# Chase Termination Beyond Polynomial Time

PHILIPP HANISCH, Knowledge-Based Systems Group, TU Dresden, Germany

MARKUS KRÖTZSCH, Knowledge-Based Systems Group, TU Dresden, Germany

The *chase* is a widely implemented approach to reason with tuple-generating dependencies (tgds), used in data exchange, data integration, and ontology-based query answering. However, it is merely a semi-decision procedure, which may fail to terminate. Many decidable conditions have been proposed for tgds to ensure chase termination, typically by forbidding some kind of "cycle" in the chase process. We propose a new criterion that explicitly allows some such cycles, and yet ensures termination of the standard chase under reasonable conditions. This leads to new decidable fragments of tgds that are not only syntactically more general but also strictly more expressive than the fragments defined by prior acyclicity conditions. Indeed, while known terminating fragments are restricted to PTime data complexity, our conditions yield decidable languages for any $k$-ExpTime. We further refine our syntactic conditions to obtain fragments of tgds for which an optimised chase procedure decides query entailment in PSpace or $k$-ExpSpace, respectively.

CCS Concepts: • **Theory of computation → Constraint and logic programming**; **Database constraints theory**; *Database query languages (principles)*; *Logic and databases*.

Additional Key Words and Phrases: tuple-generating dependencies, chase termination, restricted chase

## 1 INTRODUCTION

The *chase* [1, 6, 33] is an essential method for reasoning with constraints in databases, with application areas including data exchange [24], constraint implication [5], data cleansing [26], query optimization [2, 7, 36], query answering under constraints [13, 38], and ontological reasoning [4, 12]. The basis for this wide applicability is the chase's ability to compute a *universal model* for a set of constraints [21], which is either used directly (e.g., as a repaired database) or indirectly (e.g., for deciding query entailment).

However, universal models can be infinite, and chase termination is undecidable on a single database [5] as well as in its stricter *uniform* version over arbitrary databases [27, 29]. The root of this problem are a form of constraints known as *tuple-generating dependencies* (tgds) – or *existential rules*: Horn logic rules with existential quantification in conclusions –, since they may require additional domain elements (represented by *nulls*) to be satisfied.

A large body of research is devoted to finding decidable cases for which termination can be guaranteed, mainly by analysing the data flow, i.e., the propagation of nulls in the chase [14, 18, 22, 24, 31, 34].[1] Here we can distinguish graph-based abstractions, such as *weak acyclicity* [24], from materialisation-based approaches, such as MFA [18]. In general, decidable criteria are sufficient but not necessary for termination, but recent breakthroughs established decidability of termination for the linear, guarded, and sticky classes of tgds [9, 10, 28, 37].

---

[1]As a notable exception, control flow is analysed in the *graph of rule dependencies* [4].

Meanwhile, another productive line of recent research clarified the computational power of the chase, characterising the query functions that can be expressed by conjunctive queries under tgd constraints. This expressive power is bounded by data complexity, but can be lower: famously, Datalog does not capture all queries in P, not even those *closed under homomorphisms*[2] [20]. Results are more satisfying for tuple-generating dependencies. Not only do arbitrary tgds capture the recursively enumerable homomorphism-closed queries [39], but, surprisingly, the decidable homomorphism-closed queries can be captured by tgds for which the standard (a.k.a. *restricted*) chase terminates [8]. In other words, the standard chase over tgds without any extensions is a universal computational paradigm for database query answering. This is highly encouraging since a majority of chase implementations already support this version of the chase procedure [3, 7, 26, 36, 38, 41], often with favourable performance [6].

Naturally, tractable data complexity is often desirable in database applications, and has therefore been the focus of many works in the area. Unfortunately, the potential of chase-based computation for more complex computations has meanwhile been neglected. To our knowledge, the only line of research where the chase was used for harder-than-P computations relies on a method for modelling finite sets with tgds [15]. Practical feasibility was shown for ExpTime-complete ontology-based query answering [15], set-terms in answer set programming [25], complex values in Datalog [35], and why-provenance [23]. In essence, all of these works are based on a single set of tgds for which uniform termination is shown directly. Known chase termination criteria fail for this case, since they only recognise queries in P. In fact, most criteria are even known to describe fragments that do not increase upon the expressive power of Datalog [32, 42], This limitation is common to all tgd sets on which the semi-oblivious chase uniformly terminates. We are aware of only one approach so far that studies the more complicated standard chase at all [14], but which also remains in P.

We tackle this challenge with a new method that extends the graph-based termination criterion of *joint acyclicity* [31] with new decidable conditions that allow for some kinds of cycles. These conditions are specific to the standard chase, and enforce that tgd applications within a cycle are eventually blocked, possibly only after (double) exponentially many loops. Our conditions detect the uniform termination of the set-modelling tgds of Carral et al. [15], but significantly extend upon this baseline: even a single strongly connected component in the data-flow graph can correspond to a 2-ExpTime-complete query (whereas Carral et al. deal with exponentially many sets).

Leveraging again the graph-based view, we can further analyse the data flows between cyclic strongly connected components to describe decidable fragments of arbitrary multi-exponential data complexity. Our resulting decidable fragment of *saturating tgds* therefore has non-elementary complexity. This is already very general, especially when considering that capturing all decidable homomorphism-closed queries can in principle only be achieved with tgd fragments that are not even recursively enumerable [8].

The structure of the data-flow graph enables us to compute more precise $\kappa$-ExpTime bounds for any saturating tgd set. Refining this analysis further, we then identify cases where the complexity of query answering drops to $(\kappa - 1)$-ExpSpace. In particular, we therefore obtain a decidable fragment of uniformly terminating tgd sets with PSpace-complete query entailment. To the best of our knowledge, this is the first such fragment.

All of our results establish uniform termination of the standard chase under all chase strategies that prioritise Datalog rules (Krötzsch et al. call this the *Datalog-first* strategy [32]). By a recent result, this is a stronger requirement than termination under *some* strategies [16]. As of today,

---

[2]Such queries are monotone, i.e., intuitively do not need negation (but cf. [40]).

however, all known termination criteria imply Datalog-first termination, and preferring Datalog rules is also a common heuristic in practice.

In summary, our main contributions are as follows:

- In Section 3, we refine the *dependency graph* of Krötzsch and Rudolph [31] to capture data flow in more detail.
- In Section 4, we study the propagation of inferences between nulls related to strongly connected components in the extended dependency graph. We find conditions that suffice to prevent infinite repetitions of inference cycles, define the language of *saturating tgds*, and show uniform standard chase termination for this fragment.
- In Section 5, we analyse the exact complexity (and size) of the standard chase over saturating sets by assigning *ranks* to strongly connected components in the extended dependency graph. We differentiate the case of single and double exponential complexity for individual components, and we establish matching lower bounds to show that query entailment is $\kappa$-ExpTime-complete for tgd sets of rank $\kappa$.
- In Section 6, we describe conditions that reduce query entailment for tgd sets of rank $\kappa$ to $(\kappa - 1)$-ExpSpace. To this end, we discover a tree-like structure within the chase, and we define a syntactic condition called *path guardedness* that allows us to use an optimised chase procedure. Again, we establish matching lower bounds.

Detailed proofs are included in the appendix.

## 2 PRELIMINARIES

We consider a signature based on mutually disjoint, countably infinite sets of *constants* C, *variables* V, *predicates* P, and *nulls* N. Each predicate name $p \in$ P has an *arity* $\text{ar}(p) \geq 0$. *Terms* are elements of $V \cup N \cup C$. We use $t$ to denote a list $t_1, \ldots, t_{|t|}$ of terms, and similarly for special types of terms. An *atom* is an expression $p(t)$ with $p \in$ P, $t$ a list of terms, and $\text{ar}(p) = |t|$. An *interpretation* $\mathcal{I}$ is a set of atoms without variables. A *database* $\mathcal{D}$ is a finite interpretation without nulls, i.e., a finite set of *facts* (variable-free, null-free atoms). For an interpretation $\mathcal{I}$, we use $N(\mathcal{I}) = \{n \in N \mid p(t) \in \mathcal{I}, n \in t\}$ to denote the nulls used in $\mathcal{I}$.

*Rules.* A *tuple-generating dependency* (tgd) $\rho$ is a formula

$$\rho = \forall x, y. B[x, y] \rightarrow \exists v. H[y, v], \tag{1}$$

where $B$ and $H$ are conjunctions of atoms using only terms from C or from the mutually disjoint lists of variables $x, y, v \subseteq$ V. We call $B$ the *body* (denoted $\text{body}(\rho)$), $H$ the *head* (denoted $\text{head}(\rho)$), and $y$ the *frontier* of $\rho$. We may treat conjunctions of atoms as sets, and we omit universal quantifiers in tgds. We require that all variables in $y$ do really occur in $B$ (*safety*). A tgd without existential quantifiers is a *Datalog rule*.

*Renamings and Substitutions.* Without loss of generality, we require that variables in tgd sets $\Sigma$ are *renamed apart*, i.e., each variable $x \in$ V in $\Sigma$ is bound by a unique quantifier in a unique tgd $\rho_x \in \Sigma$. A *substitution* is a partial mapping $\sigma : V \rightarrow C \cup V \cup N$. As usual, $x\sigma = \sigma(x)$ if $\sigma$ is defined on $x$, and $x\sigma = x$ otherwise. For a logical expression $\alpha[x]$, the expression $\alpha\sigma$ is obtained by simultaneously replacing each variable $x$ by $x\sigma$. Given a list of terms $t = t_1, \ldots, t_{|t|}$, we write $\alpha[x/t]$ for the expression $\alpha\{x_1 \mapsto t_1, \ldots, x_{|x|} \mapsto t_{|x|}\}$. If $x$ is clear from the context and no confusion is likely, we write $\alpha[t]$ for $\alpha[x/t]$.

*Semantics.* We consider a standard first-order semantics. A *match* of a tgd $\rho$ as in (1) in an interpretation $\mathcal{I}$ is a substitution $\sigma$ that maps $x \cup y$ to terms in $\mathcal{I}$, such that $B\sigma \subseteq \mathcal{I}$. A match is *satisfied* if it can be extended to a substitution $\sigma'$ over $x \cup y \cup v$ such that $H\sigma' \subseteq \mathcal{I}$. A tgd $\rho$

is *satisfied* by $\mathcal{I}$, written $\mathcal{I} \models \rho$, if all matches of $\rho$ on $\mathcal{I}$ are satisfied. A fact $\alpha$ is satisfied in $\mathcal{I}$, written $\mathcal{I} \models \alpha$, if $\alpha \in \mathcal{I}$. Satisfaction extends to sets of tgds and facts as usual. A tgd or fact $\alpha$ is *entailed* by tgd set $\Sigma$ and database $\mathcal{D}$ if $\mathcal{I} \models \alpha$ for all $\mathcal{I}$ with $\mathcal{I} \models \Sigma \cup \mathcal{D}$.

*Reasoning with the Chase.* An important reasoning task for tgds is conjunctive query (CQ) answering, which can further be reduced to the entailment of Boolean CQs (BCQs), which are formulas $\exists z.Q[z]$ with $Q$ a conjunction of null-free atoms. This task is undecidable in general. A sound and complete (but not always terminating) class of reasoning procedures is the *chase*, which exists in many variants. We are interested in the *standard chase* (a.k.a. *restricted chase*) under *Datalog-first* strategies.

**Definition 1.** *A* (standard) chase sequence *for a database $\mathcal{D}$ and a tgd set $\Sigma$ is a potentially infinite sequence of interpretations $\mathcal{D}^0, \mathcal{D}^1, \ldots$ such that*

(1) $\mathcal{D}^0 = \mathcal{D}$;

(2) *for every $\mathcal{D}^{i+1}$ with $i \geq 0$, there is a match $\sigma$ for some tgd $\rho = B[x,y] \rightarrow \exists v.H[y,v] \in \Sigma$ in $\mathcal{D}^i$ such that both of the following hold true:*

   (a) *$\sigma$ is an unsatisfied match in $\mathcal{D}^i$ (i.e., $\sigma$ cannot be extended to a substitution $\sigma^+$ with $H\sigma^+ \subseteq \mathcal{D}^i$),*

   (b) *$\mathcal{D}^{i+1} = \mathcal{D}^i \cup H[\sigma^+(y), \sigma^+(v)]$, where $\sigma^+$ is such that $\sigma^+(y) = \sigma(y)$ for all $y \in y$, and for all $v \in v$, $\sigma^+(v) \in \mathbf{N}$ is a distinct null not occurring in $\mathcal{D}^i$;*

   *we then say that $\rho$ was applied in step $i$, and we define $\mathrm{tgd}[i] := \rho$, $\sigma[i] := \sigma$, and $\sigma^+[i] := \sigma^+$;*

(3) *if a tgd with existential variables is applied in step $i$, then $\mathcal{D}^i$ must satisfy all Datalog rules in $\Sigma$;*

(4) *if $\sigma$ is a match for a tgd $\rho \in \Sigma$ and $\mathcal{D}^i$ ($i \geq 0$), then there is $j > i$ such that $\sigma$ is satisfied in $\mathcal{D}^j$.*

*Item* (3) *requires rule applications to follow a* Datalog-first *strategy, and item* (4) *ensures* fairness. *The* (standard) chase *for such a chase sequence then is* $\mathrm{chase}(\Sigma, \mathcal{D}) = \bigcup_{i \geq 0} \mathcal{D}^i$.

A BCQ $q$ is entailed by $\Sigma$ and $\mathcal{D}$ if and only if $\mathrm{chase}(\Sigma, \mathcal{D}) \models q$. If the chase terminates, this can be determined from the (finite) $\mathrm{chase}(\Sigma, \mathcal{D})$. Termination may depend on the chosen order of tgd applications. The Datalog-first strategy (3) is a common heuristic that tends to improve termination in practice, although one can construct examples where this is not the case [16]. Since Datalog rules can only be applied finitely many times, Datalog-first does not impair fairness (4).

## 3 THE LABELLED DEPENDENCY GRAPH

The *existential dependency graph* is used to analyse the data flow between existential variables in a tgd set [31]. In particular, a tgd set is *jointly acyclic* if this graph does not have cycles. In this section, we recall and slightly extend this approach to better suit our needs.

The following definition mostly follows Krötzsch and Rudolph [31], but adds variables as labels to the edges in the graph. Recall that we assume tgd sets to be renamed apart, so that variables $x$ can be used to identify tgds $\rho_x$.

**Definition 2.** *Let $\Sigma$ be a tgd set. A* predicate position *is a pair $\langle p, i \rangle \in \mathbf{P} \times \mathbb{N}$ with $1 \leq i \leq \mathrm{ar}(p)$. For a variable $x$ in $\Sigma$, let $\mathrm{Pos}_x^B$ (resp. $\mathrm{Pos}_x^H$) be the set of all predicate positions where $x$ occurs in the body (resp. head) of its unique tgd $\rho_x \in \Sigma$.*

*For an existential variable $v$ in $\Sigma$, let $\Omega_v$ be the smallest set of positions such that (i) $\mathrm{Pos}_v^H \subseteq \Omega_v$, and (ii) for every universal variable $x$, $\mathrm{Pos}_x^B \subseteq \Omega_v$ implies $\mathrm{Pos}_x^H \subseteq \Omega_v$.*

*The* labelled existential dependency graph $\mathrm{LXG}(\Sigma)$ *of $\Sigma$ is a directed graph with the existentially quantified variables of $\Sigma$ as its vertices and an edge $v \xrightarrow{y} w$ for every tgd $\rho_w \in \Sigma$ with a frontier variable $y$ such that $\mathrm{Pos}_y^B \subseteq \Omega_v$.*

**Example 3.** *This running example illustrates several notions below, hence is not minimal. Consider a database $\mathcal{D}_{lvl} = \{first(1), last(\ell)\} \cup \{next(i, i+1) \mid 1 \leq i < \ell\}$, which defines a total strict order of "levels" $1, \ldots, \ell$. $\Sigma$ consists of the following tgds with constants $F$ (false) and $T$ (true):*

$$first(z) \rightarrow lvl(F, z) \wedge lvl(T, z) \tag{2}$$

$$lvl(\bar{x}_1, z) \wedge lvl(\bar{x}_2, z) \rightarrow \exists v.cat(\bar{x}_1, \bar{x}_2, z, v) \tag{3}$$

$$cat(\bar{x}_1, \bar{x}_2, z, x) \rightarrow part(\bar{x}_1, x) \wedge part(\bar{x}_2, x) \tag{4}$$

$$cat(\bar{x}_1, \bar{x}_2, z, x) \wedge next(z, z_+) \rightarrow \exists \bar{w}.up(x, z_+, \bar{w}) \tag{5}$$

$$cat(\bar{x}_1, \bar{x}_2, z, x) \wedge next(z, z_+) \wedge up(x, z_+, \bar{x}) \rightarrow lvl(\bar{x}, z_+) \tag{6}$$

*Overlined variables denote sequences of $F$ and $T$, where $lvl(\bar{x}, z)$ says that sequence $\bar{x}$ has length $2^z$. Initial sequences have length 1 (2). Longer sequences emerge from concatenations (3) with equal-length parts (4). The tgd (5) promotes concatenations $x$ to sequences $\bar{w}$ on the next level $z_+$. We mark such new sequences separately (6), since this will later be useful. The tgds (2)–(6) on $\mathcal{D}_{lvl}$ therefore construct all $F$-$T$-sequences of length $2^\ell$, i.e., $2^{2^\ell}$ distinct nulls.*

*Then $\Omega_v = \{\langle cat, 4\rangle, \langle part, 2\rangle, \langle up, 1\rangle\}$ and $\Omega_{\bar{w}} = \{\langle up, 3\rangle, \langle lvl, 1\rangle, \langle cat, 1\rangle, \langle cat, 2\rangle\}$. The three edges of $\mathrm{LXG}(\Sigma)$ are $v \xrightarrow{x(5)} \bar{w}$, $\bar{w} \xrightarrow{\bar{x}_1(3)} v$, and $\bar{w} \xrightarrow{\bar{x}_2(3)} v$, where we disambiguate some variables with the numbers of their tgds.*

Definition 2 slightly sharpens the original definition by introducing edges only based on frontier variables. Moreover, by adding labels, a single edge in the original graph now corresponds to one or more edges in $\mathrm{LXG}(\Sigma)$. Therefore, if the existential dependency graph contains no cycles (i.e., $\Sigma$ is jointly acyclic), then the labelled existential dependency graph does not contain cycles either.

$\mathrm{LXG}(\Sigma)$ is useful since it over-estimates the possible data flow in the computation of $\mathrm{chase}(\Sigma, \mathcal{D})$. To make this precise, note that any null $n$ in $\mathrm{chase}(\Sigma, \mathcal{D})$ is introduced in a chase step $i$ as a fresh value $n = \sigma^+[i](v)$ of some existential variable $v$ in $\mathrm{tgd}[i]$. We write $\mathrm{var}(n)$ for this $v$, and $t \overset{y}{\twoheadrightarrow} n$ to indicate that $\sigma[i](y) = t$ for some term $t \in \mathbf{C} \cup \mathbf{N}$ and frontier variable $y$ of $\mathrm{tgd}[i]$.

**Lemma 4.** *If $n_1 \overset{y}{\twoheadrightarrow} n_2$ in $\mathrm{chase}(\Sigma, \mathcal{D})$ for nulls $n_1, n_2 \in \mathbf{N}$, then there is an edge $\mathrm{var}(n_1) \xrightarrow{y} \mathrm{var}(n_2)$ in $\mathrm{LXG}(\Sigma)$.*

**Lemma 5.** *If $\mathrm{chase}(\Sigma, \mathcal{D})$ is infinite, then $\mathrm{chase}(\Sigma, \mathcal{D})$ contains an infinite chain $n_0 \overset{y_1}{\twoheadrightarrow} n_1 \overset{y_2}{\twoheadrightarrow} \cdots$.*

By Lemma 4, the infinite path of Lemma 5 corresponds to an infinite path in $\mathrm{LXG}(\Sigma)$, which must therefore, since it is finite, contain a cycle. Hence, if $\mathrm{LXG}(\Sigma)$ is acyclic, $\mathrm{chase}(\Sigma, \mathcal{D})$ is finite.

**Example 6.** *For $\Sigma$ as in Example 3, $\mathrm{LXG}(\Sigma)$ is cyclic. Indeed, there are databases $\mathcal{D}$ for which $\mathrm{chase}(\Sigma, \mathcal{D})$ is infinite, e.g., for $\mathcal{D} = \{first(1), last(1), next(1, 1)\}$.*

## 4 TERMINATION WITH CYCLES

In this section, we establish decidable criteria to show that the chase on a tgd set is guaranteed to be finite, for all input databases, even though the existential dependency graph has some cycles.

Whenever a tgd $B[\boldsymbol{x}, \boldsymbol{y}] \rightarrow \exists v.H[\boldsymbol{y}, v]$ was applied in step $i$ of $\mathrm{chase}(\Sigma, \mathcal{D})$, the conjunctive query $(B \wedge H)[\boldsymbol{x}, \boldsymbol{y}, v]$ has the answer $\sigma^+[i]$ over $\mathrm{chase}(\Sigma, \mathcal{D})$. Likewise, a chain of chase steps (or path in $\mathrm{LXG}(\Sigma)$) leads to a match for a larger query, defined next.

**Definition 7.** *Consider a path $\mathfrak{p} = v_0 \xrightarrow{y_1} v_1 \xrightarrow{y_2} \cdots \xrightarrow{y_k} v_k$ in $\mathrm{LXG}(\Sigma)$ for variables $v_0, v_i, y_i \in \mathbf{V}$ $(1 \leq i \leq k)$. Let $\rho_i : B_i[\boldsymbol{x_i}, \boldsymbol{y_i}] \rightarrow \exists v_i.H_i[\boldsymbol{y_i}, v_i]$ be a variant of the tgd $\rho_{v_i}$ of variable $v_i$, where variables have been bijectively renamed such that tgds for different steps do not share variables. For $1 \leq i \leq k$, let $\tilde{y}_i$ (and $\tilde{v}_i$) denote the renamed version of $y_i$ (and $v_i$), and let $\tilde{y}_{k+1}$ denote a fresh variable.*
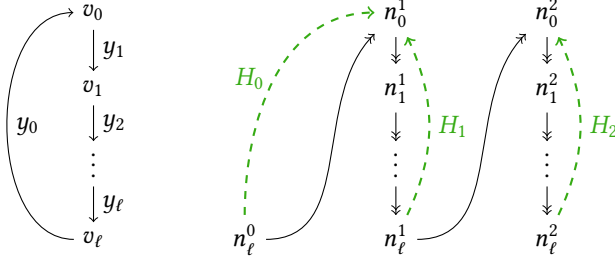
Fig. 1. Dependency cycle in LXG($\Sigma$) (left) and corresponding chain of derived nulls in chase($\Sigma, \mathcal{D}$) (right); $H_0, H_1$, and $H_2$ denote variants of a tgd head to illustrate propagation

We define Path($\mathfrak{p}$) := $\bigwedge_{i=1}^{k}(B_i \wedge H_i[\tilde{v}_i/\tilde{y}_{i+1}])$. Given $1 \leq i \leq k$, we write Path($\mathfrak{p}$)$_{B,i}$ := $B_i$ and Path($\mathfrak{p}$)$_{H,i}$ := $H_i[\tilde{v}_i/\tilde{y}_{i+1}]$ for the $i$-th body and head conjunction, respectively, with the renamed variables.

**Example 8.** LXG($\Sigma$) of Example 3 has a path $\mathfrak{p} = v \xrightarrow{x} \bar{w} \xrightarrow{\bar{x}_1} v$, where we omit the numbers of the tgds from variables for simplicity. Then Path($\mathfrak{p}$) = $cat(\bar{x}'_1, \bar{x}'_2, z', x) \wedge next(z', z'_+) \wedge up(x, z'_+, \bar{x}_1) \wedge lvl(\bar{x}_1, z'') \wedge lvl(\bar{x}''_2, z'') \wedge cat(\bar{x}_1, \bar{x}''_2, z'', y_3)$. Variables marked with ′ and ″ stem from renamed variants of tgds (3) and (5), respectively. Similarly, Path($\mathfrak{p}$)$_{H,1}$ = $up(x, z'_+, \bar{x}_1)$.

**Lemma 9.** Let $\mathfrak{c} = n_0 \xrightarrow{y_1} n_1 \xrightarrow{y_2} \cdots \xrightarrow{y_k} n_k$ be a chain in chase($\Sigma, \mathcal{D}$), with $n_k$ derived in the $m$-th chase step $\mathcal{D}^m$. There is a path $\mathfrak{p} = \text{var}(v_0) \xrightarrow{y_1} \text{var}(n_1) \xrightarrow{y_2} \cdots \xrightarrow{y_k} \text{var}(n_k)$ in LXG($\Sigma$), and $\mathcal{D}^m \models$ Path($\mathfrak{p}$)$\theta$ holds for substitution $\theta$ defined as follows: for each $n_i$ derived in chase step $j$, and each variable $x$ in tgd[$j$] that was renamed to $\tilde{x}$ in Path($\mathfrak{c}$), $\theta(\tilde{x}) = \sigma^+[j](x)$.

Lemma 9 ensures that every chain of nulls in the chase is accompanied by facts of the form Path($\mathfrak{p}$), connecting all nulls of the chain. Figure 1 sketches this situation for a cyclic path (left). A corresponding chain of nulls $n_\ell^0 \xrightarrow{y_0} n_0^1 \xrightarrow{y_1} \ldots$ is sketched on the right, where vertical positions indicate existential variables, i.e., var($n_i^j$) = $v_i$. Moreover, a solid arrow like $n_\ell^0 \twoheadrightarrow n_0^1$ also corresponds to facts of the form Path($n_\ell^0 \xrightarrow{y_0} n_0^1$) that connect the respective nulls. The dashed arrows are explained below.

We can use the facts of Path($\mathfrak{p}$) to infer additional information that may help with chase termination for cyclic paths. Indeed, additional information may prevent tgd applications if the head of a tgd is already entailed (Definition 1 (2.a)).

To understand this better, let's denote the tgd for edge $v_\ell \xrightarrow{y_0} v_0$ in Figure 1 as $\rho = B[\boldsymbol{x}, \boldsymbol{y}] \rightarrow \exists \boldsymbol{v}.H[\boldsymbol{y}, \boldsymbol{v}]$ with $y_0 \in \boldsymbol{y}$ and $v_0 \in \boldsymbol{v}$. If $i$ is the chase step that introduced $n_0^1$, then $H\sigma^+[i] \subseteq \mathcal{D}^{i+1}$. In Figure 1, $H_0$ represents the facts $H\sigma^+[i]$, which involve (among others) the terms $n_\ell^0$ and $n_0^1$, as suggested by the dashed arrow.

Clearly, we cannot apply $\rho$ twice for the same substitution of its frontier: for all $i' > i$ and substitutions $\sigma$ with $\boldsymbol{y}\sigma = \boldsymbol{y}\sigma[i]$, we have $\mathcal{D}^{i'} \models \exists \boldsymbol{v}.H\sigma$ due to the presence of $H_0$. However, a cycle in LXG($\Sigma$) can lead to a chain of tgd applications as in Figure 1, where $\rho$ is applied to different frontiers in several steps. Indeed, if $j$ is the chase step that introduced $n_0^2$, then $y_0\sigma[i] \neq y_0\sigma[j]$, so the matches differ on at least this frontier variable. If we write $\boldsymbol{y}^- := \boldsymbol{y} \setminus \{y_0\}$ for the frontier of $\rho$ without $y_0$, we can think of $\boldsymbol{y}^-\sigma[i]$ as the "context" for which $\rho$ was applied to $n_\ell^0$ in step $i$.

Our goal is that $\rho$ can never be applied twice to the same context within a single chain. To ensure this, we would like the chase to derive additional facts $H_1$ and $H_2$ as indicated in Figure 1. Fixing a variable order $H[y_0, \boldsymbol{y}^-, \boldsymbol{v}]$, these facts are $H_1 := H[n_\ell^1, \boldsymbol{y}^- \sigma[i], \boldsymbol{v}\sigma^+[i]]$ and $H_2 := H[n_\ell^2, \boldsymbol{y}^- \sigma[i], \boldsymbol{v}\sigma^+[j]]$. Chains like those in Figure 1, even if finite, can become rather long, and we need facts $H_p$ to be propagated to all nulls $n_\ell^p$. We therefore require two kinds of conditions: (i) a base case that turns recently derived (forward) $H_0$ into a (backwards) $H_1$, and (ii) an inductive step that propagates a (backwards) $H_p$ to another (backwards) $H_{p+1}$. The next definition spells out the two conditions. We generalise slightly by allowing that, instead of applying a single tgd $\rho$ several times, the propagation can occur between different tgds along a path.

**Definition 10.** *For $p \in \{*, a, b\}$, let $u^p \xrightarrow{y^p} v^p \in \mathrm{LXG}(\Sigma)$ be such that the corresponding tgd $\rho^p$ has a head $\exists v^p, \boldsymbol{w}^p . H^p[y^p, \boldsymbol{z}^p, v^p, \boldsymbol{w}^p]$. Moreover, let $\Sigma_{\mathrm{DL}}$ denote the set of Datalog rules in $\Sigma$.*

*A path $\mathfrak{p} = u^* \xrightarrow{y^*} v^* \xrightarrow{z_2} \dots \xrightarrow{z_\ell} w$ is* base-propagating *if*

$$\Sigma_{\mathrm{DL}} \models \mathrm{Path}(\mathfrak{p}) \rightarrow H^*[\tilde{y}_{\ell+1}, \tilde{z}_1, \tilde{y}_2, \tilde{w}_1] \tag{7}$$

*where the conclusion is the path's first head conjunction $\mathrm{Path}(\mathfrak{p})_{H,1} = H^*[\tilde{y}_1, \tilde{z}_1, \tilde{y}_2, \tilde{w}_1]$ with $\tilde{y}_1$ replaced by the variable $\tilde{y}_{\ell+1}$ that occurs in $\mathrm{Path}(\mathfrak{p})_{H,\ell} = H_\ell[\tilde{y}_\ell, \tilde{z}_\ell, \tilde{y}_{\ell+1}, \tilde{w}_\ell]$.*

*A composite path $\mathfrak{p}$ of the form $\mathfrak{p}^a \mathfrak{p}^b$ with*

$$\mathfrak{p}^a = u^a \xrightarrow{y^a} v^a \xrightarrow{z_2} \dots \xrightarrow{z_\ell} u^b \text{ and } \mathfrak{p}^b = u^b \xrightarrow{y^b} v^b \xrightarrow{z'_2} \dots \xrightarrow{z'_k} w$$

*is* step-propagating *for $\rho^*$ (or, equivalently, for $u^* \xrightarrow{y^*} v^*$) if*

$$\Sigma_{\mathrm{DL}} \models \mathrm{Path}(\mathfrak{p}) \wedge H^*[\tilde{y}_{\ell+1}, \boldsymbol{x}_z, \tilde{y}_2, \boldsymbol{x}_w] \rightarrow H^*[\tilde{y}_{\ell+k+1}, \boldsymbol{x}_z, \tilde{y}_{\ell+2}, \boldsymbol{x}_w] \tag{8}$$

*where $\boldsymbol{x}_z$ and $\boldsymbol{x}_w$ are lists of fresh variables of length $|\boldsymbol{x}_z| = |\boldsymbol{z}^*|$ and $|\boldsymbol{x}_w| = |\boldsymbol{w}^*|$, respectively; $\tilde{y}_{\ell+k+1}$ is the renamed version of the final variable $w$ in $\mathrm{Path}(\mathfrak{p})$; and the other mentioned variables stem from the atom sets $\mathrm{Path}(\mathfrak{p})_{H,1} = H^a[\tilde{y}_1, \tilde{z}_1, \tilde{y}_2, \tilde{w}_1]$ and $\mathrm{Path}(\mathfrak{p})_{H,\ell+1} = H^b[\tilde{y}_{\ell+1}, \tilde{z}_{\ell+1}, \tilde{y}_{\ell+2}, \tilde{w}_{\ell+1}]$.*

Intuitively speaking, considering Figure 1, base propagation allows us to infer $H_1$ from $H_0$, and step propagation allows us to infer $H_2$ from $H_1$. In contrast to the simplified situation in the figure, the definition clarifies that the propagated head $H^*$ is not necessarily the head of the tgds $\rho^a$ and $\rho^b$ that occur in the current piece of chain we consider.

**Example 11.** *We extend $\Sigma$ from Example 3 by two additional tgds:*

$$up(y_1, z, \bar{y}_2) \wedge part(\bar{y}_2, y_3) \rightarrow up(y_3, z, \bar{y}_2) \tag{9}$$

$$up(y_3, z, \bar{y}_2) \wedge part(\bar{y}_2, y_3) \wedge up(y_3, z', \bar{y}_4) \wedge part(\bar{y}_4, y_5) \rightarrow up(y_5, z, \bar{y}_4) \tag{10}$$

$\mathrm{LXG}(\Sigma)$ *remains as in Example 3. The path $\mathfrak{p} = v \xrightarrow{x} \bar{w} \xrightarrow{\bar{x}_1} v$ from Example 8 is base-propagating. Indeed, tgds (4) and (9) entail the required tgd $\mathrm{Path}(\mathfrak{p}) \rightarrow up(y_3, z'_+, \bar{x}_1)$, in fact even the stronger tgd $up(x, z'_+, \bar{x}_1) \wedge cat(\bar{x}_1, \bar{x}''_2, z'', y_3) \rightarrow up(y_3, z'_+, \bar{x}_1)$. Similarly, for $\mathfrak{p}^a = v \xrightarrow{x} \bar{w} \xrightarrow{\bar{x}_1} v$ and $\mathfrak{p}^b = v \xrightarrow{x} \bar{w} \xrightarrow{\bar{x}_2} v$, the path $\mathfrak{p}^a \mathfrak{p}^b$ is step-propagating for (5) (i.e., for $v \xrightarrow{x} \bar{w}$) using tgds (4) and (10). Note that variables $y_i$ in (10) match the path labels in Definition 10.*

*As we will see below, the extended example achieves universal termination, and in particular also terminates for the database of Example 6. For the case that next is a strict order, tgd (6) ensures that each sequence is assigned a unique level, even with the additional up-facts from tgds (9) and (10).*

The complexity of checking Definition 10 is dominated by the ExpTime complexity of Datalog.

**Lemma 12.** *Checking whether a path $\mathfrak{p}$ is base-propagating (or step-propagating) is ExpTime-complete, and P-complete with respect to the length of $\mathfrak{p}$.*

The conditions of Definition 10 have the desired impact on the chase: if we find a sequence of nulls that corresponds (by Lemma 4) to a path that is propagating, then satisfied head atoms in the chase are propagated accordingly. Moreover, since propagation is defined based on $\Sigma_{DL}$, the Datalog-first chase ensures that the propagation happens before further nulls are introduced. To ensure termination, we require that the cycles in each strongly connected component in $LXG(\Sigma)$ can be broken by removing a set of edges $E$ that are connected by propagating paths:

**Definition 13.** *Let $C$ be a strongly connected component in $LXG(\Sigma)$, and let $E$ be a set of edges in $C$. An $\bar{E}$-path is a path in $C$ that (i) contains no edges from $E$, (ii) starts in some $s \in \{v \mid u \xrightarrow{y} v \in E\}$, and (iii) ends in some $t \in \{u \mid u \xrightarrow{y} v \in E\}$. Moreover, $C$ is $E$-saturating if*

*(1) $C$ without the edges of $E$ is acyclic;*

*(2) for all $u \xrightarrow{x} v, u' \xrightarrow{x'} v \in E$, we have $x = x'$;*

*(3) all paths $e\mathfrak{p}$ in $C$, such that $e \in E$ and $\mathfrak{p}$ is an $\bar{E}$-path, are base-propagating; and*

*(4) all paths $e^a \mathfrak{p}_1 e^b \mathfrak{p}_2$ in $C$, such that $e^a, e^b \in E$ and $\mathfrak{p}_1, \mathfrak{p}_2$ are $\bar{E}$-paths, are step-propagating for every $e \in E$.*

$\Sigma$ is *saturating* if all strongly connected components in $LXG(\Sigma)$ are $E$-saturating for some $E$.

**Example 14.** *The extended set of tgds from Example 11 is saturating. $LXG(\Sigma)$ has a single strongly connected component (cf. Example 3) where Definition 13 is satisfied for $E = \{v \xrightarrow{x} \bar{w}\}$. For $e = v \xrightarrow{x} \bar{w}$, (3) applies to two paths $e\mathfrak{p}$ and (4) to four paths $e\mathfrak{p}e\mathfrak{p}'$, where $\mathfrak{p}, \mathfrak{p}' \in \{\bar{w} \xrightarrow{\bar{x}_i} v \mid i \in \{1,2\}\}$. Propagation follows as in Example 11.*

**Theorem 15.** *Deciding if $C$ is $E$-saturating is* ExpTime*-complete in the size of $C$. The same complexity applies to deciding if a tgd set $\Sigma$ is saturating.*

In practice, the exponential factors in Theorem 15 might be small, already because strongly connected components are often small. The next result states that tgds in the edge sets $E$ can only be applied at most once for each substitution of their "context" variables – an essential ingredient for termination.

**Lemma 16.** *Let $C$ be an $E$-saturating strongly connected component in $LXG(\Sigma)$, and let $u \xrightarrow{y} v \in E$ have the corresponding tgd $\rho_v$ with head $\exists v, w.H[y, z, v, w]$. For every database $\mathcal{D}$ and chain of nulls $n_0 \xrightarrow{y} n_1 \xrightarrow{x_2} \ldots \xrightarrow{x_{\ell-1}} n_{\ell-1} \xrightarrow{y} n_\ell$ in chase$(\Sigma, \mathcal{D})$ such that $n_0 = y\sigma[a]$, $n_1 = v\sigma^+[a]$, $n_{\ell-1} = y\sigma[b]$, and $n_\ell = v\sigma^+[b]$ for chase steps $a < b$, we have $z\sigma[a] \neq z\sigma[b]$.*

The main result of this section is as follows. A more detailed analysis follows in the next section.

**Theorem 17.** *If $\Sigma$ is saturating, then* chase$(\Sigma, \mathcal{D})$ *is finite for all databases $\mathcal{D}$.*

**Example 18.** *Theorem 17 also subsumes previous termination results by Carral et al. [15]. The following tgds captures the essence of their approach of simulating finite sets in tgds:*

$$elem(x) \wedge set(S) \rightarrow \exists v. \, set(S) \wedge su(x, S, v) \wedge su(x, v, v) \tag{11}$$

$$su(x, S, T) \wedge su(y, S, S) \rightarrow su(y, T, T) \tag{12}$$

*The database can provide elem-facts that define a domain of elements, and a fact set($\emptyset$). The tgd (11) constructs new "sets" by creating facts $su(x, S, T)$, which can be read as $\{x\} \cup S = T$. In particular, $su(x, S, S)$ means $x \in S$. The tgd (12) propagates memberships $x \in S$ from $S$ to a direct superset $T$, which we recognise as a special case of step propagation. Indeed, the only dependency here is $v \xrightarrow{S} v$, and all paths $\mathfrak{p}_{(1/2)}$ in Definition 13 are empty. Base propagation is achieved by the atom $su(x, v, v)$ in (11) without requiring any Datalog rules.*

**Example 19.** *Definition 13 (2) excludes "confluent" edges with distinct labels: being based on single frontier variables, our propagation conditions do not suffice for this case. For example, consider the tgd set $\Sigma$ that extend the tgd set of Example 18:*

$$set(S) \rightarrow elem(S) \tag{13}$$

$$su(x, S, T) \rightarrow su(T, S, T) \tag{14}$$

$$su(x, S, T) \wedge su(y, x, x) \rightarrow su(y, T, T) \tag{15}$$

$$su(x, S, T) \wedge su(S, y, S) \rightarrow su(T, y, T) \tag{16}$$

$$su(x, S, T) \wedge su(x, y, x) \rightarrow su(T, y, T) \tag{17}$$

*We allow "sets" to be used as "elements" (13), thereby creating a confluent E-edge. The tgd (14) ensures that the new E-edge is base-propagating, and tgds (15) − (17) ensure that the additional pairs of E-edges are step-propagating. For $\mathcal{D} = \{set(\emptyset)\}$, $\mathrm{chase}(\Sigma, \mathcal{D})$ is infinite.*

However, strongly connected components may have arbitrary confluent edges that are not in $E$, as in Example 14, and $E$ itself may have confluent edges that are using the same label.

## 5 COMPLEXITY OF THE SATURATING CHASE

Next, we refine Theorem 17 by deriving specific bounds for the size of the chase over saturating tgd sets $\Sigma$, based on the structure of $\mathrm{LXG}(\Sigma)$. For a vertex $v$ of $\mathrm{LXG}(\Sigma)$, we write $\mathrm{SCC}(v)$ for the strongly connected component that contains $v$, and $\mathrm{SCC}(\mathrm{LXG}(\Sigma))$ for the set of all strongly connected components. An edge $u \xrightarrow{y} w$ is *incoming* for $\mathrm{SCC}(v)$ if $u \notin \mathrm{SCC}(v)$ and $w \in \mathrm{SCC}(v)$; in this case we write $\mathrm{SCC}(u) \prec \mathrm{SCC}(w)$ (where $\mathrm{SCC}(w) = \mathrm{SCC}(v)$). The transitive reflexive closure of $\prec$ is the usual induced partial order on $\mathrm{SCC}(\mathrm{LXG}(\Sigma))$.

**Definition 20.** *For vertex $v$ of $\mathrm{LXG}(\Sigma)$, we define the label set $\lambda(v) = \{x \mid u \in \mathrm{SCC}(v), u \xrightarrow{x} v \text{ in } \mathrm{LXG}(\Sigma)\}$, and* confluence $\mathrm{conf}(v) = |\lambda(v)|$. *For $C \in \mathrm{SCC}(\mathrm{LXG}(\Sigma))$, we define $\mathrm{conf}(C) = \max\{\mathrm{conf}(u) \mid u \in C\}$. $C$ is* homogeneously confluent *if $\mathrm{conf}(C) = 1$.*

Note that $\mathrm{conf}(C) = 0$ implies that $C$ is trivial, i.e., a singleton set without any cycle (self loop).

**Definition 21.** *Consider a database $\mathcal{D}$ and $C \in \mathrm{SCC}(\mathrm{LXG}(\Sigma))$. A term $t$ in $\mathrm{chase}(\Sigma, \mathcal{D})$ is a $C$-input if (i) $t$ is a constant, or (ii) $t$ is a null and $C$ has an incoming edge $\mathrm{var}(t) \xrightarrow{x} w$. A null $n$ with $\mathrm{var}(n) \in C$ has $C$-depth $k$ if there is a maximal chain of nulls $m_0 \xrightarrow{y_1} \ldots \xrightarrow{y_k} m_k$ in $\mathrm{chase}(\Sigma, \mathcal{D})$ with $m_k = n$ and $\mathrm{var}(m_i) \in C$ for $0 \le i \le k$. The $C$-depth of $n$ is undefined if the length of such chains is unbounded.*

The next result limits the number of bounded-$C$-depth nulls based on the number of $C$-inputs, thereby also clarifying the significance of homogeneous confluence. Note that this general insight does not restrict to saturating tgd sets.

**Lemma 22.** *Consider a database $\mathcal{D}$ and $C \in \mathrm{SCC}(\mathrm{LXG}(\Sigma))$. If $i$ is the number of $C$-inputs in $\mathrm{chase}(\Sigma, \mathcal{D})$, then, for any $d \ge 0$, the number of nulls at $C$-depth $\le d$ is at most doubly exponential in $d$ and polynomial in $i$. If $C$ is homogeneously confluent, then this number is at most exponential in $d$.*

**Definition 23.** *Let $\Sigma$ be saturating with $E_C$ denoting the set $E$ of Definition 13 for $C \in \mathrm{SCC}(\mathrm{LXG}(\Sigma))$. Let $C_1, \ldots, C_k$ be a list of all $C_i \in \mathrm{SCC}(\mathrm{LXG}(\Sigma))$ topologically ordered w.r.t. the induced order $\prec$.*

*We define $\mathrm{rank}(C_i)$ iteratively for $i = 1, \ldots, k$ as follows, where we assume $\max\{\} = 0$. First, let $r_{\mathrm{in}}^i := \max\{\mathrm{rank}(C_j) \mid C_j \prec C_i\}$ and let $r_{\mathrm{cxt}}^i := \max\{\mathrm{rank}(C) \mid C \in C_{\mathrm{cxt}}^i\}$ where $C_{\mathrm{cxt}}^i$ denotes the set*

$$\{C_j \mid v, w \in C_i, u \in C_j, w \xrightarrow{x} v \in E_{C_i}, u \xrightarrow{y} v \in \mathrm{LXG}(\Sigma), x \ne y\}.$$

Then the rank of $C_i$ is

$$\operatorname{rank}(C_i) := \begin{cases} r_{\mathsf{in}}^i & \text{if } \operatorname{conf}(C_i) = 0, \\ \max\{r_{\mathsf{in}}^i, r_{\mathsf{cxt}}^i + 1\} & \text{if } \operatorname{conf}(C_i) = 1, \\ \max\{r_{\mathsf{in}}^i, r_{\mathsf{cxt}}^i + 2\} & \text{if } \operatorname{conf}(C_i) \geq 2. \end{cases}$$

The rank of $\Sigma$ is $\operatorname{rank}(\Sigma) = \max\{\operatorname{rank}(C) \mid C \in \operatorname{SCC}(\operatorname{LXG}(\Sigma))\}$.

**Theorem 24.** *Let $\Sigma$ be saturating and $\mathcal{D}$ a database. For every existential variable $v$ in $\Sigma$, the number of nulls $n$ with $\operatorname{var}(n) = v$ in $\operatorname{chase}(\Sigma, \mathcal{D})$ is at most $\operatorname{rank}(\operatorname{SCC}(v))$-exponential in the size of $\mathcal{D}$.*

Example 14 (and the discussion in Example 3) shows that the upper bound of Theorem 24 can be reached. We can further strengthen this into a hardness result:

**Theorem 25.** *Let $\Sigma$ be saturating. For every database $\mathcal{D}$ the size of $\operatorname{chase}(\Sigma, \mathcal{D})$ is at most $\operatorname{rank}(\Sigma)$-exponential in the size of $\mathcal{D}$, and BCQ entailment is $\operatorname{rank}(\Sigma)$-ExpTime-complete for data complexity.*

## 6 THE CHASE IN THE FOREST

In this section, we refine Theorem 25 by identifying cases where BCQ answering over saturating tgd sets $\Sigma$ is not $\operatorname{rank}(\Sigma)$-ExpTime-complete but merely $(\operatorname{rank}(\Sigma) - 1)$-ExpSpace-complete, and we design a chase procedure that runs within these complexity bounds. To simplify presentation, we consider tgd sets with a single rank-maximal SCC $\hat{C}$ in $\operatorname{LXG}(\Sigma)$ (generalisations are possible; see concluding remarks). If $\operatorname{conf}(\hat{C}) = 1$, we can establish a tree-like search space within the chase that follows the edges of $\operatorname{LXG}(\Sigma)$. A new syntactic restriction on tgds, called *path guardedness*, ensures that a chase that follows this tree-like structure remains complete for conjunctive query answering.

**Definition 26.** *A tgd set $\Sigma$ is* arboreous *if it is saturating, and has a unique $\hat{C} \in \operatorname{SCC}(\Sigma)$ with $\operatorname{rank}(\hat{C}) = \operatorname{rank}(\Sigma)$, which satisfies $\operatorname{conf}(\hat{C}) \leq 1$. For such $\Sigma$ and some $\operatorname{chase}(\Sigma, \mathcal{D})$, the* null forest *is the directed graph $\langle \hat{N}, \hat{\twoheadrightarrow} \rangle$ with $\hat{N} = \{n \in \operatorname{N}(\operatorname{chase}(\Sigma, \mathcal{D})) \mid \operatorname{var}(n) \in \hat{C}\}$ the nulls of variables in $\hat{C}$, and $n \hat{\twoheadrightarrow} m$ if $n \overset{x}{\twoheadrightarrow} m$ for some $x$.*

**Lemma 27.** *For every arboreous $\Sigma$ and $\operatorname{chase}(\Sigma, \mathcal{D})$, the null forest is indeed a forest (set of trees).*

Next, we use the special tgds that correspond to set $E$ of Definition 13 to partition the null forest into sub-forests. The intuition is that special tgd applications start a new sub-forest (their fresh nulls being the roots), whereas other tgd applications remain within their current sub-forest. By placing all remaining terms of the chase in an additional root node, we obtain a tree structure whose nodes are sets of terms that partition the terms of the chase:

**Definition 28.** *Let $\Sigma$, $\hat{C}$, $\hat{N}$, and $\hat{\twoheadrightarrow}$ be as in Definition 26, let $\hat{E}$ be the edge set $E$ of Definition 13 for $\hat{C}$, and let $T$ be the set of all terms in $\operatorname{chase}(\Sigma, \mathcal{D})$. An $\hat{E}$-tgd is a tgd with an existential variable that is the target of an edge in $\hat{E}$. The $\hat{E}$-variables $V_{\hat{E}}$ are all existential variables in $\hat{C}$ that occur in some $\hat{E}$-tgd.*

*For every chase step $i$ where an $\hat{E}$-tgd was applied, let $R[i]$ be the set of fresh nulls introduced in step $i$. Let $\bar{R} \subseteq \hat{N}$ be the set of nulls that are not in any such $R[i]$. Then, for each $R[i]$, let $F[i] \subseteq \hat{N}$ be the least set that contains $R[i]$ and all $m \in \bar{R}$ for which there is $n \in F[i]$ with $n \hat{\twoheadrightarrow} m$. Let $\mathcal{F}$ be the set of all such sets $F[i]$. For $F_1, F_2 \in \mathcal{F}$ with $F_1 \neq F_2$, we write $F_1 \tilde{\twoheadrightarrow} F_2$ if there are $n_i \in F_i$ $(i = 1, 2)$ such that $n_1 \hat{\twoheadrightarrow} n_2$.*

*Let $\mathsf{L}_0 = T \setminus \bigcup_{F \in \mathcal{F}} F$. The* term tree *of $\operatorname{chase}(\Sigma, \mathcal{D})$ is the graph $\langle N_\sim, \tilde{\twoheadrightarrow} \rangle$ with $N_\sim = \mathcal{F} \cup \{\mathsf{L}_0\}$, and $\tilde{\twoheadrightarrow}$ extended to $N_\sim$ by setting $\mathsf{L}_0 \tilde{\twoheadrightarrow} F$ for all $F \in \mathcal{F}$ that have no $\tilde{\twoheadrightarrow}$-predecessor in $\mathcal{F}$. The reflexive transitive closure of $\tilde{\twoheadrightarrow}$ is denoted $\tilde{\twoheadrightarrow}^*$.*
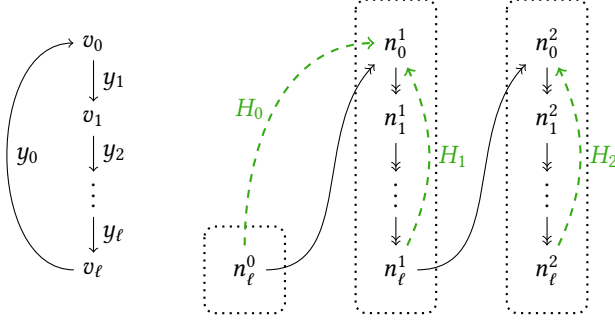
Fig. 2. Example graph LXG($\Sigma$) (left) and chain of derived nulls in chase($\Sigma, \mathcal{D}$) (right) from Figure 1, with nulls partitioned according to Definition 28 for the set $\hat{E} = \{v_\ell \xrightarrow{y_0} v_0\}$

**Lemma 29.** *For every arboreous $\Sigma$ and* chase($\Sigma, \mathcal{D}$)*, $N_\sim$ is a partition of the terms in* chase($\Sigma, \mathcal{D}$)*, and the term tree is indeed a tree.*

For every term $t$, Lemma 29 allows us to use $\mathsf{L}(t)$ to denote the unique set $\mathsf{L} \in N_\sim$ with $t \in \mathsf{L}$.

**Example 30.** *Figure 2 revisits the abstract example from Figure 1, which we assume to be saturating according to Definition 13 using $E = \{v_\ell \xrightarrow{y_0} v_0\}$. If this is the unique maximal SCC $\hat{C}$, then $\hat{E} = E$, and we obtain three partitions of nulls that are illustrated in Figure 2: $\mathsf{L}(n_\ell^0)$, $\mathsf{L}(n_0^1) = \mathsf{L}(n_1^1) = \cdots = \mathsf{L}(n_\ell^1)$, and $\mathsf{L}(n_0^2) = \mathsf{L}(n_1^2) = \cdots = \mathsf{L}(n_\ell^2)$. These partitions are part of a path in the term tree, which is overlaid on the original null forest.*

*The motivation for defining such a coarser tree structure is that we intend to use this tree to guide the computation of facts during the chase. We will limit its space-complexity by storing, at each particular moment during the chase, only those facts that can be represented using terms on a single path of this tree. The coarser the tree structure, the more facts can be considered at any moment, the more cases can be handled by this limited form of chase.*

*Besides this general intuition, the factorisation also plays a crucial role in deriving a syntactic criterion to recognise cases where such a tree-based chase can safely be applied, which we will consider next. The difficulty for this endeavour is that any such syntactic condition eventually has to rely on the facts that have induced the tree-like structure in the first place, such as $H_0$ in Figure 2. But these very facts also occur in backwards direction to ensure saturation, as indicated by $H_1$ and $H_2$ in the figure. In the coarser tree structure, such "backward edges" merely lead to the same node of the tree, rather than to a predecessor node from which we could enter parallel branches in forward direction.*

The tree structure of terms as such does not constrain the structure of inferred facts in chase($\Sigma, \mathcal{D}$), which may relate nulls from arbitrary positions in the null forest. We seek syntactic restrictions that ensure that the chase respects the term tree in the sense that the terms of any fact are on a common path and impose an order on terms that matches their position in the path. To this end, we derive relationships $\langle p, i \rangle \preceq \langle p, j \rangle$ on predicate positions such that, for all $p(t) \in$ chase($\Sigma, \mathcal{D}$) with $t_i, t_j \in \hat{N}$, we have that $\mathsf{L}(t_i)$ is an ancestor of (or possibly equal to) $\mathsf{L}(t_j)$ in the term tree. The next definition once again uses our assumption that distinct tgds do not share variables.

**Definition 31.** *Let $\Sigma$ be arboreous with $\hat{C}$, $\hat{E}$, and $V_{\hat{E}}$ as in Definition 28. We will define a (not necessarily transitive) binary relation $\preceq$ on predicate positions $\langle p, i \rangle$. Any such $\preceq$ induces a relation $\trianglelefteq$ on variables of $\Sigma$ as the reflexive, transitive closure of the set of all $x \trianglelefteq y$ such that $x$ and $y$ occur at positions $\langle p, i \rangle$ and $\langle p, j \rangle$ in a single body atom of some tgd in $\Sigma$, and $\langle p, i \rangle \preceq \langle p, j \rangle$.*

---

**Algorithm 1:** Non-deterministic chase for BCQ entailment

---

**In:** tgd set $\Sigma$, database $\mathcal{D}$, BCQ $q$, variable set $V_{\hat{E}}$, integer $M$
**Out:** $\Sigma \cup \mathcal{D} \models q$

1  $\mathcal{I} := \mathcal{D}$
2  $\mathcal{T} := \langle C_0 \rangle$ with $C_0$ the set of all constants in $\Sigma$ and $\mathcal{D}$
3  **for** $i \in \{1, \ldots, |q|\}$ **do**
4      **for** $j \in \{1, \ldots, M\}$ **do**
5          **choose** $\rho \in \Sigma$ with frontier $\boldsymbol{y}$, and match $\sigma$ applicable to $\mathcal{I}$ in Datalog-first chase
           **or break** (go to (L3))
6          **while** $|\mathcal{T}| > 1$ **and** $\boldsymbol{y}\sigma \cap \mathcal{T}.\mathsf{last}() = \emptyset$ **do**
7              $\mathcal{T}.\mathsf{pop}()$
8          **if** $v \cap V_{\hat{E}} \neq \emptyset$ **then**
9              $\mathcal{T}.\mathsf{push}(v\sigma^+)$
10         **else**
11             $\mathcal{T}.\mathsf{push}(\mathcal{T}.\mathsf{pop}() \cup v\sigma^+)$
12         $\mathcal{I} := \{p(\boldsymbol{t}) \in (\mathcal{I} \cup \mathrm{head}(\rho)\sigma^+) \mid \boldsymbol{t} \subseteq \bigcup_{T \in \mathcal{T}} T\}$
13     $\mathcal{T} := \langle \bigcup_{T \in \mathcal{T}} T \rangle$
14 **return** $\mathcal{I} \models q$

---

Now, concretely, let $\leq$ denote the largest relation on predicate positions such that, for all head atoms $p(\boldsymbol{t})$ in tgds of $\Sigma$ with universally quantified variables $\boldsymbol{y}$ and existentially quantified variables $\boldsymbol{v}$:

(1) if $t_i \in \boldsymbol{v} \cap \hat{C}$ and $t_j \in \boldsymbol{v} \setminus \hat{C}$ then $\langle p, i \rangle \nleq \langle p, j \rangle$,
(2) if $t_i \in \boldsymbol{v} \cap V_{\hat{E}}$ and $t_j \in \boldsymbol{y}$ then $\langle p, i \rangle \nleq \langle p, j \rangle$,
(3) if $t_i \in \boldsymbol{v} \cap \hat{C}$ and $t_j \in \boldsymbol{y}$ with $t_j \notin \lambda(t_i)$ then $\langle p, i \rangle \nleq \langle p, j \rangle$,
(4) if $t_i, t_j \in \boldsymbol{y}$ and $t_i \ntrianglelefteq t_j$ then $\langle p, i \rangle \nleq \langle p, j \rangle$.

One can construct $\leq$ in polynomial time with a simple greatest fixed point computation. Note that such a construction is anti-monotonic in $\Sigma$: more tgds lead to fewer constraints $\leq$.

**Lemma 32.** *Let $\Sigma$ be arboreous. For every $\mathcal{D}$, if $p(t_1, \ldots, t_n) \in \mathrm{chase}(\Sigma, \mathcal{D})$ with $\langle p, i \rangle \leq \langle p, j \rangle$, then $\mathsf{L}(t_i) \mathbin{\tilde{\twoheadrightarrow}}^* \mathsf{L}(t_j)$.*

**Definition 33.** *Let $\Sigma$ be arboreous with $\hat{C}$ and $\leq$ as in Definition 31. The set of $\hat{C}$-affected positions $\hat{\Omega}$ in $\Sigma$ is $\bigcup_{v \in \hat{C}} \Omega_v$, with $\Omega_v$ as in Definition 2. A body variable $x$ is $\hat{C}$-affected if $\mathrm{Pos}_x^B \subseteq \hat{\Omega}$.*

*A tgd $\rho \in \Sigma$ is path-guarded if all $\hat{C}$-affected body variables in $\rho$ are mutually comparable with respect to the relation $\leq$, i.e., form a chain in $\leq$. $\Sigma$ is path-guarded if all of its tgds are.*

Any node $\mathsf{L}$ of the term tree induces a unique upwards path $\mathrm{path}(\mathsf{L})$ that consists of all nodes $\mathsf{L}'$ with $\mathsf{L}' \mathbin{\tilde{\twoheadrightarrow}}^* \mathsf{L}$. For term $t$, we write $\mathrm{path}(t)$ for $\mathrm{path}(\mathsf{L}(t))$. Inferences of path-guarded tgds are situated on such paths:

**Lemma 34.** *Let $\Sigma$ be arboreous and path-guarded, and $\mathcal{D}$ some database. For every step $i$ in $\mathrm{chase}(\Sigma, \mathcal{D})$, with $\mathrm{tgd}[i] = B[\boldsymbol{x}, \boldsymbol{y}] \to \exists \boldsymbol{v}.H[\boldsymbol{y}, \boldsymbol{v}]$, and $T_B = (\boldsymbol{x} \cup \boldsymbol{y})\sigma[i]$ and $T_H = (\boldsymbol{y} \cup \boldsymbol{v})\sigma^+[i]$:*

(1) *if $T_B \neq \emptyset$ then $T_B \subseteq \bigcup \mathrm{path}(t)$ for some $t \in T_B$, and*
(2) *if $T_H \neq \emptyset$ then $T_H \subseteq \bigcup \mathrm{path}(t)$ for some $t \in T_H$.*

Algorithm 1 specifies a non-deterministic chase procedure to check the entailment of a BCQ. It is intended for arboreous, path-guarded tgd sets $\Sigma$ with variables $V_{\hat{E}}$ as in Definition 28. The input

$M$ bounds the length of the search: we will determine a rank$(\Sigma)$-exponential value for $M$ for which the algorithm decides query entailment in $(\text{rank}(\Sigma) - 1)$-NExpSpace, showing the problem to be in $(\text{rank}(\Sigma) - 1)$-ExpSpace by Savitch's Theorem.

Algorithm 1 maintains a set of inferences $\mathcal{I}$ over terms in $\mathcal{T}$, which is a list of sets of terms that corresponds to a path in the term tree. We use operations push, pop, and last, respectively, to add, remove, or read $\mathcal{T}$'s last element. The algorithm performs a search for each atom in $q$ (L3), adding one current path to $\mathcal{T}$'s root node after each run (L13). The inner loop (L4) non-deterministically chooses (L5) a tgd to apply to $\mathcal{I}$, or to break the iteration early (even if tgds are applicable). The current path $\mathcal{T}$ is pruned so that its last element contains a frontier term of the tgd (L6), before we either add a new term set (L9) or augment the last term set (L11), depending on whether $\rho$ is an $\hat{E}$-tgd. Finally, we add the inferred head and restrict $\mathcal{I}$ to atoms with terms in $\mathcal{T}$ (L12), where $\sigma^+$ extends $\sigma$ to existentially quantified variables using globally fresh nulls (not used in the algorithm before). Finally, we check if $\mathcal{I}$ entails $q$ (L14).

Given a run of Algorithm 1, we write $\mathcal{I}_i^j$ (resp. $\mathcal{T}_i^j$) to denote the value of $\mathcal{I}$ (resp. $\mathcal{T}$) after executing (L12) in the $i$th iteration of loop (L3) and the $j$th iteration of (L4). Although Algorithm 1 can forget inferences and repeat the same tgd application with distinct fresh nulls, its computations are correct in the following sense:

**Lemma 35.** *Let $\mathcal{I}^*$ be the union of all sets $\mathcal{I}_i^j$ of some run of Algorithm 1. There is a homomorphism $\tau : \mathcal{I}^* \to \text{chase}(\Sigma, \mathcal{D})$, and therefore $\text{chase}(\Sigma, \mathcal{D}) \models q$ whenever Algorithm 1 returns* true.

**Lemma 36.** *If $\Sigma$ is arboreous and path-guarded with rank$(\Sigma) > 0$, and $M \leq f(|\mathcal{D}|)$ for some rank$(\Sigma)$-exponential function $f$, then there is a $(\text{rank}(\Sigma) - 1)$-exponential function $g$ such that Algorithm 1 runs in space $g(|\mathcal{D}|)$, where 0-exponential means polynomial.*

It remains to show that, whenever $\Sigma, \mathcal{D} \models q$, Algorithm 1 admits a run that returns true and is bounded by an $M$ as in Lemma 36. Any run corresponds to a sequence of choices for (L5), which consists of $|q|$ sequences $\mathbf{s}_j$ of tgd applications. Let $\rho$ and $\sigma$ define the tgd application used to compute $\mathcal{I}_{i+1}^j$ from $\mathcal{I}_i^j$. We say that this tgd application *corresponds to chase step $s$* if $\text{tgd}[s] = \rho$, the restriction of $\tau$ to terms in $\mathcal{T}_i^j$ is injective, and $\sigma(z) = \tau^-(\sigma[s](z))$. In this case, we *canonically extend* $\tau$ to the fresh nulls by $\tau(v\sigma^+) := v\sigma^+[s]$ for all existential variables $v$ in $\rho$. This makes $\tau$ locally injective:

**Lemma 37.** *If all tgd applications in a run of Algorithm 1 correspond to chase steps, and $\tau$ in each iteration is the canonical extension for the respective step, then $\tau$ is a homomorphism $\mathcal{I}^* \to \text{chase}(\Sigma, \mathcal{D})$ that is injective on all term sets $\bigcup \mathcal{T}_i^j$ that occur during the run.*

For completeness of the algorithm, we are interested in runs where all tgd applications correspond to chase steps. Such runs can be specified as a list of $|q|$ sequences of chase steps, where repetitions of steps are allowed (and sometimes necessary).

We obtain a suitable choice sequence by scheduling *tasks* $\langle d, i \rangle$, to be read as "perform chase step $i$ under the assumption that $\mathcal{I}$ already contains (isomorphic copies of) all atoms that can be expressed using terms from the first $d - 1$ elements of the path of Lemma 34 (1)." Such tasks may require other tasks to be completed first, since $\mathcal{I}$ may not yet contain the whole premise of $\text{tgd}[i]$:

**Definition 38.** *For chase step $i$, let $\text{path}_B(i)$ denote the path of Lemma 34 (1). For an atom $\alpha \in \text{chase}(\Sigma, \mathcal{D})$, let $\text{path}(\alpha)$ be the smallest path in the term tree that contains all terms of $\alpha$ (which are on a path by induction over Lemma 34 (2)).*

*The subtasks of a task $\langle d, i \rangle$ are all tasks $\langle e, j \rangle$ with $e \geq d$ a depth, and $j < i$ the largest chase step that produced an atom $\alpha \in \mathcal{D}^{j+1} \setminus \mathcal{D}^j$ with $|\text{path}(\alpha)| = e$ and $\text{path}(\alpha) \subseteq \text{path}_B(i)$. The task tree*

*for $\langle d, i \rangle$ has a root node with label $\langle d, i \rangle$ and the task trees for all subtasks of $\langle d, i \rangle$ as its children. Children of a single parent are ordered by the depth in their label: $\langle e, j \rangle < \langle e', j' \rangle$ if $e < e'$.*

The atom $\alpha$ in Definition 38 ensures that the application of tgd$[j]$ in Algorithm 1 will not delete any previous inferences up to depth $e$ (through (L6) and (L12)). The order of subtasks ensures that inferences at smaller depths are computed first. Now the required sequence to successfully perform chase step $i$ in the inner loop of Algorithm 1 is obtained by traversing the task tree for $\langle 1, i \rangle$ in a topological, order-respecting way (children before parents, smaller sibling nodes before larger ones), extracting the sequence of chase steps from the second component of the sequence of tasks.

**Lemma 39.** *If $\Sigma$ is arboreous and path-guarded, $\boldsymbol{s}$ is the sequence of chase steps obtained from the task tree with root $\langle 1, i \rangle$, then the length $|\boldsymbol{s}|$ of $\boldsymbol{s}$ is bounded by a $\mathrm{rank}(\Sigma)$-exponential function.*

**Lemma 40.** *If $\Sigma$ is arboreous and path-guarded, $\boldsymbol{s}$ is the sequence of chase steps obtained from the task tree for $\langle 1, i \rangle$, and $M \geq |\boldsymbol{s}|$, then Algorithm 1 can choose tgd applications according to $\boldsymbol{s}$.*

Combining these results, we obtain the completeness of our tree-based chase. Indeed, whenever $q\theta \subseteq \mathrm{chase}(\Sigma, \mathcal{D})$ for some match $\theta$, there are chase steps $s_1, \ldots, s_{|q|}$ that produce the atoms of $q\theta$. A suitable strategy then executes Algorithm 1 for the choice sequences $\boldsymbol{s}_i$ obtained from the task trees for $\langle 1, s_i \rangle$ with $i = 1, \ldots, |q|$. Lemma 39 ensures that we can use a suitably small bound $M$ to use the complexity results of Lemma 36, whereas Lemma 40 ensures that the final atom set $\mathcal{I}_{|s_{|q|}|}^{|q|}$ contains all atoms of $\tau^-(q\theta)$.

**Theorem 41.** *BCQ entailment for arboreous and path-guarded tgd sets $\Sigma$ with $\mathrm{rank}(\Sigma) = \kappa \geq 1$ is $(\kappa - 1)$-ExpSpace-complete for data complexity, where 0-ExpSpace = PSpace.*

Hardness is shown by reduction from the word problem of $\kappa$-exponentially time-bounded alternating Turing machines. PSpace-hardness is illustrated with a simpler reduction from TrueQBF:

**Example 42.** *Let $\phi = \mho_1 p_1, \ldots, \mho_\ell p_\ell. \psi$ be a quantified Boolean formula with $\mho_i \in \{\exists, \forall\}$ for $1 \leq i \leq \ell$ and $\psi = \bigwedge_{j=1}^{k} C^j$ in CNF. Let $\tilde{p}_i$ denote $p_i$ or its negation $\bar{p}_i$; the clauses $C^i$ are sets of such literals.*

*The facts $\mathcal{D}_0 = \{empty(\emptyset), new(\vdash, \emptyset), nxt(\vdash, p_1), nxt(\vdash, \bar{p}_1)\} \cup \{nxt(\tilde{p}_i, \tilde{p}_{i+1}) \mid 1 \leq i < \ell\}$ encode an order over the literals. The tgds $\Sigma_0$ construct literal sets (truth assignments) similar to Example 18:*

$$getSU(x, S) \rightarrow \exists v.\, su(x, S, v) \land su(x, v, v) \tag{18}$$

$$su(x, S, T) \land su(y, S, S) \rightarrow su(y, T, T) \tag{19}$$

$$new(x, S) \land nxt(x, y) \rightarrow getSU(y, S) \tag{20}$$

$$new(x, S) \land nxt(x, y) \land su(y, S, T) \rightarrow new(y, T) \tag{21}$$

*Facts $new(l, S)$ mark the latest literal $l$ added to a set $S$. The tgd (20) triggers the creation of a suitable set, and tgd (21) marks the newly added literal. Note that we never have $su(p_i, s, s)$ and $su(\bar{p}_i, s, s)$.*

*Now $\mathcal{D}_1 = \{in(\tilde{p}_i, C) \mid \tilde{p}_i \in C^j, 1 \leq j \leq k\} \cup \{next(C^{j-1}, C^j) \mid 1 < j \leq k\} \cup \{first(C^1), last(C^k)\}$ encodes the clauses. We use tgds $\Sigma_1$ to evaluate $\psi$ under a given truth assignment:*

$$su(x, S, S) \land in(x, c) \land first(c) \rightarrow sat(c, S) \tag{22}$$

$$sat(c, S) \land next(c, d) \land su(x, S, S) \land in(x, d) \rightarrow sat(d, S) \tag{23}$$

$$sat(c, S) \land last(c) \rightarrow sat(S) \tag{24}$$

*Here, $sat(c, s)$ means "set $s$ satisfies clause $c$". The tgds $\Sigma_1$ then propagate satisfaction along the order of $\mathcal{D}_1$. Therefore, $\mathrm{chase}(\Sigma_0 \cup \Sigma_1, \mathcal{D}_0 \cup \mathcal{D}_1) \models \exists v.\, sat(v)$ if and only if $\psi$ is satisfiable.*

*Finally, we use $\mathcal{D}_2 = \{ex(p_i) \mid \mathsf{Q}_i = \exists, 1 \le i \le \ell\} \cup \{pos(p_i), neg(\bar{p}_i) \mid 1 \le i \le \ell\}$ and the following tgds $\Sigma_2$ to evaluate $\phi$:*

$$su(x, S, T) \wedge ex(x) \wedge sat(T) \rightarrow sat(S) \tag{25}$$

$$su(x, S, T) \wedge pos(x) \wedge sat(T) \rightarrow sat_+(S) \tag{26}$$

$$su(x, S, T) \wedge neg(x) \wedge sat(T) \rightarrow sat_-(S) \tag{27}$$

$$sat_+(S) \wedge sat_-(S) \rightarrow sat(S) \tag{28}$$

*We check satisfaction by evaluating the tree of literal sets by propagating satisfaction from leafs towards the root. The tgd (25) handles existential quantification, and tgds (26)–(28) handle universal quantification. Handling the successors of universal states separately ensures that the tgds are path-guarded. Indeed, $\le$ of Definition 31 for the used tgds contains $\langle su, 2\rangle \le \langle su, 3\rangle$. Note that $\mathsf{chase}(\Sigma_0 \cup \Sigma_1 \cup \Sigma_2, \mathcal{D}_0 \cup \mathcal{D}_1 \cup \mathcal{D}_2) \models \exists v.\, empty(v) \wedge sat(v)$ if and only if $\phi$ is true.*

## 7 CONCLUSIONS AND OUTLOOK

We have established new criteria for chase termination, which advance the state of the art in two important ways: (1) they can take advantage of the standard chase, and (2) they yield new decidable tgd classes with data complexities that are complete for $\kappa$-ExpTime and $\kappa$-ExpSpace, for any $\kappa \ge 0$. This is obviously too high for transactional DBMS loads, but it allows us to address a much larger range of complex computational tasks over databases with the chase. Practical problems of this kind include ontology reasoning [15], database provenance computation [23], and querying databases with complex values [35]. We also note that checking our criteria is always dominated by Datalog reasoning, which is practically feasible and of lower complexity than some established criteria [18].

Our work brings up many follow-up questions. First, the new tgd classes are candidates for capturing their respective complexity classes (at least from PSpace upwards), but known proof techniques rely on non-saturating tgds [8]. Second, our techniques require a Datalog-first chase strategy, which is avoidable for sets and complex values [35]. It is open if similar approaches could apply in our setting. Third, our criteria can be broadened, e.g., the restriction to a single maximal-rank component in Section 6 can be relaxed.

Taking a wider view, a central methodological contribution of our work is the labelled dependency graph and its extensive use for analysing the internal structure of the standard chase. It can be seen as a surrogate for the more syntactic "lineage" of nulls that is available in the semi-oblivious chase – best exposed through the use of skolem terms [34] –, which has been extremely useful in studying that chase variant. It is exciting to ask how our method can be similarly useful in further studying the standard chase, e.g., to detect non-termination or to decide termination in new cases, and whether it can be refined in the style of termination checks based on materialisation or control flow analysis.

## REFERENCES
[1] Alfred V. Aho, Catriel Beeri, and Jeffrey D. Ullman. 1979. The Theory of Joins in Relational Databases. *ACM Trans. Database Syst.* 4, 3 (1979), 297–314.

[2]  Alfred V. Aho, Yehoshua Sagiv, and Jeffrey D. Ullman. 1979. Efficient Optimization of a Class of Relational Expressions. *ACM Trans. Database Syst.* 4, 4 (1979), 435–454. https://doi.org/10.1145/320107.320112

[3]  Jean-François Baget, Michel Leclère, Marie-Laure Mugnier, Swan Rocher, and Clément Sipieter. 2015. Graal: A Toolkit for Query Answering with Existential Rules. In *Proc. 9th Int. Web Rule Symposium (RuleML'15) (LNCS, Vol. 9202)*, Nick Bassiliades, Georg Gottlob, Fariba Sadri, Adrian Paschke, and Dumitru Roman (Eds.). Springer, 328–344.

[4]  Jean-François Baget, Michel Leclère, Marie-Laure Mugnier, and Eric Salvat. 2011. On rules with existential variables: Walking the decidability line. *Artificial Intelligence* 175, 9–10 (2011), 1620–1654.

[5]  Catriel Beeri and Moshe Y. Vardi. 1984. A Proof Procedure for Data Dependencies. *J. ACM* 31, 4 (1984), 718–741.

[6]  Michael Benedikt, George Konstantinidis, Giansalvatore Mecca, Boris Motik, Paolo Papotti, Donatello Santoro, and Efthymia Tsamoura. 2017. Benchmarking the Chase. In *Proc. 36th Symposium on Principles of Database Systems (PODS'17)*, Emanuel Sallinger, Jan Van den Bussche, and Floris Geerts (Eds.). ACM, 37–52.

[7]  Michael Benedikt, Julien Leblay, and Efthymia Tsamoura. 2014. PDQ: Proof-driven Query Answering over Web-based Data. *Proc. VLDB Endowment* 7, 13 (2014), 1553–1556.

[8]  Camille Bourgaux, David Carral, Markus Krötzsch, Sebastian Rudolph, and Michaël Thomazo. 2021. Capturing Homomorphism-Closed Decidable Queries with Existential Rules. In *Proc. 18th Int. Conf. on Principles of Knowledge Representation and Reasoning (KR'21)*, Meghyn Bienvenu, Gerhard Lakemeyer, and Esra Erdem (Eds.). 141–150.

[9]  Marco Calautti, Georg Gottlob, and Andreas Pieris. 2022. Non-Uniformly Terminating Chase: Size and Complexity. In *Proc. 41st Symposium on Principles of Database Systems (PODS'22)*, Leonid Libkin and Pablo Barceló (Eds.). ACM, 369–378.

[10]  Marco Calautti, Mostafa Milani, and Andreas Pieris. 2023. Semi-Oblivious Chase Termination for Linear Existential Rules: An Experimental Study. *Proc. VLDB Endow.* 16, 11 (2023), 2858–2870. https://doi.org/10.14778/3611479.3611493

[11]  Andrea Calì, Georg Gottlob, and Andreas Pieris. 2010. Query Answering under Non-guarded Rules in Datalog+/-. In *Proc. 4th Int. Conf. on Web Reasoning and Rule Systems (RR 2010) (LNCS, Vol. 6333)*, Pascal Hitzler and Thomas Lukasiewicz (Eds.). Springer, 1–17.

[12]  Andrea Calì, Georg Gottlob, and Andreas Pieris. 2012. Towards more expressive ontology languages: The query answering problem. *J. of Artif. Intell.* 193 (2012), 87–128.

[13]  Andrea Calì, Domenico Lembo, and Riccardo Rosati. 2003. On the decidability and complexity of query answering over inconsistent and incomplete databases. In *Proc. 22nd Symposium on Principles of Database Systems (PODS'03)*, Frank Neven, Catriel Beeri, and Tova Milo (Eds.). ACM, 260–271. https://doi.org/10.1145/773153.773179

[14]  David Carral, Irina Dragoste, and Markus Krötzsch. 2017. Restricted Chase (Non)Termination for Existential Rules with Disjunctions. In *Proc. 26th Int. Joint Conf. on Artificial Intelligence (IJCAI'17)*, Carles Sierra (Ed.). ijcai.org, 922–928.

[15]  David Carral, Irina Dragoste, Markus Krötzsch, and Christian Lewe. 2019. Chasing Sets: How to Use Existential Rules for Expressive Reasoning. In *Proc. 28th Int. Joint Conf. on Artificial Intelligence (IJCAI'19)*, Sarit Kraus (Ed.). ijcai.org, 1624–1631.

[16]  David Carral, Lucas Larroque, Marie-Laure Mugnier, and Michaël Thomazo. 2022. Normalisations of Existential Rules: Not so Innocuous!. In *Proc. 19th Int. Conf. on Principles of Knowledge Representation and Reasoning (KR'22)*, Gabriele Kern-Isberner, Gerhard Lakemeyer, and Thomas Meyer (Eds.). 10 pages.

[17]  Ashok K. Chandra, Dexter C. Kozen, and Larry J. Stockmeyer. 1981. Alternation. *J. ACM* 28, 1 (1981), 114–133.

[18]  Bernardo Cuenca Grau, Ian Horrocks, Markus Krötzsch, Clemens Kupke, Despoina Magka, Boris Motik, and Zhe Wang. 2013. Acyclicity Notions for Existential Rules and Their Application to Query Answering in Ontologies. *J. of Artificial Intelligence Research* 47 (2013), 741–808.

[19]  Evgeny Dantsin, Thomas Eiter, Georg Gottlob, and Andrei Voronkov. 2001. Complexity and Expressive Power of Logic Programming. *Comput. Surveys* 33, 3 (2001), 374–425.

[20]  Anuj Dawar and Stephan Kreutzer. 2008. On Datalog vs. LFP. In *Proc. 35th Int. Colloquium on Automata, Languages, and Programming (ICALP'08); Part II (LNCS, Vol. 5126)*, Luca Aceto, Ivan Damgård, Leslie Ann Goldberg, Magnús M. Halldórsson, Anna Ingólfsdóttir, and Igor Walukiewicz (Eds.). Springer, 160–171.

[21]  Alin Deutsch, Alan Nash, and Jeffrey B. Remmel. 2008. The Chase Revisited. In *Proc. 27th Symposium on Principles of Database Systems (PODS'08)*, Maurizio Lenzerini and Domenico Lembo (Eds.). ACM, 149–158.

[22]  Alin Deutsch and Val Tannen. 2003. Reformulation of XML Queries and Constraints. In *Proc. 9th Int. Conf. on Database Theory (ICDT'03) (LNCS, Vol. 2572)*, Diego Calvanese, Maurizio Lenzerini, and Rajeev Motwani (Eds.). Springer, 225–241.

[23]  Ali Elhalawati, Markus Krötzsch, and Stephan Mennicke. 2022. An Existential Rule Framework for Computing Why-Provenance On-Demand for Datalog. In *Proc. 2nd Int. Joint Conf. on Rules and Reasoning (RuleML+RR'22) (LNCS, Vol. 13752)*, Guido Governatori and Anni-Yasmin Turhan (Eds.). Springer, 146–163.

[24]  Ronald Fagin, Phokion G. Kolaitis, Renée J. Miller, and Lucian Popa. 2005. Data exchange: semantics and query answering. *Theoretical Computer Science* 336, 1 (2005), 89–124.

[25] Sarah Alice Gaggl, Philipp Hanisch, and Markus Krötzsch. 2022. Simulating Sets in Answer Set Programming. In *Proc. 31st Int. Joint Conf. on Artificial Intelligence (IJCAI'22)*, Luc De Raedt (Ed.). ijcai.org, 2634–2640. https://doi.org/10.24963/ijcai.2022/365

[26] Floris Geerts, Giansalvatore Mecca, Paolo Papotti, and Donatello Santoro. 2014. That's All Folks! LLUNATIC Goes Open Source. *PVLDB* 7, 13 (2014), 1565–1568.

[27] Tomasz Gogacz and Jerzy Marcinkowski. 2014. All-Instances Termination of Chase is Undecidable. In *Proc. 41st Int. Colloquium on Automata, Languages, and Programming (ICALP'14); Part II (LNCS, Vol. 8573)*, Javier Esparza, Pierre Fraigniaud, Thore Husfeldt, and Elias Koutsoupias (Eds.). Springer, 293–304.

[28] Tomasz Gogacz, Jerzy Marcinkowski, and Andreas Pieris. 2023. Uniform Restricted Chase Termination. *SIAM J. Comput.* 52, 3 (2023), 641–683. https://doi.org/10.1137/20M1377035

[29] Gösta Grahne and Adrian Onet. 2018. Anatomy of the Chase. *Fundam. Inform.* 157, 3 (2018), 221–270.

[30] Philipp Hanisch and Markus Krötzsch. 2024. Chase Termination Beyond Polynomial Time. In *Proc. 43rd Symposium on Principles of Database Systems (PODS'24)*. ACM.

[31] Markus Krötzsch and Sebastian Rudolph. 2011. Extending Decidable Existential Rules by Joining Acyclicity and Guardedness. In *Proc. 22nd Int. Joint Conf. on Artificial Intelligence (IJCAI'11)*, Toby Walsh (Ed.). AAAI Press/IJCAI, 963–968.

[32] Markus Krötzsch, Maximilian Marx, and Sebastian Rudolph. 2019. The Power of the Terminating Chase. In *Proc. 22nd Int. Conf. on Database Theory (ICDT'19) (LIPIcs, Vol. 127)*, Pablo Barceló and Marco Calautti (Eds.). Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, 3:1–3:17.

[33] David Maier, Alberto O. Mendelzon, and Yehoshua Sagiv. 1979. Testing implications of data dependencies. *ACM Transactions on Database Systems* 4 (1979), 455–469. Issue 4.

[34] Bruno Marnette. 2009. Generalized Schema-Mappings: from Termination to Tractability. In *Proc. 28th Symposium on Principles of Database Systems (PODS'09)*, Jan Paredaens and Jianwen Su (Eds.). ACM, 13–22.

[35] Maximilian Marx and Markus Krötzsch. 2022. Tuple-Generating Dependencies Capture Complex Values. In *Proc. 25th Int. Conf. on Database Theory (ICDT'22) (LIPIcs, Vol. 220)*, Dan Olteanu and Nils Vortmeier (Eds.). Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 13:1–13:20. https://doi.org/10.4230/LIPIcs.ICDT.2022.13

[36] Michael Meier. 2014. The backchase revisited. *VLDB J.* 23, 3 (2014), 495–516. https://doi.org/10.1007/S00778-013-0333-Y

[37] Michaël Thomazo Michel Leclére, Marie-Laure Mugnier and Federico Ulliana. 2019. On Chase Termination for Linear Existential Rules. In *Proc. 22nd Int. Conf. on Database Theory (ICDT'19)*. Schloss Dagstuhl – Leibniz-Zentrum fuer Informatik.

[38] Yavor Nenov, Robert Piro, Boris Motik, Ian Horrocks, Zhe Wu, and Jay Banerjee. 2015. RDFox: A Highly-Scalable RDF Store. In *Proc. 14th Int. Semantic Web Conf. (ISWC'15), Part II (LNCS, Vol. 9367)*, Marcelo Arenas, Óscar Corcho, Elena Simperl, Markus Strohmaier, Mathieu d'Aquin, Kavitha Srinivas, Paul T. Groth, Michel Dumontier, Jeff Heflin, Krishnaprasad Thirunarayan, and Steffen Staab (Eds.). Springer, 3–20.

[39] Sebastian Rudolph and Michaël Thomazo. 2015. Characterization of the Expressivity of Existential Rule Queries. In *Proc. 24th Int. Joint Conf. on Artificial Intelligence (IJCAI'15)*, Qiang Yang and Michael Wooldridge (Eds.). AAAI Press, 3193–3199.

[40] Sebastian Rudolph and Michaël Thomazo. 2016. Expressivity of Datalog Variants - Completing the Picture. In *Proc. 25th Int. Joint Conf. on Artificial Intelligence (IJCAI'15)*, Subbarao Kambhampati (Ed.). AAAI Press, 1230–1236.

[41] Jacopo Urbani, Markus Krötzsch, Ceriel J. H. Jacobs, Irina Dragoste, and David Carral. 2018. Efficient Model Construction for Horn Logic with VLog: System description. In *Proc. 9th Int. Joint Conf. on Automated Reasoning (IJCAR'18) (LNCS, Vol. 10900)*, Didier Galmiche, Stephan Schulz, and Roberto Sebastiani (Eds.). Springer, 680–688.

[42] Heng Zhang, Yan Zhang, and Jia-Huai You. 2015. Existential Rule Languages with Finite Chase: Complexity and Expressiveness. In *Proc. 29th AAAI Conf. on Artificial Intelligence (AAAI'15)*, Blai Bonet and Sven Koenig (Eds.). AAAI Press.

## A  PROOFS FOR SECTION 3

**Lemma 5.** *If* $\text{chase}(\Sigma, \mathcal{D})$ *is infinite, then* $\text{chase}(\Sigma, \mathcal{D})$ *contains an infinite chain* $n_0 \overset{y_1}{\twoheadrightarrow} n_1 \overset{y_2}{\twoheadrightarrow} \cdots$.

PROOF. If $\text{chase}(\Sigma, \mathcal{D})$ is infinite, it must contain infinitely many nulls. Let $\prec$ be some total order on nulls of $\text{chase}(\Sigma, \mathcal{D})$ such that $n \prec m$ holds whenever $n$ was introduced at an earlier chase step than $m$. Let $G$ be the directed graph that has all such nulls as its vertices and that contains an edge $n \to m$ if $n$ is the $\prec$-largest null with $n \overset{y}{\twoheadrightarrow} m$ in $\text{chase}(\Sigma, \mathcal{D})$ for some $y$. Then $G$ is acyclic (since $n \overset{y}{\twoheadrightarrow} m$ implies $n \prec m$) and a forest (since each vertex has a unique predecessor by definition). Root vertices in $G$ correspond to nulls introduced by tgd applications to non-null frontier variables; since there are only finitely many such root nulls, $G$ (being infinite) contains an infinite tree $T$. Moreover, $T$ is finitely branching: indeed, $n \to m \in T$ implies that $m$ was produced by a tgd application to frontier terms that were introduced in $\text{chase}(\Sigma, \mathcal{D})$ no later than $n$; there are only finitely many such terms; hence there are only finitely many such tgd applications. As a finitely branching, infinite tree, $T$ must contain an infinite path by Kőnig's lemma, and this path corresponds to the required chain. □

## B  PROOFS FOR SECTION 4

**Lemma 9.** *Let* $\mathfrak{c} = n_0 \overset{y_1}{\twoheadrightarrow} n_1 \overset{y_2}{\twoheadrightarrow} \cdots \overset{y_k}{\twoheadrightarrow} n_k$ *be a chain in* $\text{chase}(\Sigma, \mathcal{D})$, *with* $n_k$ *derived in the* $m$-*th chase step* $\mathcal{D}^m$. *There is a path* $\mathfrak{p} = \text{var}(v_0) \overset{y_1}{\to} \text{var}(n_1) \overset{y_2}{\to} \cdots \overset{y_k}{\to} \text{var}(n_k)$ *in* $\text{LXG}(\Sigma)$, *and* $\mathcal{D}^m \models \text{Path}(\mathfrak{p})\theta$ *holds for substitution* $\theta$ *defined as follows: for each* $n_i$ *derived in chase step* $j$, *and each variable* $x$ *in* $\text{tgd}[j]$ *that was renamed to* $\tilde{x}$ *in* $\text{Path}(\mathfrak{c})$, $\theta(\tilde{x}) = \sigma^+[j](x)$.

PROOF. The existence of $\mathfrak{p}$ is guaranteed by Lemma 4. $\text{Path}(\mathfrak{p})\theta \subseteq \mathcal{D}^m$ follows from Definitions 1 and 7, and the definition of $\theta$. □

**Lemma 12.** *Checking whether a path* $\mathfrak{p}$ *is base-propagating (or step-propagating) is* EXPTIME-*complete, and* P-*complete with respect to the length of* $\mathfrak{p}$.

PROOF. To check the relevant entailments of the form $\Sigma_{\text{DL}} \models B \to H$ as in (7) and (8), we can check the entailment $\Sigma_{\text{DL}}, B' \models H'$, where $B'$ and $H'$ are sets of atoms obtained by uniformly replacing variables in $B$ and $H$ with fresh constants. The claimed complexities are that of Datalog entailment [19]. □

**Theorem 15.** *Deciding if* $C$ *is* $E$-*saturating is* EXPTIME-*complete in the size of* $C$. *The same complexity applies to deciding if a tgd set* $\Sigma$ *is saturating.*

PROOF. Hardness follows from Lemma 12. For inclusion, note that Definition 13 (1) can be checked in polynomial time for a given $E$. If it holds true, there are at most exponentially many $\bar{E}$-paths in $C$, leading to exponentially many checks for (3) and (4), which are each in EXPTIME by Lemma 12. The last part of the claim follows since there are a polynomial number of strongly connected components, each with at most exponentially many candidate sets for $E$. □

**Lemma 16.** *Let* $C$ *be an* $E$-*saturating strongly connected component in* $\text{LXG}(\Sigma)$, *and let* $u \overset{y}{\to} v \in E$ *have the corresponding tgd* $\rho_v$ *with head* $\exists v, \mathbf{w}.H[y, z, v, \mathbf{w}]$. *For every database* $\mathcal{D}$ *and chain of nulls* $n_0 \overset{y}{\twoheadrightarrow} n_1 \overset{x_2}{\twoheadrightarrow} \ldots \overset{x_{\ell-1}}{\twoheadrightarrow} n_{\ell-1} \overset{y}{\twoheadrightarrow} n_\ell$ *in* $\text{chase}(\Sigma, \mathcal{D})$ *such that* $n_0 = y\sigma[a]$, $n_1 = v\sigma^+[a]$, $n_{\ell-1} = y\sigma[b]$, *and* $n_\ell = v\sigma^+[b]$ *for chase steps* $a < b$, *we have* $z\sigma[a] \neq z\sigma[b]$.

PROOF. Let $e := u \overset{y}{\to} v$. By Lemma 4, the given chain $\mathfrak{c}$ corresponds to a path $\mathfrak{p}$ in $\text{LXG}(\Sigma)$. The first and last edge of $\mathfrak{p}$ is $e$, so all edges of $\mathfrak{p}$ are in $C$. We can view $\mathfrak{p}$ as a path of the form $e_1 \mathfrak{c}_1 e_2 \cdots e_{k-1} \mathfrak{c}_{k-1} e_k$ with $e_1 = e_k = e$; $e_2, \ldots, e_k \in E$; and $\mathfrak{c}_i$ an $\bar{E}$-path in the sense of Definition 13

for $1 \leq i < k$. For $1 \leq i \leq k$, let $n_s^i \overset{x(i)}{\twoheadrightarrow} n_t^i$ be the part of $\mathfrak{c}$ that corresponds to $e_i$, and let $s(i)$ be the respective chase step (in particular, $s(1) = a$ and $s(k) = b$).

We show by induction that, for all $2 \leq i \leq k$, we have $\mathcal{D}^{s(i)} \models H[n_s^i, z\sigma[a], n_t^{i-1}, w\sigma^+[a]]$. In particular, for $i = k$ we get $\mathcal{D}^b \models H[n_{\ell-1}, z\sigma[a], n_t^{k-1}, w\sigma^+[a]]$. This shows the claim, since it means that Definition 1 (2.a) would not be satisfied if $z\sigma[a] = z\sigma[b]$.

The base case $i = 2$ follows since $e_1 \mathfrak{c}_1$ is base-propagating by Definition 13 (3). For the induction step, suppose the claim holds for $i$. The claim follows for $i + 1$ since $e_i \mathfrak{c}_i e_{i+1} \mathfrak{c}_{i+1}$ is step-propagating for $e$ by Definition 13 (4). □

**Theorem 17.** *If $\Sigma$ is saturating, then* chase$(\Sigma, \mathcal{D})$ *is finite for all databases $\mathcal{D}$.*

PROOF. Suppose for a contradiction that chase$(\Sigma, \mathcal{D})$ is infinite. Then there are infinitely many nulls, and a tgd that is applied infinitely often. The *tgd of* an edge $w' \overset{x}{\to} w$ in LXG$(\Sigma)$ is $\rho_w$, the unique tgd with variable $w$. Now let $C$ be a strongly connected component of LXG$(\Sigma)$ such that (1) $C$ contains an edge whose tgd is applied infinitely often, and (2) the tgd of every edge $w' \overset{x}{\to} w$ with $w' \notin C$ and $w \in C$ is applied only finitely often.

Using a similar argument as for the proof of Lemma 5, we find that there is an infinite chain $n_0 \overset{y_1}{\twoheadrightarrow} n_1 \overset{y_2}{\twoheadrightarrow} \cdots$ in chase$(\Sigma, \mathcal{D})$ such that $\text{var}(n_i) \in C$ for all $i \geq 0$. By Lemma 4, this chain corresponds to an infinite path $\mathfrak{p}$ in $C$. By assumption, $C$ is $E$-saturating for some set $E$, so, by Definition 13 (1), some edge $u \overset{y}{\to} v \in E$ occurs infinitely often in $\mathfrak{p}$. Let $\rho_v$ be its tgd, and let head$(\rho_v) = \exists v, w.H[y, z, v, w]$.

By Definition 13 (2) and Lemma 4, applications of $\rho_v$ can only involve null values $n$ for variables of $z$ if (i) $\text{var}(n) \notin C$, and (ii) there is an edge from $\text{var}(n)$ to $C$ in LXG$(\Sigma)$. By our choice of $C$, the number of such nulls $n$ is finite, as is the number of constants in $\Sigma$ and $\mathcal{D}$, so there are only finitely many possible instantiations of $z$ in applications of $\rho_v$. Together with Lemma 16, this implies that $\rho_v$ is applied only finitely many times – a contradiction. □

## C  PROOFS FOR SECTION 5

**Lemma 22.** *Consider a database $\mathcal{D}$ and $C \in \text{SCC}(\text{LXG}(\Sigma))$. If $i$ is the number of $C$-inputs in* chase$(\Sigma, \mathcal{D})$*, then, for any $d \geq 0$, the number of nulls at $C$-depth $\leq d$ is at most doubly exponential in $d$ and polynomial in $i$. If $C$ is homogeneously confluent, then this number is at most exponential in $d$.*

PROOF. Let $\alpha = \max\{|y| : B[x, y] \to \exists v.H[y, v] \in \Sigma, v \cap C \neq \emptyset\}$, and let $\beta = \text{conf}(C)$. The claim is easy to see for trivial components (case $\beta = 0$), so we only consider the case $\beta \geq 1$. We define upper bounds $g(d)$ for the number of nulls of $C$-depth $\leq d$.

Nulls of $C$-depth $d + 1$ are created by applying $\rho_v$ with $v \in C$ to a frontier of at most $\alpha$ terms that are $C$-inputs and nulls of $C$-depth $\leq d$, where at most $\beta$ of the terms are nulls of $C$-depth $\leq d$. An upper bound for nulls in $C$ produced by such tgd applications therefore is $|C|i^\alpha g(d)^\beta$. The total number of nulls at $C$-depth $\leq d + 1$ is therefore bounded by $|C|i^\alpha g(d)^\beta + g(d) \leq (|C|i^\alpha + 1)g(d)^\beta$ (where we use $\beta \geq 1$). We can therefore define $g(0) := |C|i^\alpha + 1$ and $g(d + 1) := (|C|i^\alpha + 1)g(d)^\beta = g(0)g(d)^\beta$ to obtain an upper bound for the number of nulls of $C$-depth $\leq d$.

If $\beta = 1$, then $g(d) = g(0)^{d+1}$ is exponential in $d$, as claimed. If $\beta > 1$, we use a relaxed upper bound $f(d) := g(0)^{(\beta+1)^d}$. Indeed, $g(d) \leq f(d)$ follows by induction: the base case follows from $f(0) = g(0)$; and the step follows from

$$g(d + 1) = g(0)g(d)^\beta \leq g(d)g(d)^\beta = g(d)^{\beta+1}$$

$$\leq f(d)^{\beta+1} = \left(g(0)^{(\beta+1)^d}\right)^{\beta+1} = g(0)^{(\beta+1)^{d+1}} = f(d + 1)$$

where the second "$\leq$" uses the induction hypothesis. Hence, $f$ is the claimed doubly exponential bound.                                                                                                                                    □

**Theorem 24.** *Let $\Sigma$ be saturating and $\mathcal{D}$ a database. For every existential variable $v$ in $\Sigma$, the number of nulls $n$ with $\mathrm{var}(n) = v$ in $\mathrm{chase}(\Sigma, \mathcal{D})$ is at most $\mathrm{rank}(\mathrm{SCC}(v))$-exponential in the size of $\mathcal{D}$.*

PROOF. Let $C_0, \ldots, C_k$ be a topological order as in Definition 23, and let $\iota(v)$ denote the number of nulls $n$ with $\mathrm{var}(n) = v$ in $\mathrm{chase}(\Sigma, \mathcal{D})$. Moreover, let $\kappa$ denote the number of constants in $\mathcal{D}$ and $\Sigma$. We show the claim for all existential variables $v$ with $\mathrm{SCC}(v) = C_i$ by induction over $i = 0, \ldots, k$.

Consider $C_i$ and assume that the claim holds for all $C_j$ with $j < i$. Let $C_1^-, \ldots, C_\ell^-$ with $C_j^- \prec C_i$ be the direct $\prec$-predecessors of $C_i$. Hence, for all $1 \leq j \leq \ell$, there is $m < i$ with $C_j^- = C_m$, and, for every $v \in C_j^-$, $\iota(v)$ is at most $\mathrm{rank}(C_j^-)$-exponential by induction ($*$). Therefore, the number of $C_i$-inputs is at most $\mathrm{in}_i := \kappa + \sum_{j=1}^\ell \sum_{v \in C_j^-} \iota(v)$, which (by ($*$)) is $r_{\mathrm{in}}^i$-exponential for $r_{\mathrm{in}}^i$ as in Definition 23. Analogously, the number $\mathrm{cxt}_i$ of $C_i$-inputs that are nulls $n$ with a $C_i$-incoming edge $\mathrm{var}(n) \xrightarrow{y} v$ such that there is an edge $w \xrightarrow{x} v \in E_{C_i}$ with $x \neq y$ is $r_{\mathrm{cxt}}^i$-exponential.

Case $\mathrm{conf}(C_i) = 0$. Then $C_i = \{v\}$ and $\rho_v = B[x, y] \to \exists v.H[y, v]$ with $v \in v$. The number of instantiations of $y$ and nulls $n$ with $\mathrm{var}(n) = v$ is bounded by $(\mathrm{in}_i)^{|y|}$. This bound is $r_{\mathrm{in}}^i$-exponential and $\mathrm{rank}(C_i) = r_{\mathrm{in}}^i$, so the claim holds.

Case $\mathrm{conf}(C_i) \geq 1$. Let $E$ denote the set $E_{C_i}$. The $C_i$-depth of nulls has an upper bound that is polynomial in $\mathrm{cxt}_i$. Indeed, consider an arbitrary chain $\mathfrak{c}$ of nulls $m_j$ with $\mathrm{var}(m_j) \in C_i$ as in Definition 21. By Lemma 16, for every $u \xrightarrow{y} v \in E$ with head $\exists v, w.H[y, z, v, w]$ of the corresponding tgd $\rho_v$, all applications of $\rho_v$ in $\mathfrak{c}$ use a different instantiation of $z$. By Definition 13 (2), variables in $z$ can only be instantiated with values from the $\leq \kappa + \mathrm{cxt}_i$ many $C_i$-inputs for which $C_i$ has an incoming edge to $v$. Hence, there are at most $(\kappa + \mathrm{cxt}_i)^{|z|}$ applications of $\rho_v$ in $\mathfrak{c}$. Using $\alpha$ to denote the maximal number of frontier variables $z$ in any tgd of $E$, there are at most $|E| \cdot (\kappa + \mathrm{cxt}_i)^\alpha$ applications of $E$-tgds in $\mathfrak{c}$.

Now by Definition 13 (1), the number of consecutive tgd applications in $\mathfrak{c}$ that correspond to $\bar{E}$-edges is at most $|\bar{E}|$. Hence, the overall length of $\mathfrak{c}$ is bounded by $|\bar{E}| \cdot (|E| \cdot (\kappa + \mathrm{cxt}_i)^\alpha + 1)$. This bound is polynomial in $\mathrm{cxt}_i$, since $|\bar{E}|$, $|E|$ and $\alpha$ are fixed by $\Sigma$. Therefore, the $C_i$-depth of nulls is at most polynomial in $\mathrm{cxt}_i$, and hence bounded by an $r_{\mathrm{cxt}}^i$-exponential function.

The claim now follows from Lemma 22, using that the number of $C_i$-inputs $\mathrm{in}_i$ is $r_{\mathrm{in}}^i$-exponential as noted above, where the single and double exponential dependency on the $C_i$-depth corresponds to the use of $r_{\mathrm{cxt}}^i + 1$ and $r_{\mathrm{cxt}}^i + 2$ in Definition 23.                                                                      □

**Theorem 25.** *Let $\Sigma$ be saturating. For every database $\mathcal{D}$ the size of $\mathrm{chase}(\Sigma, \mathcal{D})$ is at most $\mathrm{rank}(\Sigma)$-exponential in the size of $\mathcal{D}$, and BCQ entailment is $\mathrm{rank}(\Sigma)$-ExpTime-complete for data complexity.*

PROOF. Theorem 24 yields a $\mathrm{rank}(\Sigma)$-exponential bound on the number of terms in $\mathrm{chase}(\Sigma, \mathcal{D})$. The number of atoms in $\mathrm{chase}(\Sigma, \mathcal{D})$, for fixed $\Sigma$, is polynomial in the number of terms. Since BCQ entailment can be decided over $\mathrm{chase}(\Sigma, \mathcal{D})$, the claimed $\mathrm{rank}(\Sigma)$-ExpTime upper bound follows.

The lower bound can be shown by reduction from the word problem of $k$-exponentially time bounded Turing machines (TMs). The simulation of TMs with Datalog rules is standard [19], using a strict total order for time steps and tape cells. To construct such an order of the required length, we expand the construction of Example 14 (i.e., the combined tgds from Examples 3 and 11) with

the following additional tgds:

$$first(z) \rightarrow min(\bot, z) \wedge max(\top, z) \wedge succ(\bot, \top, z) \tag{29}$$

$$cat(\bar{x}_1, \bar{x}_2, z, x) \wedge next(z, z_+) \wedge up(x, z_+, \bar{x}) \rightarrow cnu(\bar{x}_1, \bar{x}_2, \bar{x}, z, z_+) \tag{30}$$

$$cnu(\bar{x}_1, \bar{x}_2, \bar{x}, z, z_+) \wedge cnu(\bar{x}_1, \bar{x}_2', \bar{x}', z, z_+) \wedge succ(\bar{x}_2, \bar{x}_2', z) \rightarrow succ(\bar{x}, \bar{x}', z_+) \tag{31}$$

$$\begin{aligned} cnu(\bar{x}_1, \bar{x}_2, \bar{x}, z, z_+) \wedge cnu(\bar{x}_1', \bar{x}_2', \bar{x}', z, z_+) \wedge \\ succ(\bar{x}_1, \bar{x}_1', z) \wedge max(\bar{x}_2, z) \wedge min(\bar{x}_2', z) \rightarrow succ(\bar{x}, \bar{x}', z_+) \end{aligned} \tag{32}$$

$$cnu(\bar{x}_1, \bar{x}_1, \bar{x}, z, z_+) \wedge min(\bar{x}_1, z) \rightarrow min(\bar{x}, z_+) \tag{33}$$

$$cnu(\bar{x}_1, \bar{x}_1, \bar{x}, z, z_+) \wedge max(\bar{x}_1, z) \rightarrow max(\bar{x}, z_+) \tag{34}$$

Reading $\bot$-$\top$-sequences as binary numbers, these tgds compute the usual numeric successor order. In particular, for each $\ell$, facts of the form $succ(\bar{x}, \bar{x}', \ell)$, $min(\bar{x}, \ell)$, and $max(\bar{x}, \ell)$ describe a total order of length $2^{2^{\ell}}$. The approach is similar to a technique first introduced by Calì et al. [11], generalised to make the depth of the construction data-dependent.

We can iterate this construction by using another instance of the above tgds with each predicate $p$ replaced with a fresh name $p'$, and using tgds like $last(z) \wedge succ(\bar{x}, \bar{x}', z) \rightarrow next'(\bar{x}, \bar{x}')$ to derive an initial order of levels. The new set of tgds leads to another strongly connected component whose rank is increased by 2, since $next'$ provides context terms for the renamed version of tgd (5). This allows us to construct $k$-exponential total orders for all even numbers $k$.

To cover odd ranks as well, we can use the known simulation of nested finite sets with tgds [35]. The construction of single exponentially long chains that has been given for data complexity in this case [35, Theorem 6]. Using this construction instead of the above, we obtain exponential chains for strongly connected components $C$ with $conf(C) = 1$.

In summary, we can therefore find, for any number $r \geq 0$, tgds $\Sigma$ with $rank(\Sigma) = r$ that gives rise to an $r$-exponential chain, and for which the $r$-ExpTime-hard word problem for $k$-exponential time TMs reduces to BCQ entailment. The special case $r = 0$, where 0-ExpTime denotes PTime, agrees with the known data complexity for Datalog [19]. □

## D   PROOFS FOR SECTION 6

**Lemma 27.** *For every arboreous $\Sigma$ and chase$(\Sigma, \mathcal{D})$, the null forest is indeed a forest (set of trees).*

PROOF. We use the notation as in Definition 26. By definition, the graph of all edges $t \overset{x}{\twoheadrightarrow} n$ is acyclic, so the null forest $\langle \hat{N}, \twoheadrightarrow \rangle$ is too. Suppose for contradiction that there is $m \in \hat{N}$ with in-degree $\geq 2$, i.e., that there are $n_1, n_2 \in \hat{N}$ with $n_i \overset{y_i}{\twoheadrightarrow} m$ for $i \in \{1, 2\}$. Then $y_1 \neq y_2$, since different nulls $n_1, n_2$ must match different variables to be used in one tgd application. By definition, $var(n_1), var(n_2) \in \hat{C}$, and therefore $conf(\hat{C}) \geq 2$. A contradiction. □

**Lemma 29.** *For every arboreous $\Sigma$ and chase$(\Sigma, \mathcal{D})$, $N_\sim$ is a partition of the terms in chase$(\Sigma, \mathcal{D})$, and the term tree is indeed a tree.*

PROOF. We use the notation as in Definition 28. By definition, the sets $R[i]$ are disjoint, since each null is created by a unique tgd application. The sets $F[i]$ are disjoint from all $R[j]$ with $j \neq i$ by construction. Since $\langle \hat{N}, \twoheadrightarrow \rangle$ is a tree, the sets $F[i]$ are therefore mutually disjoint. Since $\mathsf{L}_0$ contains all remaining terms by construction, $N_\sim$ is indeed a partition of the terms of the chase.

The relation $\tilde{\twoheadrightarrow}$ is acyclic on $\mathcal{F}$. Indeed, suppose that there were a $\tilde{\twoheadrightarrow}$-cycle $C$ in $\mathcal{F}$. Each $F[i] \overset{}{\tilde{\twoheadrightarrow}} F[j]$ in $C$ corresponds to a relation $n_1 \overset{}{\twoheadrightarrow} n_2$ for nulls $n_1, n_2$ as in Definition 28. Since $F[i] \neq F[j]$, $n_2 \in R[j]$. Stemming from a single tgd application, all elements in $R[j]$ have the same predecessors, hence there is a $\tilde{\twoheadrightarrow}$-path from $n_1$ to every $n \in F[j]$. Since we obtain such paths for every edge in

the $\tilde{\twoheadrightarrow}$-cycle $C$, we find a $\tilde{\twoheadrightarrow}$-cycle in the null forest. This contradicts the acyclicity of the null forest (Lemma 27).

With the additional root element $L_0$, $\mathcal{F}$ therefore becomes a tree as required.   □

We make another small observation that was not included in the main text of the paper, but that is relevant to avoid Definition 28 from requiring further special cases.

**Lemma 43.** *Let $\Sigma$ be arboreous with related notation as in Definition 28. Then $V_{\hat{E}} \subseteq \hat{C}$, i.e., all existential variables of all $\hat{E}$-tgds are contained in $\hat{C}$.*

PROOF. Consider some edge $w \xrightarrow{y} v$ in $\hat{E}$. Then $v \in \hat{C}$. Let $\boldsymbol{v}$ be the set of existential variables in the tgd that contains $v$. Then every $v' \in \boldsymbol{v}$ also satisfies $v' \in \hat{C}$, since otherwise $v'$ would be part of a distinct strongly connected component that would have the same or a greater rank than $\hat{C}$, contradicting the requirement that $\hat{C}$ is the unique SCC of maximal rank (Definition 26).   □

The previous result clarifies possible uncertainty about the sets $R[i]$ in Definition 28: even when applied to a match that only includes terms that are not in the null forest $\hat{N}$, the resulting set of fresh nulls is fully contained in $\hat{N}$.

**Lemma 32.** *Let $\Sigma$ be arboreous. For every $\mathcal{D}$, if $p(t_1, \ldots, t_n) \in \text{chase}(\Sigma, \mathcal{D})$ with $\langle p, i \rangle \preceq \langle p, j \rangle$, then $L(t_i) \tilde{\twoheadrightarrow}^* L(t_j)$.*

PROOF. Let $\mathcal{D}^0, \mathcal{D}^1, \ldots$ be the chase sequence of $\text{chase}(\Sigma, \mathcal{D})$. We show the claim for $p(\boldsymbol{t}) \in \mathcal{D}^c$ by strong induction on $c > 0$. For $\mathcal{D}^0 = \mathcal{D}$, $t_i \notin \hat{N}$ so $t_i \in L_0$, and the claim holds since $L_0$ is the root of the term tree.

For the induction step $\mathcal{D}^{c+1}$, we only consider the case $t_i \in \hat{N}$ (the case $t_i \notin \hat{N}$ works as before), and we show that $L(t_i) \tilde{\twoheadrightarrow}^* L(t_j)$. Note that this implies $t_j \in \hat{N}$, since there is no edge from $L(t_i)$ to $L_0$.

Therefore, let $\rho := \text{tgd}[c] = B[\boldsymbol{x}, \boldsymbol{y}] \rightarrow \exists \boldsymbol{v}.H[\boldsymbol{y}, \boldsymbol{v}]$ and let $\sigma^+ := \sigma^+[c]$. Moreover, assume that $\langle p, i \rangle \preceq \langle p, j \rangle$ and $p(t_1, \ldots, t_k) = p(z_1 \sigma^+, \ldots, z_k \sigma^+) \in H\sigma^+$ with $t_i \in \hat{N}$. We have $z_i, z_j \in \boldsymbol{y} \cup \boldsymbol{v}$. We distinguish the possible cases:

(1) Case $z_i, z_j \in \boldsymbol{v}$. Since $t_i \in \hat{N}$, $z_i \in \hat{C}$. By Definition 31 (1), $z_j \in \hat{C}$ and therefore $t_j \in \hat{N}$. By Definition 28, $t_i, t_j \in R[c]$, so $L(t_i) = L(t_j)$ and $L(t_i) \tilde{\twoheadrightarrow}^* L(t_j)$.

(2) Case $z_i, z_j \in \boldsymbol{y}$. Since $\langle p, i \rangle \preceq \langle p, j \rangle$, Definition 31 (4) implies that there is a chain of body variables $z_i = x_1 \trianglelefteq \ldots \trianglelefteq x_\ell = z_j$ in $\rho$. By definition of $\trianglelefteq$, each pair of adjacent variables $x_m, x_{m+1}$ in that chain occurs at $\preceq$-comparable body positions $\langle q, i_m \rangle \preceq \langle q, i_{m+1} \rangle$ in a body atom $q(\boldsymbol{s})$ such that $q(\boldsymbol{s})\sigma^+ \in \bigcup_{d \leq c} \mathcal{D}^d$. By our induction hypothesis, the claim holds for $x_m \sigma^+$ and $x_{m+1}\sigma^+$. Since $\tilde{\twoheadrightarrow}^*$ is transitive, this shows the claim for $t_i$ and $t_j = x_\ell \sigma^+$.

(3) Case $z_i \in \boldsymbol{y}$ and $z_j \in \boldsymbol{v}$. Then $t_i \xrightarrow{z_i} t_j$ in $\text{chase}(\Sigma, \mathcal{D})$ and $\text{var}(t_i) \xrightarrow{z_i} z_j$ in $\text{LXG}(\Sigma)$ by Lemma 4. $\text{LXG}(\Sigma)$ has a unique rank-maximal strongly connected component $\hat{C}$ by Definition 28, and $\text{var}(t_i) \in \hat{C}$ by $t_i \in \hat{N}$. Therefore $\text{var}(t_i) \xrightarrow{z_i} z_j$ implies $z_j = \text{var}(t_j) \in \hat{C}$, hence $t_j \in \hat{N}$. By Definition 28, $L(t_i) = L(t_j)$ or $L(t_i) \tilde{\twoheadrightarrow} L(t_j)$, so $L(t_i) \tilde{\twoheadrightarrow}^* L(t_j)$.

(4) Case $z_i \in \boldsymbol{v}$ and $z_j \in \boldsymbol{y}$. By Definition 31 (2), $z_i \notin V_{\hat{E}}$, and by Definition 31 (3) $z_j \in \lambda(z_i)$. The latter implies that $\text{var}(t_j) \in \hat{C}$, and therefore $t_j \in \hat{N}$. Since $z_i \notin V_{\hat{E}}$, $\rho$ is not an $\hat{E}$-tgd as in Definition 28, and therefore $t_i \in L(t_j)$, i.e., $L(t_i) = L(t_j)$ and $L(t_i) \tilde{\twoheadrightarrow}^* L(t_j)$.   □

**Lemma 34.** *Let $\Sigma$ be arboreous and path-guarded, and $\mathcal{D}$ some database. For every step $i$ in $\text{chase}(\Sigma, \mathcal{D})$, with $\text{tgd}[i] = B[\boldsymbol{x}, \boldsymbol{y}] \rightarrow \exists \boldsymbol{v}.H[\boldsymbol{y}, \boldsymbol{v}]$, and $T_B = (\boldsymbol{x} \cup \boldsymbol{y})\sigma[i]$ and $T_H = (\boldsymbol{y} \cup \boldsymbol{v})\sigma^+[i]$:*

*(1) if $T_B \neq \emptyset$ then $T_B \subseteq \bigcup \text{path}(t)$ for some $t \in T_B$, and*

(2) if $T_H \neq \emptyset$ then $T_H \subseteq \bigcup \text{path}(t)$ for some $t \in T_H$.

PROOF. Item (1). Since $\text{tgd}[i]$ is path-guarded, the $\hat{C}$-affected variables $\boldsymbol{x}' \cup \boldsymbol{y}' \subseteq \boldsymbol{x} \cup \boldsymbol{y}$ form a chain in $\trianglelefteq$. By definition of $\trianglelefteq$, we can order $\boldsymbol{x}' \cup \boldsymbol{y}'$ as a sequence $z_1, \ldots, z_\ell$ such that, for every $j \in \{1, \ldots, \ell - 1\}$, there is a body atom $p(\boldsymbol{s})$ in $\text{tgd}[i]$, such that $z_j$ and $z_{j+1}$ occur at positions $\langle p, a \rangle \trianglelefteq \langle p, b \rangle$.

Applying Lemma 32 inductively to the sequence $z_1, \ldots, z_\ell$, we find that, for all $j \in \{1, \ell - 1\}$, $\mathsf{L}(z_j) \overset{*}{\tilde{\twoheadrightarrow}} \mathsf{L}(z_{j+1})$. For remaining variables $z \in \boldsymbol{x} \cup \boldsymbol{y}$ that are not $\hat{C}$-affected, we have that $z\sigma[i] \notin \hat{N}$ by Lemma 4 and the definition of $\hat{N}$; hence $\mathsf{L}(z) = \mathsf{L}_0$ in these cases. In summary, if $T_B \neq \emptyset$, then $T_B \subseteq \bigcup \text{path}(z_\ell \sigma[i])$.

Item (2). The claim is obvious for tgds with $\boldsymbol{v} = \emptyset$, since it follows from item (1) in this case. For $\boldsymbol{v} \neq \emptyset$ and $\boldsymbol{v} \cap \hat{C} = \emptyset$, the claim is also obvious since in this case $T_H \subseteq \mathsf{L}_0$. Next, consider $\boldsymbol{v} \cap \hat{C} \neq \emptyset$. Since $\text{conf}(\hat{C}) = 1$, there is at most one $y \in \boldsymbol{y}$ such that $y\sigma[i] \in \hat{N}$. If there is no such $y$, then $\boldsymbol{y}\sigma[i] \subseteq \mathsf{L}_0$.

- Case 1: $\text{tgd}[i]$ is an $\hat{E}$-tgd. Then $\boldsymbol{v} \cap \hat{C} \neq \emptyset$ implies $\boldsymbol{v} \cap V_{\hat{E}} \neq \emptyset$, and therefore $\boldsymbol{v} \subseteq V_{\hat{E}}$ by Lemma 43. Hence, $R[i] = \boldsymbol{v}\sigma^+[i]$ and $F[i]$ has no predecessors in $\mathcal{F}$, so $\mathsf{L}_0 \tilde{\twoheadrightarrow} F[i]$. This shows that $T_H \subseteq \bigcup \text{path}(n)$ for any $n \in \boldsymbol{v}\sigma^+[i]$.
- Case 2: $\text{tgd}[i]$ is not an $\hat{E}$-tgd. Then $T_H \subseteq \mathsf{L}_0$, so the claim holds too.

Alternatively, assume that there is $y$ with $y\sigma[i] \in \hat{N}$.

- Case 1: $\text{tgd}[i]$ is an $\hat{E}$-tgd. Then $R[i] = \boldsymbol{v}\sigma^+[i]$ follows as in Case 1 above. Since $y\sigma[i] \overset{\cdot}{\tilde{\twoheadrightarrow}} \boldsymbol{v}\sigma^+[i]$ holds for any $v \in \boldsymbol{v}$, we have $\mathsf{L}(y\sigma[i]) \tilde{\twoheadrightarrow} \mathsf{L}(v\sigma^+[i]) = F[i]$. We obtain $T_H \subseteq \bigcup \text{path}(v\sigma^+[i])$ since all $y' \in \boldsymbol{y}$ with $y' \neq y$ are such that $y'\sigma[i] \in \mathsf{L}_0$ (since $\text{conf}(\hat{C}) = 1$).
- Case 2: $\text{tgd}[i]$ is not an $\hat{E}$-tgd. Then $\mathsf{L}(y\sigma^+[i]) = \mathsf{L}(v\sigma^+[i])$ for every $v \in \boldsymbol{v}$, and hence $T_H \subseteq \bigcup \text{path}(y\sigma^+[i])$ using the same argument as in the previous case for $y' \neq y$. □

**Lemma 35.** *Let $\mathcal{I}^*$ be the union of all sets $\mathcal{I}_i^j$ of some run of Algorithm 1. There is a homomorphism $\tau : \mathcal{I}^* \to \text{chase}(\Sigma, \mathcal{D})$, and therefore $\text{chase}(\Sigma, \mathcal{D}) \models q$ whenever Algorithm 1 returns* true.

PROOF. We iteratively define $\tau$ for fresh nulls introduced during a run of Algorithm 1, and verify the claimed homomorphism property for each step. Since the outer loop does not matter here, we use $\mathcal{I}_0, \ldots, \mathcal{I}_\ell$ to denote the entire sequence of values for $\mathcal{I}$ as they occur throughout the algorithm.

Initially, $\tau$ is the identity function on constants in $\Sigma$ and $\mathcal{D}$, which is a homomorphism from the initially empty set $\mathcal{I}_0$ to $\text{chase}(\Sigma, \mathcal{D})$.

Now by way of induction, assume that $\tau$ has been defined so that it is a homomorphism $\bigcup_{i=0}^n \mathcal{I}_i \to \text{chase}(\Sigma, \mathcal{D})$ for some $n \geq 0$, and that a further tgd application with tgd $\rho$ and match $\sigma$ is chosen in (L5). Since $\sigma$ is a match for $\rho$ on $\mathcal{I}_n$, we find a corresponding match $\sigma_\perp$ of $\rho$ on $\text{chase}(\Sigma, \mathcal{D})$ where $\sigma_\perp(z) = \tau(\sigma(z))$ for all variables $z$ in the body of $\rho$. Since this match $\sigma_\perp$ is satisfied in $\text{chase}(\Sigma, \mathcal{D})$, it can be extended to a match $\sigma_\perp^+$ such that $\text{head}(\rho)\sigma_\perp^+ \subseteq \text{chase}(\Sigma, \mathcal{D})$. Therefore, given the extended match $\sigma^+$ that is used in Algorithm 1 to apply $\rho$, we define $\tau(v\sigma^+)$ for all existential variables $v$ in $\rho$ as $\tau(v\sigma^+) = v\sigma_\perp^+$. Then $\tau$ is a homomorphism $\bigcup_{i=0}^{n+1} \mathcal{I}_i \to \text{chase}(\Sigma, \mathcal{D})$ as required. □

**Lemma 36.** *If $\Sigma$ is arboreous and path-guarded with $\text{rank}(\Sigma) > 0$, and $M \leq f(|\mathcal{D}|)$ for some $\text{rank}(\Sigma)$-exponential function $f$, then there is a $(\text{rank}(\Sigma) - 1)$-exponential function $g$ such that Algorithm 1 runs in space $g(|\mathcal{D}|)$, where 0-exponential means polynomial.*

PROOF. Using binary encoding, the numbers $i \leq M \leq f(|\mathcal{D}|)$ can be stored in $(\text{rank}(\Sigma) - 1)$-exponential space. To show that $\mathcal{I}$ can be stored in $(\text{rank}(\Sigma) - 1)$-exponential space, note that the sets of $\mathcal{T}$, other than the root, correspond to nodes $\mathsf{L}(n)$ of the term tree in the following sense: if

$\mathcal{T}[d]$ ($d \in \{1, \ldots, |\mathcal{T}|\}$) is the $(d+1)$-th element in $\mathcal{T}$ (the first $\mathcal{T}[0]$ being the root), then there is a node $\mathsf{L} \in N_\sim$ that is $d$ steps away from the root such that $\tau(\mathcal{T}[d]) \subseteq \mathsf{L}$, with $\tau$ as in Lemma 35.

The terms in $\bigcup_{d=1}^{|\mathcal{T}|} \mathcal{T}[d]$ are therefore always contained in a single path of the term tree. The length of paths of the null forest (and analogously in the term tree) are bounded by a $(\mathrm{rank}(\Sigma) - 1)$-exponential function, as shown in the proof of Theorem 24. Indeed, that proof establishes the $C_i$-depth of nulls in a strongly connected component $C_i$ is exponentially bounded in $r_{\mathrm{cxt}}^i$ and polynomially bounded in $r_{\mathrm{in}}^i$. For $\mathrm{rank}(C_i) > 0$, the latter corresponds to a $(r_{\mathrm{in}}^i - 1)$-exponential bound. The claim about $\hat{C}$ follows since path lengths in the null forest corresponds to the $\hat{C}$-depth of nulls, and since $\mathrm{rank}(\Sigma) = \mathrm{rank}(\hat{C}) \geq \max\{r_{\mathrm{in}}, r_{\mathrm{cxt}} + 1\}$ by Definition 23.

The size of the sets $\mathsf{L}(n)$ is polynomial in $|\mathcal{D}|$, since $\mathsf{L}(n)$ only contains nulls from applications of non-$\hat{E}$-tgds, which do not have a dependency cycle, so that the polynomial data complexity of jointly acyclic tgds applies [31]. Therefore, the size of $\bigcup_{d=1}^{|\mathcal{T}|} \mathcal{T}[d]$ is $(\mathrm{rank}(\Sigma) - 1)$-exponentially bounded in $|\mathcal{D}|$.

This bound also applies to the initial value of $\mathcal{T}$ from (L2), so that after $|q|$ executions of loop (L3), the bound remains $(\mathrm{rank}(\Sigma) - 1)$-exponential (note that $|q|$ is constant with respect to $\mathcal{D}$).

With the overall set of available terms restricted by a $(\mathrm{rank}(\Sigma) - 1)$-exponential bound in $|\mathcal{D}|$, this bound carries over to the possible atoms in $\mathcal{I}$ throughout the computation. The final check in (L14) can also be performed in this space bound, e.g., by iterating over all possible variable bindings with respect to $\mathcal{T}$. □

**Lemma 37.** *If all tgd applications in a run of Algorithm 1 correspond to chase steps, and $\tau$ in each iteration is the canonical extension for the respective step, then $\tau$ is a homomorphism $\mathcal{I}^* \to \mathrm{chase}(\Sigma, \mathcal{D})$ that is injective on all term sets $\bigcup \mathcal{T}_i^j$ that occur during the run.*

PROOF. Note that our requirements for "corresponding to a chase step" already include the injectivity of $\tau$ on the terms used in the premise. Nevertheless, the claim is still non-trivial, since the final iteration of each run of the inner loop in Algorithm 1 is not covered by the requirements, and since the algorithm, by virtue of being able to non-deterministically break the computation at any time, can certainly perform runs that satisfy the preconditions. In particular, the required injectivity holds for the initial term set $C_0$ for which $\tau$ was defined as the identity.

We proceed by induction. Consider the tgd application that produces $\mathcal{I}_i^{j+1}$ from $\mathcal{I}_i^j$. By assumption, it corresponds to a chase step $s$. Let $v$ be the set of existential variables in $\mathrm{tgd}[s]$.

Now suppose for a contradiction that the canonical extension of $\tau$ in this step is not injective. By our definition, $\tau$ induces a bijection $v\sigma^+ \to v\sigma^+[s]$, and it is injective on $\mathcal{T}_i^j$ (a precondition for the application corresponding to step $s$). Hence, the supposed violation of injectivity requires that there is a null $n \in v\sigma^+$ and a term $t \in \bigcup \mathcal{T}_i^j$ such that $\tau(n) = \tau(t)$. Since $\tau(n) \in v\sigma^+[s]$, $t$ must also be a null (constants are always mapped to themselves in $\tau$). By the assumption, $\tau$ has been defined through a series of canonical extensions, so the value $\tau(t)$ was assigned in a previous tgd application that also corresponded to step $s$ (since $\tau(t) \in v\sigma^+[s]$ can be a fresh null only for this one step). Let $\theta^+$ be the extended match used in this tgd application (it has to agree with $\sigma$ on universal variables, but must use different nulls), hence $t \in v\theta^+$. But then $t$ was added to $\mathcal{T}$ in (L9) or (L11). In either case, the whole set $v\theta^+$ occurs in the same set of $\mathcal{T}$ that also contains $t$, so that $v\theta^+ \subseteq \bigcup \mathcal{T}_i^j$. In this case, however, $\mathrm{tgd}[s]\theta^+ \subseteq \mathcal{I}_i^j$, so $\mathrm{tgd}[s]$ is not applicable to obtain $\mathcal{I}_i^{j+1}$. A contradiction. □

**Lemma 39.** *If $\Sigma$ is arboreous and path-guarded, $s$ is the sequence of chase steps obtained from the task tree with root $\langle 1, i \rangle$, then the length $|s|$ of $s$ is bounded by a $\mathrm{rank}(\Sigma)$-exponential function.*

Proof. Let the *knobbly term tree* be obtained from the term tree by simultaneously replacing each node set $\mathsf{L} \in N_\sim$ with the union $\mathsf{L} \cup \bigcup_{\mathsf{L} \twoheadrightarrow \mathsf{L}_c} \mathsf{L}_c$ that also includes all terms in the node's direct children. The tree structure otherwise remains the same, i.e., the term tree and the knobbly term tree are isomorphic. In particular, the length of paths in the knobbly term tree is bounded by a $(\mathrm{rank}(\Sigma) - 1)$-exponential function, as observed for the term tree in the proof of Lemma 36.

Now consider any path $\langle d_1, s_1 \rangle \cdots \langle d_\ell, s_\ell \rangle$ in the task tree, and let $\alpha_i$ denote the atoms $\alpha$ of Definition 38 for every subtask $\langle d_i, s_i \rangle$ with $1 < i \leq \ell$.

For a path $\mathfrak{p}$, let $\mathfrak{p}|_d$ denote the path of the initial $d$ nodes in $\mathfrak{p}$. We claim that for all $i \in \{2, \dots, \ell\}$, $\mathrm{path}(\alpha_i)|_{d_i-1} = \mathrm{path}(\alpha_\ell)|_{d_i-1}$, i.e., the paths of atoms $\alpha_i$ agree with the path of $\alpha_\ell$, except possibly for the lowest node (Claim ‡). This is trivial for $i = \ell$. For a task $\langle d_i, s_i \rangle$ with $1 < i < \ell$, assume by way of induction that the claim was shown for $\langle d_{i+1}, s_{i+1} \rangle$. Then $\mathrm{path}(\alpha_{i+1})|_{d_{i+1}-1} = \mathrm{path}(\alpha_\ell)|_{d_{i+1}-1}$ by induction hypothesis, and $\mathrm{path}(\alpha_{i+1})|_{d_{i+1}-1} \subseteq \mathrm{path}(\alpha_{i+1}) \subseteq \mathrm{path}_B(s_i)$ by Definition 38. Thus, $\mathrm{path}_B(s_i)|_{d_{i+1}-1} \subseteq \mathrm{path}(\alpha_\ell)|_{d_{i+1}-1}$, and, as $d_{i+1} \geq d_i$, $\mathrm{path}_B(s_i)|_{d_i-1} \subseteq \mathrm{path}(\alpha_\ell)|_{d_i-1}$ (∗). Since $|\mathrm{path}(\alpha_i)| = d_i$, there are two cases:

- $\mathrm{path}(\alpha_i) = \mathrm{path}_B(s_i)|_{d_i}$ (then $\alpha_i$ contains no new nulls, or $\mathrm{tgd}[s_i]$ is not an $\hat{E}$-tgd).
- $\mathrm{path}(\alpha_i)$ extends the path $\mathrm{path}_B(s_i)|_{d_i-1}$ by one additional child node (then $\alpha_i$ contains new nulls from the $\hat{E}$-tgd $\mathrm{tgd}[s_i]$).

In either case, $\mathrm{path}(\alpha_i)|_{d_i-1} \subseteq \mathrm{path}_B(s_i)|_{d_i-1} \subseteq \mathrm{path}(\alpha_\ell)|_{d_i-1}$ as required, where the second $\subseteq$ is (∗).

Now let $\mathrm{path}_\bullet(\alpha_\ell)$ denote the path in the knobbly term tree that corresponds to $\mathrm{path}(\alpha_\ell)$ by the isomorphism. By Claim ‡, for every $i \in \{2, \dots, \ell\}$, $\mathrm{path}(\alpha_i) \subseteq \mathrm{path}_\bullet(\alpha_\ell)$. Containment is clear for $\mathrm{path}(\alpha_i)|_{d_i-1}$ by Claim ‡. Since $|\mathrm{path}(\alpha_i)| = d_i$, the path has at most one additional final node not in $\mathrm{path}(\alpha_\ell)$, and the terms of this node, being a direct child, are included in $\mathrm{path}_\bullet(\alpha_\ell)$.

Since $|\mathrm{path}_\bullet(\alpha_\ell)|$ is bounded by a $(\mathrm{rank}(\Sigma) - 1)$-exponential function, and since the term sets that constitute the nodes are still of constant size (being unions of terms generated by jointly-acyclic sets of tgds, cf. proof of Lemma 36), the cardinality of $\bigcup \mathrm{path}_\bullet(\alpha_\ell)$ is also bounded by a $(\mathrm{rank}(\Sigma) - 1)$-exponential function. But then, given the fixed signature of $\Sigma$, there are at most $(\mathrm{rank}(\Sigma) - 1)$-exponentially many atoms that can play the role of $\alpha_i$ ($2 \leq i \leq \ell$) in the above path, and since each atom is produced in just one chase step, the path corresponds to a (strictly decreasing) sequence of at most $(\mathrm{rank}(\Sigma) - 1)$-exponentially many chase steps.

This shows that the depth of the task tree is bounded by a $(\mathrm{rank}(\Sigma) - 1)$-exponential function, so the size of the task tree (and of the induced sequence of steps) is bounded by a $\mathrm{rank}(\Sigma)$-exponential function. □

**Lemma 40.** *If $\Sigma$ is arboreous and path-guarded, $s$ is the sequence of chase steps obtained from the task tree for $\langle 1, i \rangle$, and $M \geq |s|$, then Algorithm 1 can choose tgd applications according to $s$.*

Proof. The claim refers to a single execution of the inner loop of Algorithm 1, starting from $\mathcal{I}_0 = \mathcal{D}$. Let $\ell = |s|$ be the length of $s$, let $\mathcal{I}_i$ for $1 \leq i \leq \ell$ denote the value of $\mathcal{I}$ after $i$ iterations, and let $s[i]$ be the $i$th element of $s$. Moreover, let $\mathrm{task}[i]$ be the task with label $\langle d, s[i] \rangle$ that gave rise to $s[i]$ in $s$, and let $\mathrm{depth}[i] = d$ be its depth. If index $i$ corresponds to a subtask in the term tree, let $\alpha[i]$ be the atom $\alpha$ of Definition 38 that justified its inclusion as a child node. As before, given a path $\mathfrak{p}$, we write $\mathfrak{p}|_d$ for the path of the initial $d$ nodes in $\mathfrak{p}$.

We show by induction over $i \in \{1, \dots, \ell\}$ that for all atoms $\alpha \in \mathrm{chase}(\Sigma, \mathcal{D})$ with $\alpha \in \mathcal{D}^{s[i]-1}$ and $\mathrm{path}(\alpha) \subseteq \mathrm{path}_B(s[i])$, there is a corresponding atom $\beta = \tau^-(\alpha) \in \mathcal{I}_{i-1}$, where $\tau^-$ is well-defined by Lemma 37. In particular, this shows that

(1) the match $\sigma[s[i]]$ of $\mathrm{tgd}[s[i]]$ has a corresponding match on $\mathcal{I}_{i-1}$ via $\tau^-$,

(2) if tgd[$s[i]$] contains an existential variable, then no Datalog rule in $\Sigma$ is applicable to $\mathcal{I}_{i-1}$, and

(3) Algorithm 1 can execute all non-redundant tgd applications in the given sequence, and $\mathcal{I}_i$ contains an instance of the head of tgd[$s[i]$].

Item (1) is clear. Item (2) follows since the sequence of chase steps in chase($\Sigma, \mathcal{D}$) also respects the Datalog-first condition, and since the conclusions of Datalog rules over $\mathcal{I}_{i-1}$ can only use the terms in $\mathcal{I}_{i-1}$, and in particular are in path($\beta$) $\subseteq$ path$_B(s[i])$. Together, (1) and (2) ensure that tgd[$s[i]$] is applicable at step $i$ of Algorithm 1 if its head is not already satisfied in $\mathcal{I}_{i-1}$. Note that the latter case can only occur if the same tgd application has been performed before, since earlier chase steps $< s[i]$ have not prevented the application of tgd[$s[i]$] in chase($\Sigma, \mathcal{D}$). In this situation, we ignore step $i$ and continue immediately with the next choice $i + 1$ (if any), and we let $\mathcal{I}_i = \mathcal{I}_{i-1}$. Hence we also obtain (3).

Now let $i \in \{1, \ldots, \ell\}$ and assume that the induction claim holds true for all $i' < i$. Consider an arbitrary atom $\beta$ as in the claim. Let $d_\beta := |\text{path}(\beta)|$ be the depth of $\beta$, and let $s_\beta$ be the chase step that produced $\beta \in$ chase($\Sigma, \mathcal{D}$). We claim that $\beta \in \mathcal{I}_{i-1}$.

Case (i). If $d_\beta < \text{depth}[i]$, then depth[$i$] $> 1$. Let $k$ be the ancestor node of $i$ that is closest to $i$ (i.e., lowest in the task tree), such that depth[$k$] $\leq d_\beta$. By Definition 38, $s[k] > s[i] > s_\beta$, so $k$ has a child node $j$ with label $\langle d_\beta, a \rangle$ where $a \geq s_\beta$. Then $\beta \in \mathcal{I}_j$ by the induction hypothesis. Moreover, due to the traversal order of children of $k$, all nodes between position $j$ and $i$ have depth $> d_\beta$. This ensures that $\beta \in \mathcal{I}_{i-1}$ as required (we give a more detailed account of this argument for a slightly more general situation in Case (ii)).

Case (ii). If $d_\beta \geq \text{depth}[i]$, then task[$i$] has a descendant node $j$ the task tree with label $\langle d_\beta, s_\beta \rangle$. This is easy to see for $d_\beta = 1$, since the path of depth 1 is unique, so that the condition path($\alpha$) $\subseteq$ path$_B(i)$ in Definition 38 is tautological if $|\text{path}(\alpha)| = 1$. Hence the chase steps for atoms at depth 1 appear within a single path below $i$ (with chase steps of such atoms in decreasing order, largest first).

For $d_\beta > 1$, we find a similar path below task $i$. Care is needed since the condition in Definition 38 refers to the body path of the immediate parent node, which may not be the body path of $i$, since only the nodes up to depth $d_\beta - 1$ are stable yet. We therefore make the following observation: The earliest chase step $c$ that produces an atom $\beta$ such that $|\text{path}(\beta)| = d_\beta$ and path($\beta$) $\subseteq$ path$_B(i)$ must be the step that introduced the set of nulls denoted $R[c]$ in Definition 28, i.e., that initialised the node $F[i]$ of path$_B(i)$ at depth $d_\beta$. All other chase steps $a > c$ that infer an atom $\alpha$ with $|\text{path}(\alpha)| = d_\beta$ and path($\alpha$) $\subseteq$ path$_B(i)$ have a frontier variable that is matched to a term in $F[i]$. Therefore, $F[i]$ is a node in path$_B(a)$, and every atom $\gamma$ with path($\gamma$) $\subseteq$ path$_B(i)$ also satisfies path($\gamma$) $\subseteq$ path$_B(a)$. The chase steps that produce such atoms therefore form a sequence $c < a_1 < \ldots < a_m$, and we find an according path of tasks $\langle d_\beta, a_m \rangle \rightarrow \cdots \rightarrow \langle d_\beta, a_1 \rangle \rightarrow \langle d_\beta, c \rangle$ in the task tree. This finishes the argument that we find the claimed descendant node $j$ of $i$.

Then $j < i$ since the task tree is traversed in topological order. Let $\rho_\beta := \text{tgd}[s_\beta]$, $\sigma_\beta := \sigma[s_\beta]$, and $\sigma_\beta^+ := \sigma^+[s_\beta]$. By the induction hypothesis, the tgd application for this step $s_\beta = s[j]$ succeeded with a match $\theta = \tau^-(\sigma_\beta)^{[3]}$, and was performed with an extended match $\theta^+ = \tau^-(\sigma_\beta^+)$. However, it is possible that the deletions in $\mathcal{T}$ between step $j$ and step $i$ were such that $\tau^-$ (which is only defined locally) is not the same at both steps, hence we cannot yet conclude $\beta \in \mathcal{I}_j$ but merely that $\beta' \in \mathcal{I}_j$ for some variant of $\beta$ that might use different fresh nulls.

We therefore show by induction that all intermediate steps $k$ with $j < k < i$ are such that, after Algorithm 1 has executed (L6), $\mathcal{T}$ has length $\geq d_\beta$. This shows that any fresh nulls of step

---

[3]We write $\tau^-(\sigma_\beta)$ for the function that maps $z$ to $\tau^-(z\sigma_\beta)$. This is sometimes denoted as $\tau^- \circ \sigma_\beta$ but sometimes also as $\sigma_\beta \circ \tau^-$; we avoid this confusion.

$j$ are still in $\mathcal{T}$ at step $i$, and this implies $\beta' = \beta \in \mathcal{I}_{i+1}$ as required. Since $j$ is part of the path $\langle d_\beta, a_m \rangle \rightarrow \cdots \rightarrow \langle d_\beta, a_1 \rangle \rightarrow \langle d_\beta, c \rangle$ as defined above, there are two options for nodes $k$: (i) task$[k] = \langle d_\beta, a \rangle$ for some $a \in \{a_1, \ldots, a_m\}$, or (ii) task$[k] = \langle e, a \rangle$ for some $e > d_\beta$. In case (i), as noted above, tgd$[a]$ has a frontier variable that is matched to a term in $F[i]$, which shows the claim about $k$ since $F[i]$ is the element at position $d_\beta$ in $\mathcal{T}$ by induction hypothesis. In case (ii), the claim likewise follows since tgd$[a]$, in order for $\alpha[k]$ to be at depth $e$, has a frontier variable at depth $\geq e - 1 \geq d_\beta$.

This concludes the proof that $\mathcal{I}_{i-1}$ contains all atoms that were inferred at a chase step before $s[i]$ and use terms in path$_B(s[i])$. By (1)–(3) above, this completes the proof. □

**Theorem 41.** *BCQ entailment for arboreous and path-guarded tgd sets $\Sigma$ with rank$(\Sigma) = \kappa \geq 1$ is $(\kappa - 1)$-ExpSpace-complete for data complexity, where 0-ExpSpace = PSpace.*

Proof. Membership follows from the correctness (Lemma 35) and completeness (Lemmas 36, 39, and 40) of the tree-based chase, where the algorithm follows $|q|$ step sequences based on the $|q|$ atoms in a particular query match to materialise the match in the $|q|$ iterations of the outer loop in (L3). Algorithm 1 therefore decides query entailment in $(\kappa - 1)$-NExpSpace, and we get membership in $(\kappa - 1)$-ExpSpace by Savitch's Theorem.

We show hardness via reduction from the word problem for alternating Turing Machines (ATMs) with a $(\kappa - 1)$-exponential time bound, which is $(\kappa - 1)$-ExpSpace-complete [17].

Let $\mathcal{M} = \langle Q, \Gamma, \Delta, q_0, t \rangle$ be an ATM with a finite set of states $Q$, a finite tape alphabet $\Gamma$ consisting of an input alphabet and a special symbol ␣ (blank), a transition relation $\Delta \subseteq (Q \times \Gamma) \times (Q \times \Gamma \times \{l, r\})$, an initial state $q_0 \in Q$, and a function $t : Q \rightarrow \{\exists, \forall, \text{acc}\}$ that marks states as existential, universal, or accepting. As usual, we write configurations as $w_l q w_r$ where $w_l$ denotes the tape symbols left of the read-write head, $q \in Q$ is the current state, and $w_r$ is a sequence of symbols to the right of the read-write head (the first symbol of which is underneath the head). Tape symbols right of $w_r$ are assumed to be ␣. For a configuration $C = w_l q \sigma w_r$ and a transition $\langle \langle q, \sigma \rangle, \langle q^+, \sigma^+, d \rangle \rangle \in \Delta$, there is a successor configuration

$$C^+ := \begin{cases} w_l \sigma^+ q^+ w_{r␣} & \text{if } d = r, \\ \hat{w}_l q^+ \sigma_l \sigma^+ w_r & \text{if } d = l \text{ and } w_l = \hat{w}_l \sigma_l \text{ for some } \sigma_l \in \Gamma. \end{cases}$$

A configuration $w_l q w_r$ is accepting if (1) $t(q) = \text{acc}$, (2) $t(q) = \exists$ and there is an accepting successor configuration, or (3) $t(q) = \forall$ and all successor configurations are accepting. $\mathcal{M}$ accepts a word $w$ if the initial configuration $q_0 w$ is accepting.

For a word $w = \sigma_1 \ldots \sigma_n$ with $n > 0$ (the case of the empty word is irrelevant for hardness), let $\mathcal{D}_w$ denote the database with the facts $first_0(1)$, $next_0(i, i+1)$ for $i \in \{1, \ldots, n-1\}$, $last_0(n)$, and $symbol(i, \sigma_i)$ for $i \in \{1, \ldots, n\}$. Let $\Sigma_\leq$ denote a set of tgds constructed as in the proof of Theorem 25 to entail a $(\kappa-1)$-exponentially long chain over the initial chain encoded in predicates $first_0$, $next_0$, and $last_0$. We assume that the constructed $(\kappa-1)$-exponential chain is encoded predicates $first$, $next$, and $last$. Moreover, we extend $\Sigma_\leq$ with the following rules to encode the transitive (non-reflexive) closure of $next$:

$$next(x, y) \rightarrow next^+(x, y) \tag{35}$$

$$next(x, y) \land next^+(y, z) \rightarrow next^+(x, z) \tag{36}$$

We use facts $step(c, i)$ to encode that $c$ is an $i$th configuration in a run, i.e., there is a sequence $c_0, \ldots, c_i = c$ with $c_0$ being the initial configuration, facts $q(c, q)$ to encode the state $q$ of configuration $c$, facts $hpos(c, i)$ to encode that the head of $\mathcal{M}$ is at the $i$th position of the tape in configuration $w$, and facts $tape(c, i, s)$ to encode that the $i$th position of the tape of configuration $c$ is symbol $s$.

The following tgds $\Sigma_{\text{init}}$ then establish the initial configuration $c_0$ (with constants $c_0$, $q_0$, and $\textvisiblespace$):

$$first(i) \rightarrow step(c_0, i) \wedge q(c_0, q_0) \wedge hpos(c_0, i) \tag{37}$$

$$first(i) \wedge first_0(j) \rightarrow samepos(i, j) \tag{38}$$

$$samepos(i, j) \wedge next(i, i^+) \wedge next_0(j, j^+) \rightarrow samepos(i^+, j^+) \tag{39}$$

$$samepos(i, j) \wedge symbol(j, s) \rightarrow tape(c_0, i, s) \tag{40}$$

$$samepos(i, j) \wedge last_0(j) \wedge next^+(i, i^+) \rightarrow tape(c_0, i, \textvisiblespace) \tag{41}$$

The tgd (37) initialises step, state, and head position of the initial configuration $c_0$. The tgds (38) and (39) match elements of the initial chain to the first elements of the $(\kappa - 1)$-exponential chain. This is then used to transcribe the input word to the initial tape (40), with the remaining tape filled with blanks (41).

Next, we encode how to construct a null for the successor of a configuration. For this purpose, we introduce predicates $to_\delta$ for every $\delta \in \Delta$, each encoding successor configurations for this transition in a certain step. Concretely, for all $\delta, \delta' \in \Delta$ with $\delta = \langle\langle q, \sigma \rangle, \langle q^+, \sigma^+, d \rangle\rangle$, the tgds $\Sigma_+$ contain the following tgds (with constants $q$ and $\sigma$):

$$step(C, i) \wedge next(i, i^+) \wedge hpos(C, p) \wedge tape(C, p, \sigma) \wedge q(C, q) \rightarrow \exists v_\delta. \, to_\delta(i^+, C, v_\delta) \tag{42}$$

$$to_\delta(i, C^-, C) \rightarrow to_{\delta'}(i, C, C) \tag{43}$$

$$to_\delta(i, C^-, C) \wedge to_{\delta'}(j, C^-, C^-) \rightarrow to_{\delta'}(j, C, C) \tag{44}$$

The tgd (42) creates a new null for each valid transition $\delta$ from a configuration $C$ with state $q$ and $s$ at the head position. Facts of the form $to_\delta(i^+, C, C^+)$ encode that $C^+$ is a $\delta$-successor of $C$ at depth $i^+$. The idea is that each $i$ of the linear order encoded by $next$ can be used at most once per configuration path $c_1, \ldots, c_n$ to create a successor.

$\Sigma_+$ will induces a saturating strongly connected component in the overall labeled dependency graph. Indeed, Definition 13 is satisfied for the set $E = \{v_\delta \overset{C(\delta')}{\rightarrow} v'_\delta \mid \delta, \delta' \in \Delta\}$, where we use $C(\delta')$ to disambiguate the universal variables $C$ in tgds of the form (42). Note that the empty path is the only $\bar{E}$-path in this case. Now, for tgd $\rho_{\delta'}$ with existential variable $v_{\delta'}$, tgds (43) ensure base-propagation for $e = v'_\delta \overset{C(\delta')}{\rightarrow} v_{\delta'} \in E$ and tgds (44) ensure step-propagation for $e_1, e_2 = v'_\delta \overset{C(\delta')}{\rightarrow} v_{\delta'} \in E$.

Next, the tgds $\Sigma_\rightarrow$ encode the $\delta$-successor $C^+$ of a configuration $C$ based on the encoding of $C$. For $\delta = \langle\langle q, \sigma \rangle, \langle q^+, \sigma^+, d \rangle\rangle$, $\Sigma_\rightarrow$ contains the following tgds:

$$step(C, i) \wedge next(i, i^+) \wedge to_\delta(i^+, C, C^+) \rightarrow step(C^+, i^+) \wedge succ_\delta(C, C^+) \wedge q(C^+, q^+) \tag{45}$$

$$hpos(C, p) \wedge next^+(p, p') \wedge tape(C, p', s') \wedge succ_\delta(C, C^+)$$
$$\rightarrow tape(C^+, p', s') \tag{46}$$

$$hpos(C, p) \wedge next^+(p', p) \wedge tape(C, p', s') \wedge succ_\delta(C, C^+)$$
$$\rightarrow tape(C^+, p', s') \tag{47}$$

$$hpos(C, p) \wedge succ_\delta(C, C^+) \rightarrow tape(C^+, p, \sigma^+) \tag{48}$$

Here, the tgd (45) writes the step and state of $C^+$ and introduces a notion $succ_\delta(C, C^+)$, which states that $C^+$ is the 'real' $\delta$-successor of $C$ (tgds (43) and (44) makes $C^+$ a pseudo-successors of itself for all $i \leq i^+$). The tgds (46)–(48) write the tape of configuration $C^+$. Moreover, $\Sigma_\rightarrow$ requires tgds to update the head position, and we make a case distinction based on the direction $d$ to which the head moves. If $d = l$, $\Sigma_\rightarrow$ contains the tgd:

$$hpos(C, p) \wedge next(p', p) \wedge succ_\delta(C, C^+) \rightarrow hpos(C^+, p') \tag{49}$$

and, otherwise, if $d = r$, $\Sigma_{\rightarrow}$ contains the tgd:

$$hpos(C, p) \wedge next(p, p') \wedge to_\delta(C, C^+) \rightarrow hpos(C^+, p') \tag{50}$$

Finally, we define tgds to evaluate which configurations are accepting. The set of tgds $\Sigma_{\text{eval}}$ contains the following tgds, where we introduce further predicates $acc$ and $acc_\delta$ for each $\delta \in \Delta$.

1. For every state $q_a \in Q$ with $t(q_a) = $ acc:

$$q(C, q_a) \rightarrow acc(C) \tag{51}$$

2. For every state $q_\exists \in Q$ with $t(q_\exists) = \exists$, and every transition $\delta = \langle\langle q_\exists, \sigma\rangle, \langle q^+, \sigma^+, d\rangle\rangle \in \Delta$:

$$succ_\delta(C, C^+) \wedge acc(C^+) \rightarrow acc(C) \tag{52}$$

3. For every state $q_\forall \in Q$ with $t(q_\forall) = \forall$, and every symbol $\sigma \in \Gamma$ with $\delta_1, \ldots, \delta_n$ a list of all transitions of form $\langle\langle q_\forall, \sigma\rangle, \langle q^+, \sigma^+, d\rangle\rangle \in \Delta$:

$$succ_{\delta_1}(C, C^+) \wedge acc(C^+) \rightarrow acc_{\delta_1}(C) \tag{53}$$
$$acc_{\delta_i}(C) \wedge succ_{\delta_{i+1}}(C, C^+) \wedge acc(C^+) \rightarrow acc_{\delta_{i+1}}(C) \tag{54}$$
$$acc_{\delta_i}(C) \rightarrow acc(C) \tag{55}$$

The tgds recursively mark configurations as accepting. Concretely, (51) directly marks all configuration with an accepting state, and (52) marks configurations with an existential state if they have an accepting successor. The tgds (53)–(55) mark a configuration with a universal state if all of its successors are accepting. To achieve this, the successors are traversed in an arbitrary but fixed order $\delta_1, \ldots, \delta_n$. Semantically, one could just combine these rules into one, but the above splitting ensures path-guardedness. Indeed, the relation $\leq$ of Definition 31 for the above sets of tgds contains $\langle to_\delta, 2\rangle \leq \langle to_\delta, 3\rangle$ and $\langle succ_\delta, 1\rangle \leq \langle succ_\delta, 2\rangle$. Therefore, all tgds are path-guarded.

To conclude, let $\Sigma = \Sigma_\leq \cup \Sigma_{\text{init}} \cup \Sigma_+ \cup \Sigma_{\rightarrow} \cup \Sigma_{\text{eval}}$. As $\Sigma$ constructs and evaluates the configuration tree of $\mathcal{M}$, we obtain that $\mathcal{M}$ accepts $w$ if and only if chase$(\Sigma, \mathcal{D}_w) \models acc(c_0)$. Since $\Sigma$ is saturating with rank$(\Sigma) = \kappa$, arboreous, and path-guarded, this shows the claim.      $\square$