paragraph, the result of each chase step with τ_3 and substitutions σ and σ' can always be obtained (up to isomorphism) by a chase step with τ_4 and substitutions $\sigma \circ \mu_2$ and $\sigma' \circ \mu_2$. \triangleleft

In Definition 5.3 we formalize the notion of applying Inf exhaustively up to redundancy. The definition, however, does not say how to actually do it: we discuss this and other issues in Section 6.

Definition 5.3. For Inf an inference rule and Σ a finite set of GTGDs, $\text{Inf}(\Sigma)$ is the subset of all Skolem-free Datalog rules of Σ' , where Σ' is the smallest set that contains up to redundancy each TGD/rule obtained by

- transforming Σ into head-normal form if Inf manipulates TGDs or Skolemizing Σ if Inf manipulates rules, and
- selecting an adequate number of premises in Σ' , renaming any variables shared by distinct premises, applying Inf to the renamed premises, and transforming the result into head-normal form.

5.1 The Existential-Based Rewriting

As we discussed in Section 4, each loop T_i, \dots, T_j at vertex v in a one-pass chase sequence can be seen as taking $T_i(v)$ as input and producing one fact included in $T_j(v)$ as output. Let v' be child of v introduced in T_{i+1} . The idea behind the ExbDR algorithm is to derive all GTGDs such that, for each k with $i < k \leq j$, all facts of $T_k(v')$

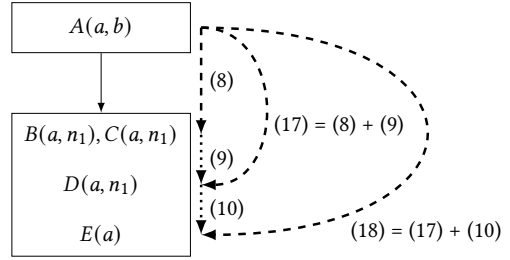


Figure 3: Deriving “shortcuts” for the loop T_0 – T_4 in ExbDR

can be derived from the input $T_i(v)$ in one step. The output of the loop can then also be derived from $T_i(v)$ in one step by full GTGD, so this GTGD provides us with the desired loop “shortcut”. Before formalizing this idea, we slightly adapt the notion of unification.

Definition 5.4. For X a set of variables, an X -unifier and an X -MGU θ of atoms A_1, \dots, A_n and B_1, \dots, B_n are defined as in Section 3, but with the additional requirement that $\theta(x) = x$ for each $x \in X$.

It is straightforward to see that an X -MGU is unique up to the renaming of variables not contained in X , and that it can be computed as usual while treating variables in X as if they were constants. We are now ready to formalize the ExbDR algorithm.

Definition 5.5. The Existential-Based Datalog Rewriting inference rule ExbDR takes two guarded TGDs

$$\begin{aligned} \tau &= \forall \vec{x} [\beta \rightarrow \exists \vec{y} \eta \wedge A_1 \wedge \dots \wedge A_n] \quad \text{with } n \geq 1 \text{ and} \\ \tau' &= \forall \vec{z} [A'_1 \wedge \dots \wedge A'_n \wedge \beta' \rightarrow H'] \end{aligned}$$

and, for θ a \vec{y} -MGU of A_1, \dots, A_n and A'_1, \dots, A'_n , if $\theta(\vec{x}) \cap \vec{y} = \emptyset$ and $\text{vars}(\theta(\beta')) \cap \vec{y} = \emptyset$, it derives

$$\theta(\beta) \wedge \theta(\beta') \rightarrow \exists \vec{y} \theta(\eta) \wedge \theta(A_1) \wedge \dots \wedge \theta(A_n) \wedge \theta(H').$$

EXAMPLE 5.6. Consider again the set Σ from Example 4.3. The idea behind the ExbDR algorithm is illustrated in Figure 3, which summarizes the steps of the loop T_0, T_1, T_2, T_3, T_4 from Figure 1. We denote the vertices by r and v_1 as in Example 4.3.

Fact $A(a, b)$ is the input to the loop, and the first step of the loop derives $B(a, n_1)$ and $C(a, n_1)$ using GTGD (8). Next, GTGD (9) evolves vertex v_1 by deriving $D(a, n_1)$. To capture this, the ExbDR inference rule combines (8), the GTGD that creates v_1 , with (9), the GTGD that evolves v_1 . This produces GTGD (17), which derives all facts of v_1 from the input fact in one step. Vertex v_1 is evolved further using GTGD (10) to derive $E(a)$. To reflect this, the ExbDR inference rule combines (17) and (10) to produce (18), which again derives all facts of v_1 from the loop’s input in one step.

$$A(x_1, x_2) \rightarrow \exists y B(x_1, y) \wedge C(x_1, y) \wedge D(x_1, y) \quad (17)$$

$$A(x_1, x_2) \rightarrow \exists y B(x_1, y) \wedge C(x_1, y) \wedge D(x_1, y) \wedge E(x_1) \quad (18)$$

Fact $E(a)$ does not contain the labeled null n_1 that is introduced when creating v_1 , so it can be propagated to the root vertex r as the output of the loop. This is reflected in (18): atom $E(x_1)$ does not contain any existential variables. Definition 5.3 requires each derived GTGD to be brought into head-normal, so (18) is broken up into (17) and (14). The latter GTGD is full, and it provides us with the desired shortcut for the loop.

Next, (12) and atom $F(x_1, y_1)$ of (11) produce (19), and transformation into head-normal form produces (11) and (15). Moreover, (8) and (13) produce (20), and transformation (20) into head-normal form produces (16) and (21).

$$A(x_1, x_2) \wedge E(x_1) \rightarrow \exists y_1, y_2 F(x_1, y_1) \wedge F(y_1, y_2) \wedge G(x_1) \quad (19)$$

$$A(x_1, x_2) \wedge G(x_1) \rightarrow \exists y B(x_1, y) \wedge C(x_1, y) \wedge H(x_1) \quad (20)$$

$$A(x_1, x_2) \wedge G(x_1) \rightarrow \exists y B(x_1, y) \wedge C(x_1, y) \quad (21)$$

GTGD (21) is subsumed by (8) so it can be dropped. No further inferences are possible after this, so all derived full GTGDs are returned as the rewriting of Σ . \triangleleft

Before proceeding, we present an auxiliary result showing certain key properties of the ExbDR inference rule.

PROPOSITION 5.7. *Each application of the ExbDR inference rule to τ , τ' , and θ as in Definition 5.5 satisfies the following properties.*

1. *Some atom A'_i with $1 \leq i \leq n$ is a guard in τ' .*
2. *For each $1 \leq i \leq n$ such that A'_i is a guard of τ' , and for σ the \bar{y} -MGU of A'_i and the corresponding atom A_i such that $\sigma(\bar{x}) \cap \bar{y} = \emptyset$, it is the case that $\text{vars}(\sigma(A'_i)) \cap \bar{y} \neq \emptyset$ for each $1 \leq j \leq n$.*
3. *The result is a GTGD whose body and head width are at most $\text{bwidth}(\Sigma)$ and $\text{hwidth}(\Sigma)$, respectively.*

In the second claim of Proposition 5.7, σ unifies only A_i and A'_i , whereas θ unifies all A_1, \dots, A_n and A'_1, \dots, A'_n ; thus, σ and θ are not necessarily the same. The third claim is needed to prove termination of ExbDR.

Proposition 5.7 can be used to guide the application of the ExbDR inference rule. Consider an attempt to apply the ExbDR inference rule to two candidate GTGDs $\tau = \beta \rightarrow \exists \bar{y} \eta$ and $\tau' = \beta' \rightarrow H'$. The first claim of Proposition 5.7 tells us that a guard of τ' will definitely participate in the inference. Thus, we can choose one such guard $G' \in \beta'$ of τ' and try to find a \bar{y} -MGU σ of G' and a counterpart atom $G \in \eta$ from the head of τ . Next, we need to check whether $\sigma(\bar{x}) \cap \bar{y} = \emptyset$; if not, there is no way for $\theta(\bar{x}) \cap \bar{y} = \emptyset$ to hold so the inference is not possible. By the second claim of Proposition 5.7, all candidates for the atoms participating in the inference will contain a variable that is mapped by σ to a member of \bar{y} ; thus, $S' = \{\sigma(A') \mid A' \in \beta' \wedge \text{vars}(\sigma(A')) \cap \bar{y} \neq \emptyset\}$ is the set of all relevant side atoms. Note that we apply σ to the atoms in S' to simplify further matching. The next step is to identify the corresponding head atoms of τ . To achieve this, for each atom $A' \in S'$ of the form $R(t_1, \dots, t_n)$, we identify the set $C[A']$ of candidate counterpart atoms as the set of atoms of the form $R(s_1, \dots, s_n) \in \sigma(\eta)$ such that, for each argument position i with $1 \leq i \leq n$, if either $t_i \in \bar{y}$ or $s_i \in \bar{y}$, then $t_i = s_i$. Finally, we consider each possible combination S of such candidates, and we try to find an MGU θ of sets S and S' . If unification succeeds, we derive the corresponding GTGD.

THEOREM 5.8. *Program ExbDR(Σ) is a Datalog rewriting of a finite set of GTGDs Σ . Moreover, the rewriting can be computed in time $O(b^{r^d \cdot (w_b + c)^{da} \cdot r^d \cdot (w_h + c)^{da}})$ for r the number of relations in Σ , a the maximum relation arity in Σ , $w_b = \text{bwidth}(\Sigma)$, $w_h = \text{hwidth}(\Sigma)$, $c = |\text{consts}(\Sigma)|$, and some b and d .*

Program ExbDR(Σ) can thus be large in the worst case. In Section 7 we show empirically that rewritings are suitable for practical

use. From a theoretical point of view, checking fact entailment via ExbDR(Σ) is worst-case optimal. To see why, let r , a , and c be as in Theorem 5.8, and consider a base instance I with c' constants. The fixpoint of ExbDR(Σ) on I contains at most $r(c + c')^a$ facts, and it can be computed in time $O(r(c + c')^a \cdot |\text{ExbDR}(\Sigma)|)$: each rule $\tau \in \text{ExbDR}(\Sigma)$ is guarded so we can apply a chase step with τ by matching a guard and then checking the remaining body atoms. Hence, we can compute ExbDR(Σ) and find its fixpoint in 2ExpTIME, in ExpTIME if the relation arity is fixed, and in PTIME if Σ is fixed (i.e., if we consider *data complexity*). These results match the lower bounds for checking fact entailment for GTGDs [33].

5.2 Using Skolemization

The ExbDR algorithm exhibits two drawbacks. First, each application of the ExbDR inference rule potentially introduces a head atom, so the rule heads can get very long. Second, each inference requires matching a subset of body atoms of τ' to a subset of the head atoms of τ ; despite the optimizations outlined after Proposition 5.7, this can be costly, particularly when rule heads are long.

We would ideally derive GTGDs with a single head atom and unify just one body atom of τ' with the head atom of τ , but this does not seem possible if we stick to manipulating GTGDs. For example, atoms $C(y)$ and $D(y)$ of GTGD (17) refer to the same labeled null (represented by variable y), and this information would be lost if we split (17) into two GTGDs. We thus need a way to refer to the same existentially quantified object in different logical formulas. This can be achieved by replacing existentially quantified variables by Skolem terms, which in turns gives rise to the SkDR algorithm from Definition 5.10. Before presenting the algorithm, in Definition 5.9 we generalize the notion of guardedness to rules.

Definition 5.9. *Rule $\forall \bar{x}[\beta \rightarrow H]$ is guarded if each function symbol in the rule is a Skolem symbol, the body β contains a Skolem-free atom $A \in \beta$ such that $\text{vars}(A) = \bar{x}$, and each Skolem term in the rule is of the form $f(\bar{t})$ where $\text{vars}(f(\bar{t})) = \bar{x}$ and \bar{t} is function-free.*

Definition 5.10. *The Skolem Datalog Rewriting inference rule SkDR takes two guarded rules*

$$\tau = \beta \rightarrow H \quad \text{and} \quad \tau' = A' \wedge \beta' \rightarrow H'$$

such that

- β is Skolem-free and H contains a Skolem symbol, and
- A' contains a Skolem symbol, or τ' is Skolem-free and A' contains all variables of τ' ,

and, for θ an MGU of H and A' , it derives

$$\theta(\beta) \wedge \theta(\beta') \rightarrow \theta(H').$$

EXAMPLE 5.11. Skolemizing GTGDs (8) and (11) produces rules (22)–(23), and (24)–(25), respectively.

$$A(x_1, x_2) \rightarrow B(x_1, f(x_1, x_2)) \quad (22)$$

$$A(x_1, x_2) \rightarrow C(x_1, f(x_1, x_2)) \quad (23)$$

$$A(x_1, x_2) \wedge E(x_1) \rightarrow F(x_1, g(x_1, x_2)) \quad (24)$$

$$A(x_1, x_2) \wedge E(x_1) \rightarrow F(g(x_1, x_2), h(x_1, x_2)) \quad (25)$$

Intuitively, rules (22) and (23) jointly represent the facts introduced by the non-full GTGD (8): functional term $f(x_1, x_2)$ allows both

rules to “talk” about the same labeled nulls. This allows the SkDR inference rule to simulate the ExbDR inference rule while unifying just pairs of atoms. In particular, SkDR combines (22) and (10) to obtain (26); it combines (23) and (9) to obtain (27); and it combines (26) and (27) to obtain the “shortcut” rule (14).

$$A(x_1, x_2) \wedge D(x_1, f(x_1, x_2)) \rightarrow E(x_1) \quad (26)$$

$$A(x_1, x_2) \rightarrow D(x_1, f(x_1, x_2)) \quad (27)$$

The rules with Skolem-free bodies derived in this way allow us to reconstruct derivations in one step analogously to Example 5.6, and the rules with Skolem symbols in body atoms capture the intermediate derivation steps. For example, rules (26) and (28) capture the result of matching the first and the second body atom, respectively, of rule (10) to facts produced by rules (22) and (27), respectively. To complete the rewriting, SkDR combines (24) with (12) to obtain (15), and it combines (22) with (13) to derive (16).

However, SkDR also combines (10) and (27) into (28), which with (22) derives (14) the second time. These inferences are superfluous: they just process the two body atoms of (10) in a different order. Also, SkDR combines (12) and (25) into rule (29), which is a “dead-end” in that it does not further contribute to a Datalog rule.

$$A(x_1, x_2) \wedge B(x_1, f(x_1, x_2)) \rightarrow E(x_1) \quad (28)$$

$$A(x_1, x_2) \wedge E(x_1) \wedge E(g(x_1, x_2)) \rightarrow G(g(x_1, x_2)) \quad (29)$$

Our HypDR algorithm in Subsection 5.3 can avoid these overheads, but at the expense of using more than two rules at a time. ◀

Proposition 5.12 and Theorem 5.13 capture the relevant properties of the SkDR algorithm.

PROPOSITION 5.12. *Each application of the SkDR inference rule to rules τ and τ' as in Definition 5.10 produces a guarded rule.*

THEOREM 5.13. *Program SkDR(Σ) is a Datalog rewriting of a finite set of GTGDs Σ . Moreover, the rewriting can be computed in time $O(b^{r^d \cdot (e + w_b + c)^{da}})$ for r the number of relations in Σ , a the maximum relation arity in Σ , e the number of existential quantifiers in Σ , w_b = bwidth(Σ), c = |consts(Σ)|, and some b and d .*

It is natural to wonder whether SkDR is guaranteed to be more efficient than ExbDR. We next show that neither algorithm is generally better: there exist families of inputs on which SkDR performs exponentially more inferences than ExbDR, and vice versa.

PROPOSITION 5.14. *There exists a family $\{\Sigma_n\}_{n \in \mathbb{N}}$ of finite sets of GTGDs such that the number of GTGDs derived by ExbDR is $O(2^n)$ times larger than the number of rules derived by SkDR on each Σ_n .*

PROOF. For each $n \in \mathbb{N}$, let Σ_n contain the following GTGDs.

$$A(x) \rightarrow \exists \vec{y} B_1(x, y_1) \wedge \cdots \wedge B_n(x, y_n) \quad (30)$$

$$B_i(x_1, x_2) \wedge C_i(x_1) \rightarrow D_i(x_1, x_2) \text{ for } 1 \leq i \leq n \quad (31)$$

On such Σ_n , ExbDR derives a GTGD of the form (32) for each subset $\{k_1, \dots, k_m\} \subseteq \{1, \dots, n\}$, and there are 2^n such GTGDs. In contrast, the Skolemization of (30) consists of n rules shown in equation (33), so SkDR derives just n rules shown in equation (34).

$$A(x) \wedge \bigwedge_{i=1}^m C_{k_i}(x) \rightarrow \exists \vec{y} \bigwedge_{i=1}^n B_i(x, y_i) \wedge \bigwedge_{i=1}^m D_{k_i}(x, y_{k_i}) \quad (32)$$

$$A(x) \rightarrow B_i(x, f_i(x)) \text{ for } 1 \leq i \leq n \quad (33)$$

$$A(x) \wedge C_i(x) \rightarrow D_i(x, f_i(x)) \text{ for } 1 \leq i \leq n \quad \square(34)$$

PROPOSITION 5.15. *There exists a family $\{\Sigma_n\}_{n \in \mathbb{N}}$ of finite sets of GTGDs such that the number of rules derived by SkDR is $O(2^n)$ times larger than the number of GTGDs derived by ExbDR on each Σ_n .*

PROOF. For each $n \in \mathbb{N}$, let Σ_n contain the following GTGDs.

$$A(x) \rightarrow \exists y B_1(x, y) \wedge \cdots \wedge B_n(x, y) \quad (35)$$

$$B_1(x_1, x_2) \wedge \cdots \wedge B_n(x_1, x_2) \rightarrow C(x_1) \quad (36)$$

On such Σ_n , ExbDR derives just GTGD (37) in one step. In contrast, the Skolemization of (35) consists of n rules of the form (38) for each $1 \leq i \leq n$. Thus, SkDR combines these with (36) to derive $2^n - 1$ rules of the form (39), one for each subset $\{k_1, \dots, k_m\} \subsetneq \{1, \dots, n\}$.

$$A(x) \rightarrow C(x) \quad (37)$$

$$A(x) \rightarrow B_i(x, f(x)) \quad (38)$$

$$A(x) \wedge B_{k_1}(x, f(x)) \wedge \cdots \wedge B_{k_m}(x, f(x)) \rightarrow C(x) \quad \square \quad (39)$$

5.3 Combining Several SkDR Steps into One

The SkDR algorithm can produce many rules with Skolem symbols in the body, which is the main reason for Proposition 5.15. We next present the HypDR algorithm, which uses the *hyperresolution* inference rule as a kind of “macro” to combine several SkDR steps into one. We show that this can be beneficial for several reasons.

Definition 5.16. *The Hyperresolution Rewriting inference rule HypDR takes guarded rules*

$$\tau_1 = \beta_1 \rightarrow H_1 \quad \dots \quad \tau_n = \beta_n \rightarrow H_n \text{ and}$$

$$\tau' = A'_1 \wedge \cdots \wedge A'_n \wedge \beta' \rightarrow H'$$

such that

- for each i with $1 \leq i \leq n$, conjunction β_i is Skolem-free and atom H_i contains a Skolem symbol, and

- rule τ' is Skolem-free,

and, for θ an MGU of H_1, \dots, H_n and A'_1, \dots, A'_n , if conjunction $\theta(\beta')$ is Skolem-free, it derives

$$\theta(\beta_1) \wedge \cdots \wedge \theta(\beta_n) \wedge \theta(\beta') \rightarrow \theta(H').$$

EXAMPLE 5.17. The HypDR inference rule simulates chase steps in the child vertex of a loop analogously to ExbDR: all body atoms matching a fact introduced in the child vertex are resolved in one step. We can see two benefits of this on our running example.

First, HypDR derives (27) from (23) and (9), and it derives (14) from (10), (22), and (27). Rule (14) is derived just once, and without intermediate rules (26) and (28). In other words, the HypDR inference rule does not resolve the body atoms of a rule in every possible order. As Proposition 5.20 below shows, this can reduce the number of derived rules by an exponential factor.

Second, HypDR derives only rules with Skolem-free bodies, and thus does not derive the “dead-end” rule (29). In other words, all consequences of HypDR derive in one step one fact in the child vertex of a loop from the loop’s input $T_i(v)$.

The downside of HypDR is that more than two rules can participate in an inference. This requires more complex unification and selection of candidates that can participate in an inference. ◀

Proposition 5.18 and Theorem 5.19 capture the properties of HypDR, and Proposition 5.20 compares it to SkDR.

PROPOSITION 5.18. *Each application of the HypDR inference rule to rules τ_1, \dots, τ_n and τ' as in Definition 5.16 produces a guarded rule.*

THEOREM 5.19. *Program HypDR(Σ) is a Datalog rewriting of a finite set of GTGDs Σ . Moreover, the rewriting can be computed in time $O(b^{r^d \cdot (e + w_b + c)^{da}})$ for r the number of relations in Σ , a the maximum relation arity in Σ , e the number of existential quantifiers in Σ , $w_b = \text{bwidth}(\Sigma)$, $c = |\text{consts}(\Sigma)|$, and some b and d .*

PROPOSITION 5.20. *There exists a family $\{\Sigma_n\}_{n \in \mathbb{N}}$ of finite sets of GTGDs such that SkDR derives $O(2^n)$ more rules than HypDR on each Σ_n .*

PROOF. For each $n \in \mathbb{N}$, let Σ_n contain the following GTGDs.

$$A(x) \rightarrow \exists y B(x, y) \quad (40)$$

$$B(x_1, x_2) \wedge C_i(x_1) \rightarrow D_i(x_1, x_2) \text{ for } 1 \leq i \leq n \quad (41)$$

$$D_1(x_1, x_2) \wedge \dots \wedge D_n(x_1, x_2) \rightarrow E(x_1) \quad (42)$$

Skolemizing (40) produces (43). Thus, SkDR combines (43) with each (41) to derive each (44), and it uses (44) and (42) to derive $2^n - 1$ rules of the form (45) for each set of indexes I satisfying $\emptyset \subsetneq I \subseteq \{1, \dots, n\}$; note that none of these rules are redundant.

$$A(x) \rightarrow B(x, f(x)) \quad (43)$$

$$A(x) \wedge C_i(x) \rightarrow D_i(x, f(x)) \text{ for } 1 \leq i \leq n \quad (44)$$

$$A(x) \wedge \bigwedge_{i \in I} C_i(x) \wedge \bigwedge_{j \in \{1, \dots, n\} \setminus I} D_j(x, f(x)) \rightarrow E(x) \quad (45)$$

In contrast, HypDR derives each (44) just like SkDR, and it combines in one step (42) and all (44) to derive (45) for $I = \{1, \dots, n\}$. \square

6 IMPLEMENTATION AND OPTIMIZATIONS

In this section, we discuss numerous issues that have to be addressed to make the computation of a rewriting practical.

Computing $\text{Inf}(\Sigma)$ in Practice. Definition 5.3 does not specify how to compute the set Σ' , and redundancy elimination makes this question nontrivial. When Inf derives a TGD/rule τ , we can apply subsumption in two ways. First, we can discard τ if τ is subsumed by a previously derived TGD/rule; this is known as *forward subsumption*. Second, if τ is not discarded, we can discard each previously derived TGD/rule that is subsumed by τ ; this is known as *backward subsumption*. The set of derived TGD/rules can thus grow and shrink, so the application of Inf has to be carefully structured to ensure that all inferences are performed eventually.

We address this problem by a variant of the *Otter loop* [35] used in first-order theorem provers. The pseudo-code is shown in Algorithm 1. The algorithm maintains two sets of TGDs/rules: the *worked-off* set \mathcal{W} contain TGDs/rules that have been processed by Inf , and the *unprocessed* set \mathcal{U} contains TGDs/rules that are still to be processed. Set \mathcal{W} is initially empty (line 1), and set \mathcal{U} is initialized to the head-normal form of Σ if Inf manipulates TGDs, or to the Skolemization of Σ if Inf manipulates rules. The algorithm then processes each $\tau \in \mathcal{U}$ until \mathcal{U} becomes empty (lines 3–10). It is generally beneficial to process shorter TGDs/rules first as that improves chances of redundancy elimination. After moving τ to \mathcal{W} (line 5),

Algorithm 1 Computing $\text{Inf}(\Sigma)$ for Σ a finite set of GTGDs

```

1:  $\mathcal{W} = \emptyset$ 
2:  $\mathcal{U}$  = the head-normal form or the Skolemization of  $\Sigma$ 
3: while  $\mathcal{U} \neq \emptyset$  do
4:   Choose some  $\tau \in \mathcal{U}$  and remove it from  $\mathcal{U}$ 
5:    $\mathcal{W} = \mathcal{W} \cup \{\tau\}$ 
6:   Let  $\mathcal{E}$  be the result of applying  $\text{Inf}$  to  $\tau$  and a subset of  $\mathcal{W}$ 
     and transforming the result into head-normal form
7:   for each  $\tau' \in \mathcal{E}$  do
8:     if  $\tau'$  is not contained in  $\mathcal{W} \cup \mathcal{U}$  up to redundancy then
9:       Remove from  $\mathcal{W}$  and  $\mathcal{U}$  each  $\tau''$  subsumed by  $\tau'$ 
10:       $\mathcal{U} = \mathcal{U} \cup \{\tau'\}$ 
11: return  $\{\tau \in \mathcal{W} \mid \tau \text{ is a Skolem-free Datalog rule}\}$ 

```

the algorithm applies Inf to τ and \mathcal{W} and transforms the results into head-normal form (line 6). The algorithm discards each resulting $\tau' \in \mathcal{E}$ that is a syntactic tautology or is forward-subsumed by an element of $\mathcal{W} \cup \mathcal{U}$ (line 8). If τ' is not discarded, the algorithm applies backward subsumption to τ' , \mathcal{W} , and \mathcal{U} (line 9) and adds τ' to \mathcal{U} (line 10). When all TGDs/rules are processed, the algorithm returns all Skolem-free Datalog rules from \mathcal{W} (line 11). The result of applying Inf to TGDs/rules in \mathcal{W} is thus contained in $\mathcal{W} \cup \mathcal{U}$ up to redundancy at all times so, upon algorithm's termination, set \mathcal{W} satisfies the condition on Σ' from Definition 5.3.

Checking Subsumption. Checking whether TGD/rule τ_1 subsumes τ_2 is NP-complete [31], and the main difficulty is in matching the variables of τ_1 to the variables of τ_2 . Thus, we use an approximate check in our implementation. First, we normalize each TGD to use fixed variables x_1, x_2, \dots and y_1, y_2, \dots : we sort the body and head atoms by their relations using an arbitrary, but fixed ordering and breaking ties arbitrarily, and then we rename all variables so that the i^{th} distinct occurrence of a universally (respectively existentially) quantified variable from left to right is x_i (respectively y_i). To see whether $\tau_1 = \beta_1 \rightarrow \exists \vec{y} \eta_1$ subsumes $\tau_2 = \beta_2 \rightarrow \exists \vec{y} \eta_2$, we determine whether $\beta_1 \subseteq \beta_2$ and $\eta_1 \supseteq \eta_2$ holds, which requires only polynomial time. We use a similar approximation for rules. Variable normalization ensures termination, and using a modified subsumption check does not affect the correctness of the rewriting: set \mathcal{W} may contain more TGDs/rules than strictly necessary, but these are all logical consequences of (the Skolemization of) Σ .

Subsumption Indexing. Sets \mathcal{W} and \mathcal{U} can be large, so we use a variant of *feature vector indexing* [43] to retrieve subsumption candidates in $\mathcal{W} \cup \mathcal{U}$. For simplicity, we consider only TGDs in the following discussion, but rules can be handled analogously. Note that a TGD τ_1 can subsume TGD τ_2 only if the set of relations occurring in the body of τ_1 (respectively the head of τ_2) is a subset of the set of relations occurring in the body of τ_2 (respectively the head of τ_1). Thus, we can reduce the problem of retrieving subsumption candidates to the problem of, given a domain set D , a set N of subsets of D , a subset $S \subseteq D$, and $\bowtie \in \{\subseteq, \supseteq\}$, retrieving each $S' \in N$ satisfying $S' \bowtie S$. The set-trie data structure [42] can address this problem. The idea is to order D in an arbitrary, yet fixed way, so that we can treat each subset of N as a word over D . We then index N by constructing a trie over the words representing the elements

of N . Finally, we retrieve all $S' \in N$ satisfying $S' \bowtie S$ by traversing the trie, where the ordering on D allows us to considerably reduce the number of vertices we visit during the traversal.

A minor issue is that retrieving TGDs that subsume a given TGD requires both subset and superset testing for body and head relations, respectively, and vice versa for retrieval of subsumed TGDs. To address this, we introduce a distinct symbol R^b and R^h for each relation R occurring in Σ , and we represent each TGD τ as a *feature vector* F_τ of these symbols corresponding to the body and head of τ . Moreover, we combine in the obvious way the subset and superset retrieval algorithms. For example, when searching for a TGD $\tau' \in \mathcal{W} \cup \mathcal{U}$ that subsumes a given TGD τ , we use the subset retrieval for the symbols R^b and the superset retrieval for symbols R^h . Finally, we order these symbols by the decreasing frequency of the order of the symbols' occurrence in the set Σ of input TGDs, and moreover we order each R^b before all R^h .

Relation Clustering. We observed that the subsumption indexes can easily get very large, so index traversal can become a considerable source of overhead. To reduce the index size, we group the symbols R^b and R^h into clusters C^b and C^h , respectively. Then, the feature vector F_τ associated with each TGD τ consists of all clusters C^b and C^h that contain a relation occurring in the body and head, respectively, of τ . We adapt the trie traversal algorithms in the obvious way to take into account this change. The number of clusters is computed using the average numbers of symbols and atoms in the input TGDs, and clusters are computed with the aim of balancing the number of TGDs stored in each leaf vertex.

Unification Indexing. We construct indexes over \mathcal{W} that allow us to quickly identify TGDs/rules that can participate in an inference with some τ . For TGDs, we maintain a hash table that maps each relation R to a set of TGDs containing R in the body, and another hash table that does the same but for TGD heads. To index rules, we use a variant of a *path indexing* [44]: each atom in a rule is represented as a sequence of relations and function symbols occurring in the atom, and such sequences are entered into two tries (one for body and one for head atoms). Then, given rule τ , we consider each body and head atom A of τ , we convert A into the corresponding sequence, and we use the sequence to query the relevant trie for all candidates participating in an inference with τ on A .

Cheap Lookahead Optimization. Consider an application of the ExbDR inference rule to GTGDs τ and τ' as in Definition 5.5, producing a GTGD τ'' where $\text{vars}(\theta(H')) \cap \vec{y} \neq \emptyset$ and the relation of H' does not occur in the body of a GTGD in Σ . In each one-pass chase sequence for some base instance and Σ , no GTGD of Σ can be applied to a fact obtained by instantiating $\theta(H')$, so deriving this fact is redundant. Consequently, we can drop such τ'' as soon as we derive it in line 6. Analogously, when the SkDR inference rule is applied to rules τ and τ' as in Definition 5.10, we can drop the resulting rule if $\theta(H')$ is not full and it contains a relation not occurring in the body of a GTGD in Σ .

7 EXPERIMENTAL EVALUATION

We implemented a system that can produce a Datalog rewriting of a set of GTGDs using our algorithms, and we conducted an empirical evaluation using a comprehensive collection of 428 synthetic and

Table 1: Input GTGDs at a Glance

Inputs	# Full TGDs				# Non-Full TGDs			
	Min	Max	Avg	Med	Min	Max	Avg	Med
428	1	171,905	11,030	789	2	156,743	5,255	283

realistic inputs. Our objectives were to show that our algorithms can indeed rewrite complex GTGDs, and that the rewriting can be successfully processed by modern Datalog systems. In Subsection 7.1 we describe the test setting. Then, in Subsection 7.2 we discuss the rewriting experiments with GTGDs obtained from ontologies, and in Subsection 7.3 we validate the usefulness of the rewriting approach end-to-end. Finally, in Subsection 7.4 we discuss rewriting GTGDs of higher arity. Due to the very large number of inputs, we can only summarize our results in this paper; however, our complete evaluation results are available online [13].

7.1 Input GTGDs, Competitors, & Test Setting

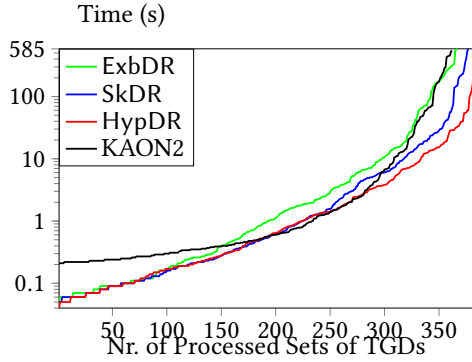
Before discussing our results, we next describe our test setting.

Input GTGDs. We are unaware of any publicly available sets of GTGDs that we could readily use in our evaluation, so we derived the input GTGDs for our evaluation from the ontologies in the Oxford Ontology Library [39]. At the time of writing, this library contained 787 ontologies, each assigned a unique five-digit identifier. After removing closely-related ontology variants, we were left with 428 core ontologies. We loaded each ontology using the parser from the Graal system [9], discarded axioms that cannot be translated into GTGDs, and converted the remaining axioms into GTGDs. We used the standard translation of description logics into first-order logic [6], where each class corresponds to a unary relation, and each property corresponds to a binary relation. We thus obtained 428 sets of input GTGDs with properties shown in Table 1.

To evaluate our algorithms on TGDs containing relations of arity higher than two, we devised a way to “blow up” relation arity. Given a set of GTGDs and a blowup factor b , our method proceeds as follows. First, in each atom of each GTGD, it replaces each variable argument with b fresh variables uniquely associated with the variable; for example, for $b = 2$, atom $A(x, y)$ is transformed into atom $A(x_1, x_2, y_1, y_2)$. Next, the method randomly introduces fresh head and body atoms over the newly introduced variables; in doing so, it ensures that the new atoms do not introduce patterns that would prevent application of the ExbDR inference rule.

Competitors. We compared the ExbDR, SkDR, and HypDR algorithms, implemented as described in Section 6. As noted in Section 2, no existing system we are aware of implements a Datalog rewriting algorithm for GTGDs. However, the KAON2 system [28, 37, 38] can rewrite GTGDs obtained from OWL ontologies, so we used KAON2 as a baseline in our experiments with OWL-based GTGDs. We made sure that all inputs to KAON2 and our algorithms include only GTGDs that all methods can process.

Test Setting. We conducted all experiments on a laptop with an Intel Core i5-6500 CPU @ 3.20 GHz and 16 GB of RAM, running Ubuntu 20.04.4 LTS and Java 11.0.15. In each test run, we loaded a set of TGDs, measured the wall-clock time required to compute



		$time(Y)/time(X) \geq 10$				X and Y both fail			
X \ Y		ExbDR	SkDR	HypDR	KAON2	ExbDR	SkDR	HypDR	KAON2
ExbDR			19	0	19	61			
SkDR		37		0	26	33	51		
HypDR		37	12		31	35	43	46	
KAON2		35	15	0		37	47	46	66

Figure 4: Results for GTGDs Derived from Ontologies

the rewriting of a set of GTGDs, and saved the produced Datalog rewriting. We used a timeout of ten minutes for each test run.

7.2 Experiments with GTGDs from Ontologies

We computed the Datalog rewriting of GTGDs obtained from OWL ontologies using our three algorithms and KAON2. Figure 4 shows the number of inputs that each algorithm processed in a given time, provides information about the inputs and outputs of each system, and compares the performance among systems. The input size for ExbDR is the number of GTGDs after transforming the input into head-normal form, and for SkDR and HypDR it is the number of rules after Skolemization. Input size is not available for KAON2 since this system reads an OWL ontology and transforms it into GTGDs internally. The output size is the number of Datalog rules in the rewriting. Finally, the blowup is the ratio of the output and the input sizes. Each input GTGD contained at most seven body atoms. Out of 428 inputs, 349 were processed within the ten minute limit by our three systems, and 334 inputs were processed by all four systems. Moreover, 32 inputs, each containing between 20,270 and 221,648 GTGD, were not processed by any system.

Discussion. As one can see in Figure 4, all algorithms were able to compute the rewriting of large inputs containing 100k+ GTGDs. Moreover, for the vast majority of inputs that were successfully

processed, the size of the rewriting and the number of body atoms in the rewriting are typically of the same order of magnitude as the input. Hence, the worst-case exponential blowup from Theorems 5.8, 5.13, and 5.19 does not appear in practice: the size of the rewriting seems to be determined primarily by the input size.

Relative Performance. No system can be identified as the best in general, but HypDR seems to offer the best performance on average. The algorithm was able to process most inputs; it was at least 35% faster than the other systems on the slowest input; it was never slower by an order of magnitude; there were only 14 inputs that could be processed by some other algorithm but not HypDR; and the output of HypDR does not differ significantly from the output of SkDR. This is in line with our motivation for HypDR outlined in Example 5.17. Specifically, HypDR derives rules with just one head atom, but it does not derive intermediate rules with functional body atoms. The main source of overhead in HypDR seems to be more complex selection of rules participating in an inference.

Impact of Subsumption. All algorithms spend a considerable portion of their running time checking TGD/rule subsumption, so it is natural to wonder whether this overhead is justified. To answer this question, we ran our three approaches using a modification of Algorithm 1: we replaced the check for containment up to redundancy in line 8 with just checking $\tau' \notin \mathcal{W} \cup \mathcal{U}$, and we removed line 9. Note that our normalization of variables described in Section 6 still guarantees termination. This change significantly increased the number of derivations: the numbers of derived GTGDs/rules increased on average by a factor of 104, 185, and 103 on ExbDR, SkDR, and HypDR, respectively. Interestingly, this increase did not affect the performance uniformly. While SkDR was able to process 12 inputs an order of magnitude faster, ExbDR and HypDR timed out on 72 and 17 additional inputs, respectively. This, we believe, is due to how different inference rules select inference candidates. The SkDR rule is applied to just pairs of rules, and candidate pairs can be efficiently retrieved using unification indexes. In contrast, ExbDR requires matching several head atoms with as many body atoms, which makes developing a precise index for candidate pair retrieval difficult; thus, as the number of derived GTGDs increases, the number of false candidates retrieved from the index increases as well. Finally, HypDR can be applied to an arbitrary number of rules, so selecting inference candidates clearly becomes more difficult as the number of derived rules increases.

Impact of Structural Transformation. KAON2 uses *structural transformation* [7] to simplify ontology axioms before translating them into GTGDs. For example, axiom $A \sqsubseteq \exists B.\exists C.D$ is transformed into $A \sqsubseteq \exists B.X$ and $X \sqsubseteq \exists C.D$ for X a fresh class. The resulting axioms have simpler structure, which is often beneficial to performance. To see how this transformation affects our algorithms, we reran our experiments while transforming the input axioms in the same way as in KAON2. This indeed improved the performance of SkDR by one order of magnitude on 22 ontologies, and it did not hurt the performance of HypDR. The main challenge is to generalize this transformation to arbitrary GTGDs: whereas description logic axioms exhibit syntactic nesting that lends itself naturally to this transformation, it is less clear how to systematically apply this

Table 2: Computing the Fixpoint of the Rewriting

Ont. ID	# Rules	# Input Facts	# Output Facts	Time (s)
00387	63,422	4,403,105	51,439,424	53
00448	67,986	5,510,444	107,235,697	110
00470	75,146	10,532,943	141,396,446	242
00471	78,977	11,077,423	128,954,126	253
00472	75,146	10,533,008	141,396,576	279
00473	78,977	11,077,459	128,954,198	291
00573	113,959	9,197,254	155,118,592	206
00682	68,461	5,183,460	105,431,952	101
00684	81,553	6,057,017	66,981,628	109
00686	124,846	10,402,324	166,366,039	238

transformation to TGDs, where heads and bodies consist of “flat” conjunctions. We leave this question for future work.

7.3 End-to-End Experiments

To validate our approach end-to-end, we selected ten inputs where ExbDR produced the largest rewritings. For each of these, we generated a large base instance using WatDiv [2], and we computed the fixpoint of the rewriting and the instance using the RDFS [45] Datalog system v5.4. Table 2 summarizes our results.

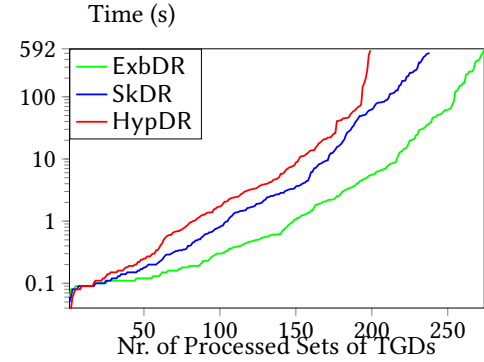
All programs used in this experiment are at least several orders of magnitude larger than what is usually encountered in practical applications of Datalog, but RDFS nevertheless computed the fixpoint of all rewritings in a few minutes. Moreover, although the fixpoints seem to be an order of magnitude larger than the base instance, this is not a problem for highly optimized systems such as RDFS. Hence, checking fact entailment via rewritings produced by our algorithms is feasible in practice.

7.4 GTGDs With Relations of Higher Arity

Finally, we computed the rewriting of GTGDs obtained by blowing up relation arity as described in Subsection 7.1 using a blowup factor of five. We did not use KAON2 since this system supports relations of arity at most two. Figure 5 summarizes our results. Out of 428 inputs, 187 were processed within the ten minute limit by our three systems, and 128 inputs were not processed by any system.

While HypDR performed best on GTGDs derived from ontologies, Figure 5 shows it to be worst-performing on higher-arity GTGDs: it successfully processed only 199 inputs within the ten minute timeout, whereas SkDR and ExbDR processed 238 and 274 inputs, respectively. This is mainly due to additional body atoms introduced by our “blowup” method: these increase the number of rules participating in an application of the HypDR inference rule, which makes selecting the participating rules harder.

This experiment proved to be more challenging, as most problems discussed in Section 6 became harder. For example, in ExbDR, higher arity of atoms increases the likelihood that an atom retrieved through a unification index does not unify with a given atom, and that the atoms of the selected GTGDs cannot be successfully matched. Subsumption indexing is also more difficult for similar reasons. However, the inputs used in this experiment consist of a large numbers of GTGDs with relations of arity ten, so they can be seen as a kind of a “stress test”. Our algorithms were able to



		ExbDR	SkDR	HypDR
# of Processed Inputs		274	238	199
Max. Processed Input Size		69,046	182,569	38,362
Max. Output Size		58,749	171,832	38,335
Max. Size Blowup		9.00	5.84	5.84
# Blowup ≥ 1.5		26	5	3
Time (s)	Min.	0.06	0.05	0.04
	Max.	591.82	504.49	557.75
	Avg.	26.70	38.39	17.05
	Med.	0.61	1.65	1.72

		$time(Y)/time(X) \geq 10$			X and Y both fail		
X	Y	ExbDR	SkDR	HypDR	ExbDR	SkDR	HypDR
ExbDR			61	87	154		
SkDR		11		21	128	190	
HypDR		6	4		148	184	229

Figure 5: Results for GTGDs with Higher-Arity Relations

process more than half of such inputs, which leads us to believe that they can also handle more well-behaved GTGDs used in practice.

8 CONCLUSION

We presented several algorithms for rewriting a finite set of guarded TGDs into a Datalog program that entails the same base facts on each base instance. Our algorithms are based on a new framework that establishes a close connection between Datalog rewritings and a particular style of the chase. In future, we plan to generalize our framework to wider classes of TGDs, such as frontier-guarded TGDs, as well as provide rewritings for conjunctive queries under certain answer semantics. Moreover, we shall investigate whether the extension of our framework to disjunctive guarded TGDs [30] can be used to obtain practical algorithms for rewriting disjunctive guarded TGDs into disjunctive Datalog programs.

ACKNOWLEDGMENTS

This work was funded by the EPSRC grants OASIS (EP/S032347/1), QUINTON (EP/T022124/1), UK FIRES (EP/S019111/1), AnaLOG (EP/P025943/1), and Concur (EP/V050869/1). For the purpose of Open Access, the author has applied a CC BY public copyright licence to any Author Accepted Manuscript (AAM) version arising from this submission.

REFERENCES

- [1] Shqiponja Ahmetaj, Magdalena Ortiz, and Mantas Simkus. 2018. Rewriting Guarded Existential Rules into Small Datalog Programs. In *ICDT*. Schloss Dagstuhl–Leibniz-Zentrum für Informatik, 4:1–4:24.
- [2] Günes Aluç, Olaf Hartig, M. Tamer Özsu, and Khuzaima Daudjee. 2014. Diversified Stress Testing of RDF Data Management Systems. In *ISWC*. Springer, 197–212.
- [3] Mario Alviano, Nicola Leone, Marco Manna, Giorgio Terracina, and Pierfrancesco Veltri. 2012. Magic-Sets for Datalog with Existential Quantifiers. In *Datalog 2.0*. Springer, 31–43.
- [4] Antoine Amarilli and Michael Benedikt. 2022. When Can We Answer Queries Using Result-Bounded Data Interfaces? *Log. Methods Comput. Sci.* 18, 2 (2022), 14:1–14:81.
- [5] Hajnal Andréka, Johan van Benthem, and István Németi. 1998. Modal Languages and Bounded Fragments of Predicate Logic. *J. Philos. Log.* 27 (1998), 217–274.
- [6] F. Baader, D. Calvanese, D. McGuinness, D. Nardi, and P. F. Patel-Schneider (Eds.). 2007. *The Description Logic Handbook: Theory, Implementation and Applications* (2nd ed.). Cambridge University Press, Cambridge, UK.
- [7] M. Baaz, U. Egly, and A. Leitsch. 2001. Normal Form Transformations. In *Handbook of Automated Reasoning*. MIT Press, Chapter 5, 273–333.
- [8] Leo Bachmair and Harald Ganzinger. 2001. Resolution Theorem Proving. In *Handbook of Automated Reasoning*. MIT Press, Chapter 2, 19–99.
- [9] J.-F. Baget, M. Leclère, M.-L. Mugnier, S. Rocher, and C. Sipieter. 2015. Graal: A Toolkit for Query Answering with Existential Rules. In *RuleML*. Springer, 328–344.
- [10] Jean-François Baget, Marie-Laure Mugnier, Sebastian Rudolph, and Michaël Thomazo. 2011. Walking the Complexity Lines for Generalized Guarded Existential Rules. In *IJCAI*. AAAI Press, 712–717.
- [11] Vince Bárány, Michael Benedikt, and Balder Ten Cate. 2013. Rewriting Guarded Negation Queries. In *MFCS*. Springer, 98–110.
- [12] Luigi Bellomarini, Emanuel Sallinger, and Georg Gottlob. 2018. The Vadalog System: Datalog-based Reasoning for Knowledge Graphs. *Proc. VLDB Endow.* 11, 9 (2018), 975–987.
- [13] Michael Benedikt, Maxime Buron, Stefano Germano, Kevin Kappelmann, and Boris Motik. 2021. *Guarded Saturation*. GitHub. Retrieved July 4, 2022 from <https://krr-oxford.github.io/Guarded-saturation/>
- [14] Michael Benedikt, George Konstantinidis, Giansalvatore Mecca, Boris Motik, Paolo Papotti, Donatello Santoro, and Efthymia Tsamoura. 2017. Benchmarking the Chase. In *PODS*. ACM, 37–52.
- [15] Andrea Cali, Domenico Lembo, and Riccardo Rosati. 2003. Query rewriting and answering under constraints in data integration systems. In *IJCAI*. Morgan Kaufmann, 16–21.
- [16] Diego Calvanese, Benjamin Cogrel, Sarah Komla-Ebri, Roman Kontchakov, Davide Lanti, Martin Rezk, Mariano Rodriguez-Muro, and Guohui Xiao. 2017. Ontop: Answering SPARQL Queries over Relational Databases. *Semantic Web* 8, 3 (2017), 471–487.
- [17] Diego Calvanese, Giuseppe De Giacomo, Domenico Lembo, Maurizio Lenzerini, and Riccardo Rosati. 2007. Tractable Reasoning and Efficient Query Answering in Description Logics: The *DL-Lite* Family. *J. Autom. Reason.* 39, 3 (2007), 385–429.
- [18] Hans de Nivelle. 1998. A Resolution Decision Procedure for the Guarded Fragment. In *CADE*. Springer, 191–204.
- [19] A. Deutsch, L. Popa, and V. Tannen. 2006. Query reformulation with constraints. *SIGMOD Rec.* 35, 1 (2006), 65–73.
- [20] R. Fagin, P. G. Kolaitis, R. J. Miller, and L. Popa. 2005. Data exchange: semantics and query answering. *Theor. Comput. Sci.* 336, 1 (2005), 89–124.
- [21] Mohamed Gaha, Arnaud Zinflou, Christian Langheit, Alexandre Bouffard, Mathieu Viau, and Luc Vouligny. 2013. An Ontology-Based Reasoning Approach for Electric Power Utilities. In *RR*. Springer, 95–108.
- [22] H. Ganzinger and H. de Nivelle. 1999. A Superposition Decision Procedure for the Guarded Fragment with Equality. In *LICS*. IEEE Computer Society, 295–305.
- [23] Georg Gottlob, Sebastian Rudolph, and Mantas Simkus. 2014. Expressiveness of Guarded Existential Rule Languages. In *PODS*. ACM, 27–38.
- [24] Alon Halevy, Anand Rajaraman, and Joann Ordille. 2006. Data Integration: The Teenage Years. In *VLDB*. ACM, 9–16.
- [25] Alon Y. Halevy. 2001. Answering Queries Using Views: A Survey. *VLDB J.* 10, 4 (2001), 270–294.
- [26] Colin Hirsch. 2002. *Guarded Logics: Algorithms and Bisimulation*. Ph.D. Dissertation. RWTH Aachen, Aachen, Germany. Retrieved July 4, 2022 from <http://www.umbrialogic.com/hirsch-thesis.pdf>
- [27] Ullrich Hustadt, Boris Motik, and Ulrike Sattler. 2004. Reducing *SHIQ* Description Logic to Disjunctive Datalog Programs. In *KR*. AAAI Press, 152–162.
- [28] Ullrich Hustadt, Boris Motik, and Ulrike Sattler. 2007. Reasoning in Description Logics by a Reduction to Disjunctive Datalog. *J. Autom. Reason.* 39, 3 (2007), 351–384.
- [29] David S. Johnson and Anthony C. Klug. 1984. Testing Containment of Conjunctive Queries under Functional and Inclusion Dependencies. *J. Comput. Syst. Sci.* 28, 1 (1984), 167–189.
- [30] Kevin Kappelmann. 2019. Decision Procedures for Guarded Logics. *CoRR* abs/1911.03679 (2019), 92.
- [31] Deepak Kapur and Paliath Narendran. 1986. NP-Completeness of the Set Unification and Matching Problems. In *CADE*. Springer, 489–495.
- [32] Alon Y. Levy. 2000. *Logic-Based Techniques in Data Integration*. Kluwer Academic Publishers, Norwell, MA, USA, 575–595.
- [33] Thomas Lukasiewicz, Andrea Cali, and Georg Gottlob. 2012. A General Datalog-Based Framework for Tractable Query Answering over Ontologies. *J. Web Semant.* 14, 0 (2012), 57–83.
- [34] Bruno Marnette. 2012. Resolution and Datalog Rewriting Under Value Invention and Equality Constraints. *CoRR* abs/1212.0254 (2012), 12.
- [35] William McCune and Larry Wos. 1997. Otter—The CADE-13 Competition Incarnations. *J. Autom. Reason.* 18, 2 (1997), 211–220.
- [36] M. Meier. 2014. The backchase revisited. *VLDB J.* 23, 3 (2014), 495–516.
- [37] Boris Motik. 2006. *Reasoning in Description Logics using Resolution and Deductive Databases*. Ph.D. Dissertation. Karlsruhe Institute of Technology, Karlsruhe, Germany. Retrieved July 4, 2022 from <http://digbib.ubka.uni-karlsruhe.de/volltexte/1000003797>
- [38] Boris Motik. 2022. *The KAON2 System*. Karlsruhe Institute of Technology. Retrieved July 4, 2022 from <http://kaon2.semanticweb.org/>
- [39] Oxford KR group. 2021. *Oxford Ontology Library*. Oxford University. Retrieved July 4, 2022 from <http://krr-nas.cs.ox.ac.uk/ontologies/>
- [40] Mike Paterson and Mark N. Wegman. 1978. Linear Unification. *J. Comput. Syst. Sci.* 16, 2 (1978), 158–167.
- [41] John Alan Robinson. 1965. A Machine-Oriented Logic Based on the Resolution Principle. *J. ACM* 12, 1 (1965), 23–41.
- [42] Iztok Savnik. 2013. Index Data Structure for Fast Subset and Superset Queries. In *CD-ARES*. Springer, 134–148.
- [43] Stephan Schulz. 2013. Simple and Efficient Clause Subsumption with Feature Vector Indexing. In *Automated Reasoning and Mathematics*. Springer, 45–67.
- [44] Mark E. Stickel. 1989. *The Path-Indexing Method for Indexing Terms*. Technical Report. SRI. Retrieved July 27, 2022 from <https://apps.dtic.mil/sti/citations/ADA460990>
- [45] Oxford Semantic Technologies. 2022. *The RDFS System*. Oxford Semantic Technologies. Retrieved July 4, 2022 from <https://www.oxfordsemantic.tech/>
- [46] Moshe Y. Vardi. 1997. Why Is Modal Logic so Robustly Decidable?. In *DIMACS Workshop*, Vol. 31. American Mathematical Society, 149–184.
- [47] Roberto De Virgilio, Giorgio Orsi, Letizia Tanca, and Riccardo Torlone. 2012. NYAYA: A System Supporting the Uniform Management of Large Sets of Semantic Data. In *ICDE*. IEEE Computer Society, 1309–1312.
- [48] Zhe Wang, Peng Xiao, Kewen Wang, Zhiqiang Zhuang, and Hai Wan. 2021. Query Answering for Existential Rules via Efficient Datalog Rewriting. In *IJCAI*. ijcai.org, 1933–1939.
- [49] Sen Zheng and Renate A. Schmidt. 2020. Deciding the Loosely Guarded Fragment and Querying Its Horn Fragment Using Resolution. In *AAAI*. AAAI Press, 3080–3087.