



HAL
open science

Ontology-Based Query Answering Over Datalog-Expressible Rule Sets Is Undecidable

Lucas Larroque

► **To cite this version:**

Lucas Larroque. Ontology-Based Query Answering Over Datalog-Expressible Rule Sets Is Undecidable. Computer Science [cs]. 2023. hal-04347020

HAL Id: hal-04347020

<https://inria.hal.science/hal-04347020>

Submitted on 15 Dec 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License



ONTOLOGY-BASED QUERY
ANSWERING OVER
DATALOG-EXPRESSIBLE RULE
SETS IS UNDECIDABLE

Lucas LARROQUE

Supervised by Michaël Thomazo and David Carral

Master Logique et Fondements de l'Informatique
Septembre 2023

Contents

1	Introduction	2
2	Preliminaries	4
2.1	Syntax	4
2.2	Semantics	5
2.3	The chase	5
2.4	Some computability theory	6
3	OBQA over Datalog-expressible rule sets is undecidable	8
3.1	The rule set	8
3.2	\mathcal{R}_M simulates M	9
3.3	\mathcal{R}_M is Datalog expressible	16
3.3.1	Skolem depth	17
3.3.2	An axiomatization of equality	18
3.3.3	Skolem-increasing derivations	22
3.3.4	Some structure — The Fail predicate	25
3.3.5	Controlling the chains	26
3.3.6	Skolem depth and Datalog-expressibility	30

Chapter 1

Introduction

To handle the ever-increasing amount of data available, and to process it efficiently, querying data has become a ubiquitous task. In some cases, one also needs to be able to reason over data. A standard use case is the harmonization of databases coming from different organizations. Indeed, to query data efficiently, it is necessary for those different organizations to agree on the vocabulary used, which may not necessarily be possible. Introducing a logical layer to simplify this harmonization is often more practical than convincing everyone to use a unique standard. This logical layer is usually called an *ontology*, giving rise to the problem of *ontology-based query answering (OBQA)*.

There are many ways to express ontologies, the most usual being description logics. However, description logics are limited to unary and binary predicates, and can only encode knowledge in tree-like structures, which restricts their expressivity. Another language that is able to express more complex structures is Datalog, but it is limited in the sense that it cannot create any new individual: all the knowledge inferred is about already known individuals. Thus, to meet both those requirements, another framework has been studied, which is known as *Existential rules*. Existential rules extend Datalog rules with the ability to create new individuals. More precisely, an existential rule is a first order formula of the form $body \rightarrow head$, where $body$ and $head$ are both conjunctions of atoms, and variables appearing only in the head are existentially quantified.

While existential rules have both the ability to express complex structures and to infer the existence of new individuals, this expressivity comes with a cost: reasoning becomes semi-decidable. The usual reasoning task with existential rules, *boolean conjunctive query (BCQ) entailment*, can be summarized as the following: given a database and an ontology expressed as a set of existential rules, is a given query a logical consequence of the ontology if we take into account the facts that compose the database? In order to solve it, two main approaches have been developed.

The first approach to solve BCQ entailment is the *chase*. The chase is a family of algorithms that works by materializing new facts in the database using the rules in the ontology, until a fixpoint is reached. Then to solve BCQ entailment, it is enough to evaluate the query on the saturated database using traditional query evaluation algorithms. While this approach enables efficient query answering (through the use of standard database query evaluation over a query that stayed the same), it is not always feasible to modify the database,

either for security reasons, or because of size constraints.

The other standard way to solve BCQ entailment is through *query rewriting*. The idea is that instead of adding inferred knowledge to the database, one can modify the query to take into account the dependencies in the rule set, and then evaluate this query over the initial database. When starting from BCQs, using standard techniques, one ends up with a union of conjunctive queries (UCQ) as a rewriting, though other target languages have been studied. In this report, we focus particularly on Datalog rewritings, as it is a very natural restriction of existential rules, that ensures termination of the chase (as Datalog rules cannot create new individuals). Let us define these two kinds of rewritings. Take note however that to end up with Datalog rewritings, one not only modifies the query, but also the rule set.

Definition 1. *The couple (\mathcal{R}', q') is a rewriting of a rule set \mathcal{R} and a query q if $(\mathcal{R}, D) \models q$ if and only if $(\mathcal{R}', D) \models q'$ for all databases D . An UCQ rewriting of a query q w.r.t. a rule set \mathcal{R} is a rewriting (\emptyset, q') where q' is a union of conjunctive queries. A Datalog rewriting of a rule set \mathcal{R} and a query q is a rewriting (\mathcal{R}', q') where all the rules in \mathcal{R}' are Datalog. A couple rule set-query is UCQ-expressible (resp. Datalog-expressible) if it admits a UCQ rewriting (resp. a Datalog rewriting).*

An interesting characteristic of UCQ and Datalog rewritings is that query answering using these rewritings is decidable. Indeed, a UCQ rewriting does not feature a rule set so usual database theory algorithms work, and for Datalog rewritings, the chase terminates. Thus, from any rule set and query, if one can compute a UCQ or a Datalog rewriting, then one can solve query entailment in finite time. Thus, the question of finding these rewritings is of paramount interest, and gives rise to the following problems.

QUERY ANSWERING OVER UCQ-EXPRESSIBLE RULE SETS

Input: A database D , a rule set \mathcal{R} and a query q that is UCQ-expressible w.r.t. \mathcal{R}

Question: Do we have $(\mathcal{R}, D) \models q$?

Theorem 2. [3] *Query answering over UCQ-expressible rule sets is decidable.*

This means that we have an effective algorithm to compute UCQ rewritings. Though, the situation is very different for Datalog rewritings.

QUERY ANSWERING OVER DATALOG-EXPRESSIBLE RULE SETS

Input: A database D , a couple rule set-query (\mathcal{R}, q) that is \mathcal{R} -Datalog-expressible

Question: Do we have $(\mathcal{R}, D) \models q$?

Theorem 3. *Query answering over Datalog-expressible rule sets is undecidable.*

Our contribution is proving Theorem 3. To do so, we present a reduction from the halting problem to the problem of query answering over Datalog-expressible rule sets.

Chapter 2

Preliminaries

2.1 Syntax

We work in a first-order setting limited to *predicates*, *constants*, and *variables*. A *term* is a variable or a constant. To each predicate we associate a natural integer, called the *arity* of the predicate. A *P-atom* (or just an atom) is a first order formula of the form $P(t_1, \dots, t_n)$ where P is a predicate of arity n and t_1, \dots, t_n are some terms.

In the following, we may also group variables into tuples (which are seen as sets when convenient). Let $\vec{x}_1, \dots, \vec{x}_n$ be pairwise disjoint tuples of variables. Then, $X[\vec{x}_1, \dots, \vec{x}_n]$ denotes a conjunction of atoms whose variables are among $\vec{x}_1 \cup \dots \cup \vec{x}_n$.

Definition 4. An (existential) rule R is a first-order formula of the following form:

$$\forall \vec{x} \forall \vec{y}. B[\vec{x}, \vec{y}] \rightarrow \exists \vec{z}. H[\vec{x}, \vec{z}]$$

where $B[\vec{x}, \vec{y}]$ and $H[\vec{x}, \vec{z}]$ are conjunctions of atoms called the body and the head of the rule, respectively. The set \vec{x} , which is shared between the body and the head, is called the frontier of the rule and is denoted by frR . An existential rule where \vec{z} is empty is a Datalog rule.

For the sake of readability, we omit universal quantifiers when writing existential rules, and we denote conjunctions of atoms with commas instead of wedges. For instance, the rule $\forall x, y. A(x) \wedge P(x, y) \rightarrow \exists z. P(y, z)$ will be written $A(x), P(x, y) \rightarrow \exists z. P(y, z)$.

An *instance* is an existentially closed conjunction of atoms, and a *database* is an instance featuring only constants. It is usual to identify an instance with the corresponding set of atoms. For instance, the instance $\exists x, y, z. A(x, y) \wedge B(x, z, y)$ is identified to the set $\{A(x, y), B(x, z, y)\}$. A *knowledge base* is a tuple $\mathcal{K} = (\mathcal{R}, D)$ where \mathcal{R} is a rule set and D is an instance.

Definition 5. A conjunctive query is a conjunction of atoms. We call the free variables of a conjunctive query its answer variables. A Boolean conjunctive query is a conjunctive query that has no answer variable.

Note that by this definition, a Boolean conjunctive query (BCQ) is an instance.

In this setting, homomorphisms are often denoted as sets of individual variable mappings. For example, the homomorphism from $\{x, y, z\}$ to $\{a, b\}$ such that $\sigma(x) = a$, $\sigma(y) = a$ and $\sigma(z) = b$ can be denoted by $\{x \mapsto a, y \mapsto a, z \mapsto b\}$.

Definition 6. A retraction r from a set of atoms D to $D' \subseteq D$ is a homomorphism from D to D' such that $r|_{D'} = id_{D'}$ and $r(D) = D'$.

2.2 Semantics

We rely on classical first-order logic semantics. Indeed, the following definitions can be derived from the classical ones by considering the knowledge base (\mathcal{R}, D) as the theory $\mathcal{R} \cup D$ (recall that an instance is a formula).

An instance D *satisfies* (or *models*) a rule $R = B \rightarrow H$, denoted with $D \models R$, if for every homomorphism π from B to D , there is an extension $\hat{\pi}$ of π to $B \cup H$ such that $\hat{\pi}(H) \subseteq D$. A *model* \mathcal{M} of a knowledge base (\mathcal{R}, D) is a set of atoms such that $D \subseteq \mathcal{M}$ that satisfies all the rules in \mathcal{R} . A model \mathcal{M} of a knowledge base \mathcal{K} is *universal* if for every model \mathcal{M}' of \mathcal{K} , there is a homomorphism from \mathcal{M} to \mathcal{M}' . An instance D_1 *entails* an instance (or a BCQ) D_2 if there is a homomorphism from D_2 to D_1 . A knowledge base (\mathcal{R}, D_1) *entails* an instance (or a BCQ) D_2 if every model of (\mathcal{R}, D_1) entails D_2 . In both cases, A entails B is denoted by $A \models B$.

Note that by this definition, if \mathcal{U} is a universal model of a knowledge base \mathcal{K} , then for every model \mathcal{M} of \mathcal{K} , $\mathcal{M} \models \mathcal{U}$. Thus, given a BCQ q , if $\mathcal{U} \models q$, then for any model \mathcal{M} of \mathcal{K} , we have $\mathcal{M} \models q$, and if $\mathcal{U} \not\models q$, then there is a model of \mathcal{K} that does not model q . Thus, $\mathcal{K} \models q$ if and only if $\mathcal{U} \models q$. This motivates the study of an algorithm to find universal models in order to solve BCQ entailment.

2.3 The chase

The chase is a family of algorithms mainly used to solve BCQ entailment, through repeated rule applications until a fixpoint is reached. The chase is indeed sound and complete [2], but since the problem of BCQ entailment is semi-decidable [1], it does not necessarily terminate. To define it properly, we first introduce the notion of a trigger.

Definition 7. A trigger t for an instance D is a pair (R, π) consisting of a rule $R = B \rightarrow H$ and a homomorphism from B to D . We then define π^R as the extension of π that maps every existential variable z in R to the fresh variable z_t that is unique for z and t . A trigger $t = (R, \pi)$ is Datalog if R is Datalog, and existential otherwise. The trigger t is applicable on D if $\pi^R(H) \not\subseteq D$. We denote $\pi(B)$ and $\pi^R(H)$ by $support(t)$ and $output(t)$, and refer to them as the support and the output of t , respectively.

Now, that we have a basic notion of applicability, we can define derivations, before refining the notion of applicability to exhibit different chase variants.

Definition 8. Let $\mathcal{K} = (\mathcal{R}, D)$ be a knowledge base. A derivation from \mathcal{K} is a (possibly infinite) sequence $\mathcal{D} = D_0, t_1, D_1, \dots$ where $D_0 = D$, and for each $i > 0$, t_i is a trigger for the instance D_{i-1} and $D_i = D_{i-1} \cup output(t_i)$. The result of a derivation is the set $res(\mathcal{D}) = \bigcup_i D_i$.

Note that $\text{res}(\mathcal{D})$ is well-defined for every derivation $\mathcal{D} = D_0, t_1, D_1, \dots$ (finite or not) because $\{D_i\}_{i>0}$ is an increasing sequence for \subset for every derivation.

In the following it will be useful to distinguish between two variants of the chase, namely the oblivious chase (\mathbb{O}) and the restricted chase (\mathbb{R}). They differ from each other in how they constrain trigger applicability, as specified in the next definition.

Definition 9. Let $\mathcal{K} = (\mathcal{R}, D)$ be a knowledge base, $\mathcal{D} = D_0, t_1, D_1, \dots$ be a derivation from \mathcal{K} , and $t = (R, \pi)$ a trigger applicable on $\text{res}(\mathcal{D})$. Then t is:

- \mathbb{O} -applicable on \mathcal{D} if there is no i such that $t = t_i$.
- \mathbb{R} -applicable on \mathcal{D} if there is no retraction from $\text{res}(\mathcal{D}) \cup \text{output}(t)$ to $\text{res}(\mathcal{D})$.

An \mathbb{X} -derivation (with $\mathbb{X} \in \{\mathbb{O}, \mathbb{R}\}$) is a derivation $\mathcal{D} = D_0, t_1, D_1, \dots$ for which every trigger t_i is \mathbb{X} -applicable on D_{i-1} . A Datalog-first \mathbb{X} -derivation is an \mathbb{X} -derivation $\mathcal{D} = D_0, t_1, D_1, \dots$ such that for all i , if t_i is existential, then there is no Datalog trigger \mathbb{X} -applicable on D_{i-1} .

Datalog-first derivations consist is a succession of closure under Datalog rules and applications of a single existential trigger. Thus, in the following, we will write Datalog-first derivations from (\mathcal{R}, D) as $\mathcal{D} = \mathcal{D}_0, t_1, \mathcal{D}_1, \dots$, where for all i , t_i is an existential trigger and \mathcal{D}_i is a derivation from $(\mathcal{R}, \text{res}(t_i))$ that features only Datalog triggers.

Theorem 10. If \mathcal{D} is a chase sequence from a knowledge base (\mathcal{R}, D) , then $(\mathcal{R}, D) \models q$ if and only if $\text{res}(\mathcal{D}) \models q$, for all queries q .

2.4 Some computability theory

Definition 11. A Turing machine M is a tuple (Q, q_0, q_f, δ) where Q is a set of states, q_0 and q_f are two special states called the initial and the final state, respectively, and $\delta : (Q \setminus \{q_f\}) \times \{0, 1, \mathbb{B}\} \rightarrow (Q \setminus \{q_0\}) \times \{0, 1, \mathbb{B}\} \times \{\leftarrow, \rightarrow\}$ is a function called the transition function.

A configuration of a Turing machine is a triple (q, h, T) where q is a state, h is an integer called the position of the head, and $T : \{0, 1, \dots, n\} \rightarrow \{0, 1, \mathbb{B}\}$ is the tape, that is such that $T(n) = \mathbb{B}$, and where n is an integer. We then say that the tape is of size n . The start configuration C_0 of M on a word w (on the alphabet $\{0, 1\}$) is $(q_0, 0, T_0)$ where T_0 is a tape of size $|w|$ such that for all i between 0 and $|w|$, we have $T_0(i) = w[i]$.

The run of M on a word w is a sequence C_0, C_1, \dots of configurations such that C_0 is the start configuration of M on w , and for all t , if $C_t = (q_t, h_t, T_t)$, then T_t is of size t and:

- If $q_t = q_f$, then C_t is the last element of the run. In this case, we say M halts on w at step t .
- If $\delta(q_t, T_t(h_t)) = (q, c, \leftarrow)$ and $h_t = 0$, then we say M halts on w at step t .
- If $\delta(q_t, T_t(h_t)) = (q, c, \leftarrow)$ and $h_t > 0$, then $C_{t+1} = (q, h_t - 1, T_{t+1})$.

- If $\delta(q_t, T_t(h_t)) = (q, c, \rightarrow)$, then $C_{t+1} = (q, h_t + 1, T_{t+1})$.

where $T_{t+1}(h_t) = c$, $T_{t+1}(t + 1) = \mathbf{B}$, and $T_{t+1}(i) = T_t(i)$ for all $i \neq h_t$ between 0 and t .

In the following, we denote the unique run from M on the empty word as $r^M = C_0^M, C_1^M, \dots$, with $C_t^M = (q_t, h_t, T_t)$ for all $t \geq 0$, and we write $M(i, t) = T_t(i)$, so $M(i, t)$ is the content of position i of M 's tape at step t .

If M halts on the empty word, we denote with $\text{time}_\varepsilon(M)$ the index at which it does. Otherwise, we define $\text{time}_\varepsilon(M) = \infty$.

Note that this definition of a Turing machine is a bit unusual, to make the representation using existential rules and the proofs a bit simpler. First, we do not consider acceptance or rejection, but only halting and not halting. In addition, instead of the usual “infinite tape” that Turing machines feature, we use a tape of unbounded size, but which is finite at each step. In addition, the size of the tape is not its real size: as a tape T of size n is a function from $\{0, 1, \dots, n\}$ to our alphabet, it actually features $n + 1$ cases. This can be explained by the fact that we ask for $T(n)$ to be \mathbf{B} , which ensures that the head of the Turing machine always has something to read (otherwise, when starting from the empty word ε , the head would not be able to the tape in the starting configuration). Though, the size of the word represented on the tape is indeed the size of the tape as we define it (the empty word has size 0).

Chapter 3

OBQA over Datalog-expressible rule sets is undecidable

We present a reduction that takes a Turing machine M and produces a rule set \mathcal{R}_M such that

1. The Turing machine M halts on the empty word if and only if $(\mathcal{R}_M, \emptyset) \models \text{Goal}$.
2. The couple rule set-query $(\mathcal{R}_M, \text{Goal})$ is Datalog-expressible.

Provided we can construct such a reduction, if we assume, by the absurd, that there is an algorithm that decides entailment for Datalog-expressible rule sets, then we can decide if M halts, which contradicts the undecidability of the halting problem. Thus, this is enough to prove Theorem 3. We then exhibit such a rule set, before proving it enjoys both these properties.

3.1 The rule set

Definition 12. For a TM M , let \mathcal{R}_M be the rule set containing the following rules:

$$\begin{aligned} & \rightarrow \exists w_0. \mathbf{q_0}(w_0, w_0) && (R_0) \\ \mathbf{q_0}(w_0, w_0) & \rightarrow \exists w_1. \mathbf{N}(w_0, w_1) && (R_1) \\ \mathbf{q_0}(w_0, w_0), \mathbf{N}^+(w_0, w_i) & \rightarrow \exists w_{i+1}. \mathbf{N}(w_i, w_{i+1}) && (R_{+1}) \\ \mathbf{N}(w_i, w_{i+1}) & \rightarrow \mathbf{N}^+(w_i, w_{i+1}) && (R_{\mathbf{N}^+}) \\ \mathbf{N}(w_i, w_{i+1}), \mathbf{N}^+(w_{i+1}, w_k) & \rightarrow \mathbf{N}^+(w_i, w_k) && (R_{\mathbf{N}^+}^{trans}) \\ \mathbf{q_0}(w_0, w_0) & \rightarrow \mathbf{B}(w_0, w_0) && (R_{\mathbf{B}}^0) \\ \mathbf{q_0}(w_0, w_0), \mathbf{N}^+(w_0, w_i) & \rightarrow \mathbf{B}(w_i, w_i) && (R_{\mathbf{B}}) \\ \mathbf{q_f}(x, y) & \rightarrow \text{Goal} && (R_{\text{Goal}}) \end{aligned}$$

For every $q, q' \in Q$, $c, c' \in \{0, 1, \mathbf{B}\}$ such that $\delta(q, c) = (q', c', \leftarrow)$,

$$\begin{aligned} N(p_1, p_2), N(t_1, t_2), \mathbf{q}(p_2, t_1), \mathbf{c}(p_2, t_1), N^+(p_2, t_2) &\rightarrow \mathbf{q}'(p_1, t_2), \mathbf{c}'(p_2, t_2) & (R_{q,c}^{\leftarrow}) \\ \mathbf{q}_0(w_0, w_0), \mathbf{q}(w_0, t), \mathbf{c}(w_0, t) &\rightarrow \mathbf{Goal} & (R_{q,c}^0) \end{aligned}$$

For every $q, q' \in Q$, $c, c' \in \{0, 1, \mathbf{B}\}$ such that $\delta(q, c) = (q', c', \rightarrow)$,

$$\begin{aligned} N(p_1, p_2), N(t_1, t_2), \mathbf{q}(p_1, t_1), \mathbf{c}(p_1, t_1), N^+(p_2, t_2) &\rightarrow \mathbf{q}'(p_2, t_2), \mathbf{c}'(p_1, t_2) & (R_{q,c}^{\rightarrow}) \\ N(w_1, w_2), \mathbf{q}(w_1, w_1), \mathbf{c}(w_1, w_1) &\rightarrow \mathbf{q}'(w_2, w_2), \mathbf{c}'(w_1, w_2) & (R_{q,c}^{max}) \end{aligned}$$

For every $q \in Q$ and $c \in \{0, 1, \mathbf{B}\}$,

$$\begin{aligned} N^+(p_i, p_k), N(t_1, t_2), \mathbf{q}(p_k, t_1), \mathbf{c}(p_i, t_2) &\rightarrow \mathbf{c}(p_i, t_2) & (R_{q,c}^L) \\ N^+(p_i, p_k), N(t_1, t_2), \mathbf{q}(p_i, t_1), \mathbf{c}(p_k, t_2) &\rightarrow \mathbf{c}(p_k, t_2) & (R_{q,c}^R) \end{aligned}$$

And the rules:

$$\begin{aligned} N(x, y), N(x, y') &\rightarrow \mathbf{Eq}(y, y') & (R_{func}^+) \\ N(x, y), N(x', y) &\rightarrow \mathbf{Eq}(x, x') & (R_{func}^-) \\ \mathbf{q}_0(x, y), \mathbf{q}_0(z, t) &\rightarrow \mathbf{Eq}(x, y), \mathbf{Eq}(y, z), \mathbf{Eq}(z, t) & (R_{q_0}) \\ N^+(x, x) &\rightarrow \mathbf{Fail} & (R_{N^+}^{cycle}) \\ N(w_{-1}, w_0), \mathbf{q}_0(w_0, w_0) &\rightarrow \mathbf{Fail} & (R_{-1}) \\ N^+(x, y), \mathbf{q}_1(x, t), \mathbf{q}_2(y, t) &\rightarrow \mathbf{Fail} \text{ for every } q_1, q_2 \in Q & (R_{q_1, q_2}^{step}) \\ \mathbf{q}_1(x, t), \mathbf{q}_2(x, t) &\rightarrow \mathbf{Fail} \text{ for every } q_1 \neq q_2 \text{ in } Q & (R_{q_1, q_2}^{cell}) \\ \mathbf{c}_1(x, t), \mathbf{c}_2(x, t) &\rightarrow \mathbf{Fail} \text{ for every } c_1 \neq c_2 \text{ in } \{0, 1, \mathbf{B}\} & (R_{c_1, c_2}^{cell}) \\ N^+(t, x), \mathbf{c}(x, t) &\rightarrow \mathbf{Fail} \text{ for every } c \in \{0, 1, \mathbf{B}\} & (R_c^{OoB}) \\ \mathbf{Fail} &\rightarrow \mathbf{Goal} & (R_{Fail}) \\ \mathbf{Eq}(x, y), \mathbf{Eq}(y, z) &\rightarrow \mathbf{Eq}(x, z) & (R_{Eq}^{trans}) \\ \mathbf{Eq}(x, y) &\rightarrow \mathbf{Eq}(y, x) & (R_{Eq}^{sym}) \end{aligned}$$

For all predicates P and integers i between 1 and $ar(P)$,

$$P(x_1, \dots, x_{ar(P)}), \mathbf{Eq}(x_i, y_i) \rightarrow P(x_1, \dots, y_i, \dots, x_{ar(P)}) \quad (R_{P,i}^{repl})$$

3.2 \mathcal{R}_M simulates M

This section is dedicated to proving the following theorem.

Theorem 13. *The Turing machine M halts on the empty word if and only if $(\mathcal{R}_M, \emptyset) \models \mathbf{Goal}$.*

Before proving this theorem, we explain the relevant parts of the rule set. For this rule set, we use an encoding derived from the usual grid representation of Turing machines. In the grid representation of Turing machines, one would

usually consider a two-way infinite grid (isomorphic to $\mathbb{N} \times \mathbb{N}$), on which the position (i, j) would represent the i -th cell of the machine's tape at step j . Then, to encode this representation in existential rules, one would have an individual for each cell and two binary predicates representing the relations “ x is right below y ” and “ x is right at the left of y ”. Finally, one would have a unary predicate for each state and for each tape symbol, and $\mathbf{q}(x)$ would mean that the head is on the cell represented by x and the Turing machine is at step q at this time, and $\mathbf{c}(x)$ means that x contains c .

Here, we use the same representation, but we encode it differently: instead of using unary predicates for the content of the tapes, we use binary predicates. This means that we only need individuals representing \mathbb{N} and not the full $\mathbb{N} \times \mathbb{N}$. In addition, we only use one successor predicate instead of two. In our case, this successor predicate is \mathbf{N} , and the rules that generate this encoding of \mathbb{N} are rules R_0 , R_1 and R_{+1} . Then, the predicate \mathbf{N}^+ is \mathbf{N} 's transitive closure, and is computed by rules $R_{\mathbf{N}^+}$ and $R_{\mathbf{N}^+}^{trans}$. In the following, we will refer to the elements of this representation of \mathbb{N} by $(z_n)_{n \in \mathbb{N}}$: z_0 is the element that does not have a predecessor through \mathbf{N} , and z_{i+1} is z_i 's successor, for all i .

Then, as we said, we model Turing machine simulation using a binary predicate for each state of the Turing machine and for each character that can appear on the tape. In more detail, we will show in Proposition 16 that if $\mathbf{q}(z_i, z_t)$ is in $chase(\mathcal{R}_M, \emptyset)$ for some state q , then the Turing machine is in state q at step t , and its head is on position i . Similarly, if $\mathbf{c}(z_i, z_t)$ is in $chase(\mathcal{R}_M, \emptyset)$, then the character at position i of the tape at step t is c . Rules R_0 , $R_{\mathbf{B}}^0$ and $R_{\mathbf{B}}$ initialize this calculation, rules $R_{q,c}^{\leftarrow}$ and $R_{q,c}^{\rightarrow}$ move the head and change the character, and rules $R_{q,c}^L$ and $R_{q,c}^R$ propagate the symbols that do not change. Finally, rules $R_{\mathbf{Goal}}^0$ and $R_{q,c}^0$ generate \mathbf{Goal} whenever the Turing machine encounters q_f or tries to move its head left when it is on position 0, that is to say whenever it halts. The remaining rules will not do much on the empty database, and will be useful later to ensure that the rule set is Datalog expressible. We will thus explain them later. We first show a quick lemma that will be useful in the proof of Theorem 13.

Lemma 14. *Let D be an instance and t a term in D . There is no trigger (R, π) that is applicable on D and such that $\mathbf{Eq}(t, t) \in \pi(\text{body}(R))$.*

Proof. Suppose for a contradiction that there is such a trigger. Then, since the only rules that feature an \mathbf{Eq} -atom in their body are rules $R_{\mathbf{Eq}}^{trans}$, $R_{\mathbf{Eq}}^{sym}$ and $R_{P,i}^{repl}$, the rule R must be one of these three.

- If $R = R_{\mathbf{Eq}}^{trans}$, then two cases are possible. Either $\pi(x) = \pi(y) = t$ and $\pi(z) = u$ such that $\mathbf{Eq}(t, u) \in D$, and (R, π) yields $\mathbf{Eq}(t, u)$ or $\pi(y) = \pi(z) = t$ and $\pi(x) = u$ with $\mathbf{Eq}(u, t) \in D$ and (R, π) yields $\mathbf{Eq}(u, t)$. In either case, $output(R, \pi) \subseteq D$ so the trigger is not applicable.
- If $R = R_{\mathbf{Eq}}^{sym}$, then $\pi(x) = \pi(y) = t$, and (R, π) yields $\mathbf{Eq}(t, t)$, which is by hypothesis in $D \cup \mathbf{Eq}(t, t)$, so this trigger is not applicable either.
- If $R = R_{P,i}^{repl}$, then there is an atom $P(t_1, \dots, t_{i-1}, t, t_{i+1}, \dots, t_n)$ (with $n = ar(P)$). Then, again, (R, π) yields $P(t_1, \dots, t_{i-1}, t, t_{i+1}, \dots, t_n)$, which is in $D \cup \{\mathbf{Eq}(t, t)\}$ by hypothesis, so (R, π) is not applicable in this case either. \square

Then, to show Theorem 13, we will specify the result of Datalog-first derivations. We first define a few sets that will be useful to express this result.

Definition 15. For all n , we define the following sets

$$\begin{aligned}
D_n^N &= \{\mathbf{N}(z_i, z_{i+1}) \mid 0 \leq i < n\} \\
D_n^{N^+} &= \{\mathbf{N}^+(z_i, z_j) \mid 0 \leq i < j \leq n\} \\
D_n^Q &= \{\mathbf{qt}(z_{h_t}, z_t) \mid 0 \leq t \leq \min(n, \text{time}_\varepsilon(M))\} \\
D_n^\Sigma &= \{\mathbf{c}(z_i, z_t) \mid 0 \leq t \leq \min(n, \text{time}_\varepsilon(M)), 0 \leq i \leq t, M(i, t) = c\} \\
D_n^{\text{Eq}} &= \{\mathbf{Eq}(z_i, z_i) \mid 0 \leq i \leq n\} \\
D_n^{\text{Goal}} &= \{\mathbf{Goal}\} \text{ if } n \geq \text{time}_\varepsilon(M) \text{ and } \emptyset \text{ otherwise.}
\end{aligned}$$

We can now formulate the following proposition.

Proposition 16. For a Datalog-first derivation $\mathcal{D} = \mathcal{D}_{\text{init}}, t_0, \mathcal{D}_0, t_1, \mathcal{D}_1, \dots$ from $(\mathcal{R}_M, \emptyset)$, we have that:

- $t_0 = ((R_0), \emptyset)$ and $t_1 = (R_1, \{w_0 \mapsto z_0\})$
- For all $n > 1$, $t_n = (R_{+1}, \{w_0 \mapsto z_0, w_i \mapsto z_{n-1}\})$
- For all $n \in \mathbb{N}$, $\text{res}(\mathcal{D}'_n) = D_n^N \cup D_n^{N^+} \cup D_n^Q \cup D_n^\Sigma \cup D_n^{\text{Eq}} \cup D_n^{\text{Goal}}$.

Proof. We proceed by induction on n .

At step $n = 0$, on the database \emptyset , only rules with an empty body are applicable. As R_0 is the only rule with an empty body, it is applicable, so $\mathcal{D}_{\text{init}}$ is empty and $t_0 = (R_0, \emptyset)$, with $\text{res}(t_0) = \{\mathbf{q_0}(z_0, z_0)\}$. Then, rules $R_{\mathbf{B}}^0$ and R_{q_0} are applicable and yield $\mathbf{B}(z_0, z_0)$ and $\mathbf{Eq}(z_0, z_0)$, respectively. Finally, if $\delta(q_0, \mathbf{B}) = (q', c, \leftarrow)$ for some q', c' , then rule $R_{q_0, \mathbf{B}}^0$ is applicable, and if $q_0 = q_f$, R_{Goal} is applicable. In both these cases, the application yields \mathbf{Goal} , and by definition of a run, $\text{time}_\varepsilon(M) = 0$. We then show that $\text{res}(\mathcal{D}_0) = D_0^N \cup D_0^{N^+} \cup D_0^Q \cup D_0^\Sigma \cup D_0^{\text{Eq}} \cup D_0^{\text{Goal}}$ by double inclusion. For the converse inclusion, note that D_0^N and $D_0^{N^+}$ are empty, and by definition of a starting configuration, we have $D_0^Q = \{\mathbf{q_0}(z_0, z_0)\}$ and $D_0^\Sigma = \{\mathbf{B}(z_0, z_0)\}$. In addition, we also have $D_0^{\text{Eq}} = \{\mathbf{Eq}(z_0, z_0)\}$. Then, $D_0^{\text{Goal}} = \{\mathbf{Goal}\}$ if $\text{time}_\varepsilon(M) = 0$ (case in which $\mathbf{Goal} \in \text{res}(\mathcal{D}_0)$), and $D_0^{\text{Goal}} = \emptyset$ otherwise. Thus, the converse inclusion holds.

Then, for the direct inclusion, we show that no Datalog rule is applicable on $D_0^N \cup D_0^{N^+} \cup D_0^Q \cup D_0^\Sigma \cup D_0^{\text{Goal}} = \{\mathbf{q_0}(z_0, z_0), \mathbf{B}(z_0, z_0)\} \cup D_0^{\text{Goal}}$. First notice that all Datalog rules except $R_{\text{Goal}}, R_{q, c}^0, R_{q_0}, R_{q_1, q_2}^{\text{cell}}, R_{c_1, c_2}^{\text{cell}}, R_{\text{Fail}}$ and the axiomatization of equality have an N -atom or an N^+ -atom in their body, so their bodies cannot be mapped to $\text{res}(t_0)$. In addition, the body of rules R_{Goal} and $R_{q, c}^0$ can only be mapped into $\{\mathbf{q_0}(z_0, z_0), \mathbf{B}(z_0, z_0)\} \cup D_0^{\text{Goal}}$ if $\text{time}_\varepsilon(M) = 0$. Indeed, if rule R_{Goal} 's body can be mapped, then there is $q_f(x, y)$ in the database, so $q_0 = q_f$, and if $R_{q, c}^0$'s body can be mapped, then $q = q_0, c = \mathbf{B}$ and $\delta(q_0, \mathbf{B}) = (q', c', \leftarrow)$. Thus, in both cases, $\text{time}_\varepsilon(M) = 0$, so $\mathbf{Goal} \in D_0^{\text{Goal}} \subseteq \text{res}(\mathcal{D}_0)$, and these rules are not applicable.

Regarding rule R_{q_0} , as it has already been applied with its body mapped to the only $\mathbf{q_0}$ -atom, it is not applicable either. Rule $R_{q_1, q_2}^{\text{cell}}$ (resp. $R_{c_1, c_2}^{\text{cell}}$) has two atoms that use distinct state (resp. character) predicates in its body. As there is

only one state (resp. character) atom in $res(\mathcal{D}'_0) \cup D_0^{Eq}$, rule R_{q_1, q_2}^{cell} (resp. R_{c_1, c_2}^{cell}) is not applicable. As $Fail \notin res(\mathcal{D}'_0)$, rule R_{Fail} is not applicable either. Finally, if $t = (R, \pi)$ is a trigger such that $R \in \{R_{Eq}^{trans}, R_{Eq}^{sym}, R_{P, i}^{repl}\}$, then there must be an Eq-atom in $\pi(body(R))$. As $Eq(z_0, z_0)$ is the only Eq-atom in $res(\mathcal{D}'_0) \cup D_0^{Eq}$, it must be that $Eq(z_0, z_0) \in \pi(body(R))$, which by Lemma 14 means that t is not applicable. Thus, the equality holds in the base case.

Before proving the inductive step, we prove that $t_1 = (R_1), \{w_0 \mapsto z_0\}$. Since $res(\mathcal{D}_0) = \{\mathbf{q}_0(z_0, z_0), \mathbf{B}(z_0, z_0)\}$, rule R_{+1} is not applicable (it has an \mathbf{N}^+ -atom in its body). Since rule R_0 has already been applied with the only homomorphism possible, the only applicable rule is R_1 , so $t_1 = (R_1), \{w_0 \mapsto z_0\}$.

Now, for the inductive step, assume that $n > 0$ and that the lemma is true for all $i < n$. We first prove that if $n > 1$, $t_n = (R_{+1}), \{w_i \mapsto z_{n-1}, w_0 \mapsto z_0\}$. By definition of a Turing machine, q_0 can only appear at step and position 0, so by induction hypothesis, the only term x such that $\mathbf{q}_0(x, x) \in res(\mathcal{D}'_{n-1})$ is z_0 . Thus, all the homomorphisms from the body of rule R_{+1} to $res(\mathcal{D}'_{n-1})$ are of the form $\{w_i \mapsto z_j, w_0 \mapsto z_0\}$ with $0 \leq j \leq n-1$. Again by induction hypothesis, all triggers $(R_{+1}), \{x_1 \mapsto z_j, w_0 \mapsto z_0\}$ for $2 \leq j \leq n-2$ have been applied at earlier steps, so the only applicable existential trigger is $(R_{+1}), \{w_i \mapsto z_{n-1}, w_0 \mapsto z_0\}$.

We now come back to the more general case $n > 0$, and we prove that $D_n^N \cup D_n^{N^+} \cup D_n^Q \cup D_n^\Sigma \cup D_n^{Eq} \cup D_n^{Goal} \subseteq res(\mathcal{D}'_n)$.

$D_n^N \subseteq res(\mathcal{D}_n)$: By induction hypothesis, $D_{n-1}^N \subseteq res(\mathcal{D}'_{n-1}) \subseteq res(\mathcal{D}'_n)$. Then, if $n = 1$, $output(t_1) = \{N(z_0, z_1)\} = \{N(z_{n-1}, z_n)\}$, and if $n > 1$, $output(t_n) = \{N(z_{n-1}, z_n)\}$, so in either case, $D_n^N \subseteq res(\mathcal{D}'_n)$.

$D_n^{N^+} \subseteq res(\mathcal{D}_n)$ By induction hypothesis, $D_{n-1}^{N^+} \subseteq res(\mathcal{D}'_{n-1}) \subseteq res(\mathcal{D}'_n)$, so it is enough to show that for all i such that $0 \leq i < n$, we have $\mathbf{N}^+(z_i, z_n) \in res(\mathcal{D}'_n)$, or equivalently that for all i such that $1 \leq i \leq n$, we have $\mathbf{N}^+(z_{n-i}, z_n) \in res(\mathcal{D}'_n)$. We prove this last result by induction on i .

We first show that if $i = 1$, then $\mathbf{N}^+(z_{n-1}, z_n) \in res(\mathcal{D}'_n)$. As $D_n^N \subseteq res(\mathcal{D}'_n)$, we have $\mathbf{N}(z_{n-1}, z_n) \in res(\mathcal{D}'_n)$. Thus, since $res(\mathcal{D}'_n)$ satisfies rule R_{N^+} , the atom $\mathbf{N}^+(z_{n-1}, z_n)$ is in $res(\mathcal{D}'_n)$.

For the induction step, we show that if $1 < i \leq n$ and $\mathbf{N}^+(z_{n-(i-1)}, z_n) \in res(\mathcal{D}'_n)$, then $\mathbf{N}^+(z_{n-i}, z_n) \in res(\mathcal{D}'_n)$. Since $D_n^N \subseteq res(\mathcal{D}'_n)$, we have $\mathbf{N}(z_{n-1}, z_{n-i+1}) \in res(\mathcal{D}'_n)$. As $n - i + 1 = n - (i - 1)$, the atoms $\mathbf{N}(z_{n-1}, z_{n-i+1})$ and $\mathbf{N}^+(z_{n-i+1}, z_n)$ are both in $res(\mathcal{D}'_n)$. Thus, since $res(\mathcal{D}'_n)$ models rule $R_{N^+}^{trans}$, $\mathbf{N}^+(z_{n-(i-1)}, z_n) \in res(\mathcal{D}'_n)$.

$D_n^Q \subseteq res(\mathcal{D}_n)$: First, notice that if $n > time_\varepsilon(M)$, then $D_n^Q = D_{n-1}^Q$, so the induction hypothesis suffices. We thus only consider the case where $n \leq time_\varepsilon(M)$ in the following. We then want to show that $\mathbf{q}_n(z_{h_n}, z_n) \in res(\mathcal{D}'_n)$. As $q_{n-1} \neq q_f$, there are q_n, c, c', d such that $c = M(h_{n-1}, n-1)$ and $\delta(q_{n-1}, c) = (q_n, c', d)$. By induction hypothesis, both the atoms $\mathbf{q}_{n-1}(z_{h_{n-1}}, z_{n-1})$ and $\mathbf{c}(z_{h_{n-1}}, z_{n-1})$ are in $res(\mathcal{D}'_n)$. In addition, as $h_{n-1} \leq n-1 < n$, we also have $\mathbf{N}^+(z_{h_{n-1}}, z_n) \in res(\mathcal{D}'_n)$, and we showed earlier that $\mathbf{N}(z_{n-1}, z_n) \in res(\mathcal{D}'_n)$. Thus:

- If $d = \rightarrow$, as $h_{n-1} < n$, we have $\mathbf{N}(z_{h_{n-1}}, z_{h_{n-1}+1}) \in res(\mathcal{D}'_n)$. Thus,

$$t_{\rightarrow} = (R_{q_{n-1}, c}^{\rightarrow}, \{p_i \mapsto z_{h_{n-1}}, p_{i+1} \mapsto z_{h_{n-1}+1}, t_i \mapsto z_{n-1}, t_{i+1} \mapsto z_n\})$$

is applicable. In this case, by definition, $h_n = h_{n-1} + 1$, so this trigger yields the atoms $\mathbf{q}_n(z_{h_n}, z_n)$ and $\mathbf{c}'(z_{h_{n-1}}, z_n)$.

- If $d = \leftarrow$, since $n - 1 < \text{time}_\varepsilon(M)$, then $h_{n-1} > 0$, so we have $\mathbf{N}(z_{h_{n-1}-1}, z_{h_{n-1}}) \in \text{res}(\mathcal{D}'_n)$. Thus, the trigger

$$t_{\leftarrow} = (R_{q_{n-1}, c}^{\leftarrow}, \{p_{i-1} \mapsto z_{h_{n-1}-1}, p_i \mapsto z_{h_{n-1}}, t_i \mapsto z_{n-1}, t_{i+1} \mapsto z_n\})$$

is applicable. In this case, by definition, $h_n = h_{n-1} - 1$, so this trigger yields the atoms $\mathbf{q}_n(z_{h_n}, z_n)$ and $\mathbf{c}'(z_{h_{n-1}}, z_n)$.

Thus, in either cases, we have $\mathbf{q}_n(z_{h_n}, z_n)$ and $\mathbf{c}'(z_{h_{n-1}}, z_n)$ in $\text{res}(\mathcal{D}'_n)$, so $D_n^Q \subseteq \text{res}(\mathcal{D}'_n)$.

$D_n^\Sigma \subseteq \text{res}(\mathcal{D}_n)$: First, as before, if $n > \text{time}_\varepsilon(M)$, then $D_n^\Sigma = D_{n-1}^\Sigma$, so the induction hypothesis suffices. In the following, we restrict ourselves to the case $n \leq \text{time}_\varepsilon(M)$. We thus need to prove that $\mathbf{c}_i(z_i, z_n) \in \text{res}(\mathcal{D}'_n)$ for all i such that $0 \leq i \leq n$ and $M(i, n) = c_i$. There are four cases:

- If $i < h_{n-1}$, then $c_i = M(i, n) = M(i, n-1)$. By induction hypothesis, the atoms $\mathbf{c}_i(z_i, z_{n-1})$, $\mathbf{q}_{n-1}(z_{h_{n-1}}, z_{n-1})$ and $\mathbf{N}^+(z_i, z_{h_{n-1}})$ are in $\text{res}(\mathcal{D}'_n)$. In addition, since $D_n^N \subseteq \text{res}(\mathcal{D}'_n)$, we have $\mathbf{N}(z_{n-1}, z_n) \in \text{res}(\mathcal{D}'_n)$. Thus, as $\text{res}(\mathcal{D}'_n)$ satisfies rule $R_{q_{n-1}, c}^L$, we have $\mathbf{c}_i(z_i, z_n) \in \text{res}(\mathcal{D}'_n)$.
- If $i = h_{n-1}$, as seen in the previous case, we have $\mathbf{c}'(z_{h_{n-1}}, z_n)$, with \mathbf{c}' such that $\delta(q_{n-1}, M(h_{n-1}, n-1)) = (q', c', d)$. By definition of a run, $c' = M(h_{n-1}, n) = c_i$, which concludes this case.
- If $h_{n-1} < i < n$, then $c_i = M(i, n) = M(i, n-1)$. As $i \leq n-1$, by induction hypothesis, the atoms $\mathbf{c}_i(z_i, z_{n-1})$, $\mathbf{q}_{n-1}(z_{h_{n-1}}, z_{n-1})$ and $\mathbf{N}^+(z_{h_{n-1}}, z_i)$ are in $\text{res}(\mathcal{D}'_n)$. In addition, since $D_n^N \subseteq \text{res}(\mathcal{D}'_n)$, we have $\mathbf{N}(z_{n-1}, z_n) \in \text{res}(\mathcal{D}'_n)$. Thus, as $\text{res}(\mathcal{D}'_n)$ satisfies rule $R_{q_{n-1}, c}^R$, we have $\mathbf{c}_i(z_i, z_n) \in \text{res}(\mathcal{D}'_n)$.
- If $i = n$, a usual property of Turing machines states that $M(n, n) = \mathbf{B}$, as a Turing machine can only move one step at a step and edit the bit its head was on at a previous step. As $D_n^{N^+} \subseteq \text{res}(\mathcal{D}'_n)$, we have $\mathbf{N}^+(z_0, z_n) \in \text{res}(\mathcal{D}'_n)$. Thus, as $\text{res}(\mathcal{D}'_n)$ satisfies rule $R_{\mathbf{B}}$, the atom $\mathbf{B}(z_n, z_n)$ is in $\text{res}(\mathcal{D}'_n)$.

$D_n^{\text{Eq}} \subseteq \text{res}(\mathcal{D}_n)$: As we saw earlier, the atom $\mathbf{N}(z_{n-1}, z_n)$ is in $\text{res}(\mathcal{D}_n)$, so rule R_{func}^+ is applicable by mapping x to z_{n-1} , and y and y' to z_n . This application yields the atom $\text{Eq}(z_n, z_n)$, so $D_n^{\text{Eq}} \subseteq \text{res}(\mathcal{D}_n)$ (as we already had $D_{n-1}^{\text{Eq}} \subseteq \text{res}(\mathcal{D}_n)$ by induction hypothesis).

$D_n^{\text{Goal}} \subseteq \text{res}(\mathcal{D}_n)$: If $n < \text{time}_\varepsilon(M)$, the result follows from the fact that D_n^{Goal} is empty in this case. Otherwise, if $n \geq \text{time}_\varepsilon(M)$, then there are two cases:

- If $q_{\text{time}_\varepsilon(M)} = q_f$, since $D_n^Q \subseteq \text{res}(\mathcal{D}'_n)$, $\mathbf{q}_f(z_{h_{\text{time}_\varepsilon(M)}}, z_{\text{time}_\varepsilon(M)}) \in \text{res}(\mathcal{D}'_n)$. Thus, as $\text{res}(\mathcal{D}'_n)$ models rule R_{Goal} , $\text{Goal} \in \text{res}(\mathcal{D}'_n)$.

- Otherwise, we have $h_{time_\varepsilon(M)} = 0$ and

$$\delta(q_{time_\varepsilon(M)}, M(0, time_\varepsilon(M))) = (q, c, \leftarrow)$$

for some q and c . In this case, as $D_n^Q \cup D_n^\Sigma \subseteq res(\mathcal{D}'_n)$, the atoms $\mathbf{q}_{time_\varepsilon(M)}(z_0, z_{time_\varepsilon(M)})$ and $\mathbf{c}_0(z_0, z_{time_\varepsilon(M)})$ are in $res(\mathcal{D}'_n)$, with $c_0 = M(0, time_\varepsilon(M))$. Thus, as $res(\mathcal{D}'_n)$ satisfies rule $R_{q,c}^0$, $\mathbf{Goal} \in res(\mathcal{D}'_n)$.

Thus, in all cases, $D_n^{\mathbf{Goal}} \subseteq res(\mathcal{D}'_n)$.

This concludes the proof of the first inclusion. To show that the other inclusion holds, we will show that $D = D_n^{\mathbf{N}} \cup D_n^{\mathbf{N}^+} \cup D_n^Q \cup D_n^\Sigma \cup D_n^{\mathbf{Eq}} \cup D_n^{\mathbf{Goal}}$ satisfies all the Datalog rules in \mathcal{R}_M .

$R_{\mathbf{N}^+}$: Let i, j be such that $\mathbf{N}(z_i, z_j) \in D$. Then, by the definition of $D_n^{\mathbf{N}}$, we have $0 \leq i < n$ and $j = i + 1$, so $0 \leq i < j \leq n$. Thus, by the definition of $D_n^{\mathbf{N}^+}$, $\mathbf{N}^+(z_i, z_j) \in D$, so $R_{\mathbf{N}^+}$ is satisfied by D .

$R_{\mathbf{N}^+}^{\mathbf{trans}}$: Assume that $\mathbf{N}(z_i, z_j)$ and $\mathbf{N}^+(z_j, z_k)$ are in D , for some i, j, k . Then, by the definition of $D_n^{\mathbf{N}}$, we have $0 \leq i < n$ and $j = i + 1$, so $0 \leq i < j \leq n$. In addition, by the definition of $D_n^{\mathbf{N}^+}$, we have $0 \leq j < k \leq n$, so $0 \leq i < j < k \leq n$. In particular, we have $0 \leq i < k \leq n$, so $\mathbf{N}^+(z_i, z_j) \in D$, and $R_{\mathbf{N}^+}^{\mathbf{trans}}$ is satisfied by D .

$R_{\mathbf{B}}^0$: The only term x such that $\mathbf{q}_0(x, x) \in D$ is z_0 , and $\mathbf{B}(z_0, z_0) \in D$, so this rule is not applicable.

$R_{\mathbf{B}}$: As stated earlier, we have $M(i, t) = \mathbf{B}$ for all $t \leq i$. In particular, $M(i, i) = \mathbf{B}$ for all $i \leq n$. Thus, by the definition of D_n^Σ , for all $i \leq n$, we have $\mathbf{B}(z_i, z_i) \in D$, so rule $R_{\mathbf{B}}$ is satisfied.

$R_{\mathbf{Goal}}$: If $\mathbf{q}_f(z_i, z_t) \in D$, for some i, t , then by definition of D_n^Q , $q_t = q_f$, and thus $time_\varepsilon(M) = t \leq n$. As such, $D_n^{\mathbf{Goal}} = \{\mathbf{Goal}\}$, and rule $R_{\mathbf{Goal}}$ is satisfied.

$R_{q,c}^0$: As stated earlier, the only term x such that $\mathbf{q}_0(x, x) \in D$ is z_0 . Assume that there are q, c and t such that $\delta(q, c) = (q', c', \leftarrow)$ for some q', c' , and $\mathbf{q}(z_0, z_t)$ and $\mathbf{c}(z_0, z_t)$ are in D . Then, by the definitions of D_n^Q and D_n^Σ , $q_t = q$, $h_t = 0$ and $M(h_t, t) = M(0, t) = c$. Thus, by definition of a run, $time_\varepsilon(M) = t \leq n$, and $D_n^{\mathbf{Goal}} = \{\mathbf{Goal}\}$, so rule $R_{q,c}^0$ is satisfied.

$R_{q,c}^{\leftarrow}$: Let q, q', c, c' and i, t be such that $\delta(q, c) = (q', c', \leftarrow)$, and the atoms $\mathbf{q}(z_i, z_t)$, $\mathbf{c}(z_i, z_t)$ are in D . Then, by definition of D_n^Q and D_n^Σ , we have $q_t = q$, $h_t = i$ and $M(h_t, t) = c$.

Now, assume that there is a homomorphism from $body(R_{q,c}^{\leftarrow})$ to D that sends p_i to z_i and t_j to z_t , and let j, t be such that p_{i-1} is mapped to z_j and t_{j+1} is mapped to z_t . The atoms $\mathbf{N}(z_j, z_i)$ and $\mathbf{N}(z_t, z_t)$ must then be present in D , so by definition of $D_n^{\mathbf{N}}$, we have $0 \leq j = i - 1$ and $t = t + 1$. To show that $R_{q,c}^{\leftarrow}$ is satisfied, we then show that $\mathbf{q}'(z_{i-1}, z_{t+1})$ and $\mathbf{c}'(z_i, z_{t+1})$ are in D . As $0 \leq j = i - 1$, we have $i > 0$, so by definition of a run, we have $q_{t+1} = q'$, $h_{t+1} = h_t - 1 = i - 1$ and $M(h_t, t + 1) = M(i, t + 1) = c'$. Thus, by definition of D_n^Q and D_n^Σ , the atoms $\mathbf{q}'(z_{i-1}, z_{t+1})$ and $\mathbf{c}'(z_i, z_{t+1})$ are in D , and rule $R_{q,c}^{\leftarrow}$ is satisfied.

$R_{q,c}^{\rightarrow}$: Let q, q', c, c' and i, t be such that $\delta(q, c) = (q', c', \leftarrow)$, and the atoms $\mathbf{q}(z_i, z_t)$, $\mathbf{c}(z_i, z_t)$ are in D . Then, by definition of D_n^Q and D_n^Σ , we have $q_t = q$, $h_t = i$ and $M(h_t, t) = c$, so by definition of a run, we have $q_{t+1} = q'$, $h_{t+1} = h_t + 1 = i + 1$ and $M(h_t, t + 1) = M(i, t + 1) = c'$.

Now, assume that there is a homomorphism from $\text{body}(R_{q,c}^{\rightarrow})$ to D that sends p_i to z_i and t_j to z_t . By definition of D_n^N (as seen in the previous case), p_{i+1} is mapped to z_{i+1} and t_{j+1} is mapped to z_{t+1} . Then, rule $R_{q,c}^{\rightarrow}$ is satisfied, as the atoms $\mathbf{q}'(z_{i+1}, z_{t+1})$ and $\mathbf{c}'(z_i, z_{t+1})$ are in D . Indeed, as $q_{t+1} = q'$ and $h_{t+1} = i + 1$, the first one is in D by definition of D_n^Q , and since $M(i, t + 1) = c'$, the second atom is in D by definition of D_n^Σ .

$R_{q,c}^L$: Assume that the atoms $\mathbf{N}^+(z_i, z_k)$, $\mathbf{N}(z_t, z_{t+1})$, $\mathbf{q}(z_k, z_t)$ and $\mathbf{c}(z_i, z_t)$ are in D , for some i, k, t . Then, we have $h_t = k$, $M(i, t) = c$, and $i < k = h_t$ by definitions of D_n^Q , D_n^Σ and $D_n^{N^+}$, respectively. Thus, by definition of a run, $M(i, t + 1) = M(i, t) = c$, so $\mathbf{c}(z_i, z_{t+1}) \in D$, and $R_{q,c}^L$ is satisfied.

$R_{q,c}^R$: Assume that the atoms $\mathbf{N}^+(z_i, z_k)$, $\mathbf{N}(z_t, z_{t+1})$, $\mathbf{q}(z_i, z_t)$ and $\mathbf{c}(z_k, z_t)$ are in D , for some i, k, t . Then, we have $h_t = i$, $M(k, t) = c$, and $h_t = i < k$ by definitions of D_n^Q , D_n^Σ and $D_n^{N^+}$, respectively. Thus, by definition of a run, $M(k, t + 1) = M(k, t) = c$, so $\mathbf{c}(z_k, z_{t+1}) \in D$, and $R_{q,c}^R$ is satisfied.

R_{func}^+ : Assume that $\mathbf{N}(z_i, z_j)$ and $\mathbf{N}(z_i, z_k)$ are in D , for some i, j, k . Then, by the definition of D_n^N , we have $j = k = i + 1$. Thus, since $\text{Eq}(z_{i+1}, z_{i+1}) \in D$, this rule is satisfied.

R_{func}^- : Assume that $\mathbf{N}(z_i, z_k)$ and $\mathbf{N}(z_j, z_k)$ are in D , for some i, j, k . Then, by the definition of D_n^N , we have $i = j = k - 1$. Thus, since $\text{Eq}(z_{k-1}, z_{k-1}) \in D$, this rule is satisfied.

$R_{N^+}^{cycle}$: By definition of $D_n^{N^+}$, there is no x such that $\mathbf{N}^+(x, x) \in D$, so this rule is not applicable.

R_{q_0} : By the definition of D_n^Q and a Turing machine, there can only be one term that appears in a $\mathbf{q_0}$ -atom, which is z_0 . Thus, if $\mathbf{q_0}(z_i, z_j)$ and $\mathbf{q_0}(z_k, z_l)$ are in D , then $i = j = k = l = 0$. Since $\text{Eq}(z_0, z_0) \in D$, this rule is satisfied.

R_{-1} : Again, the only x such that $\mathbf{q_0}(x, x) \in D$ is z_0 . Thus, by definition of \mathbf{N} , there is no i such that $\mathbf{N}(z_i, z_0) \in D$, so this rule is not applicable.

R_{q_1, q_2}^{step} : Assume that $\mathbf{q_1}(z_i, z_t)$ and $\mathbf{q_2}(z_j, z_t)$ are in D , for some q_1, q_2, i, j, t . By definition of D_n^Q , we then have $i = j = h_t$. Thus, by definition of $D_n^{N^+}$, we do not have $\mathbf{N}^+(z_{h_t}, z_{h_t})$, so this rule is never applicable (and thus satisfied).

R_{q_1, q_2}^{cell} : Assume that $\mathbf{q_1}(z_i, z_t)$ and $\mathbf{q_2}(z_j, z_t)$ are in D , for some q_1, q_2, i, t . By definition of D_n^Q , we then have $q_1 = q_2 = q_t$, so R_{q_1, q_2}^{cell} is not in \mathcal{R}_M (as it requires $q_1 \neq q_2$). Thus, by contrapositive, no rule R_{q_1, q_2}^{cell} in \mathcal{R}_M is applicable.

R_{c_1, c_2}^{cell} : Let c_1, c_2, i, t be such that $\mathbf{c}_1(z_i, z_t)$ and $\mathbf{c}_2(z_j, z_t)$ are in D . By definition of D_n^Σ , we then have $c_1 = c_2 = M(i, t)$, so R_{c_1, c_2}^{cell} is not in \mathcal{R}_M (as it requires $c_1 \neq c_2$). Thus, by contrapositive, no rule R_{c_1, c_2}^{cell} in \mathcal{R}_M is applicable.

R_c^{OoB} : Let c, i, t be such that $\mathbf{c}(z_i, z_t)$ and $\mathbf{N}^+(z_t, z_i)$ are in D . By definition of $D_n^{N^+}$, we then have $t < i$, and by definition of D_n^Σ , we have $i \leq t$, which is a contradiction.

R_{Fail} : As $Fail \notin D$ this rule is not applicable and thus satisfied.

$\{R_{Eq}^{trans}, R_{Eq}^{sym}, R_{P, i}^{repl}\}$: As all the Eq-atoms in D are of the form $\mathbf{Eq}(x, x)$, by Lemma 14, no trigger $t = (R, \pi)$ such that $body(R)$ features an Eq-atom in its body is applicable on D . As this is the case for these three rules, they are satisfied.

As D satisfies all the Datalog rules in \mathcal{R}_M , none of these rules is applicable on D . Since $D \subseteq res(\mathcal{D}'_n)$, we indeed have $res(\mathcal{D}'_n) = D$, which concludes the proof. \square

As the result of all derivations from $(\mathcal{R}_M, \emptyset)$ is characterized by this proposition, Theorem 13, which we recall here, is easy to prove.

Theorem 13. *The Turing machine M halts on the empty word if and only if $(\mathcal{R}_M, \emptyset) \models \mathbf{Goal}$.*

Proof. Consider a Datalog-first derivation $\mathcal{D} = t_0, \mathcal{D}_0, t_1, \mathcal{D}_1, \dots$ from $(\mathcal{R}_M, \emptyset)$. If $(\mathcal{R}_M, \emptyset) \models \mathbf{Goal}$, then there is $n \in \mathbb{N}$ such that $\mathbf{Goal} \in res(\mathcal{D}'_n)$. By Proposition 16, $res(\mathcal{D}'_n)$ contains \mathbf{Goal} if and only if $n \geq time_\varepsilon(M)$. Thus, $time_\varepsilon(M)$ is finite, and M halts. If M halts, then $time_\varepsilon(M)$ is finite. Thus, $res(\mathcal{D}_{time_\varepsilon(M)})$ contains \mathbf{Goal} by Proposition 16, so $(\mathcal{R}_M, \emptyset) \models \mathbf{Goal}$. \square

3.3 \mathcal{R}_M is Datalog expressible

We now show the following theorem.

Theorem 17. *The pair $(\mathcal{R}_M, \mathbf{Goal})$ is Datalog-expressible.*

If M halts on the empty word, then $(\mathcal{R}_M, \emptyset) \models \mathbf{Goal}$ by Theorem 13, and since entailment w.r.t. first-order theories is monotonic, $(\mathcal{R}_M, D) \models \mathbf{Goal}$ for all databases D . Thus, the query $(\{\rightarrow \mathbf{Goal}\}, \mathbf{Goal})$ is a rewriting of \mathcal{R}_M . Now, if M does not halt, while \mathcal{R}_M does not entail \mathbf{Goal} on \emptyset , there are databases on which \mathcal{R}_M entails goal, even if we only consider the rules of \mathcal{R}_M up to $R_{q, c}^R$, that were used for the simulation. For instance, if D already contains \mathbf{Goal} , or a \mathbf{q}_f -atom, then $(\mathcal{R}_M, \emptyset) \models \mathbf{Goal}$, but there are also more interesting (and problematic to find a Datalog rewriting for) databases on which the first part of \mathcal{R}_M would entail \mathbf{Goal} too.

We call databases on which the first part of \mathcal{R}_M entails \mathbf{Goal} even though M does not halt on the empty word *false starts*. To control the structure of the database and to detect these false starts, we introduce a new predicate \mathbf{Fail} , and the rules in the second part of \mathcal{R}_M (from R_{func}^+ onwards). Then, the argument revolves around the notion of Skolem depth of a term.

3.3.1 Skolem depth

Definition 18. Consider a database D , a rule set \mathcal{R} , and a derivation $\mathcal{D} = D_0, t_1, D_1, \dots$ from (\mathcal{R}, D) . For all n , the Skolem depth $SkD(x)$ of a term x appearing in D_n is defined inductively as the following:

- If x is a constant, then its Skolem depth is 0.
- If x is a variable introduced by a trigger t , then its Skolem depth is $\max(\{SkD(y) \mid y \in \text{support}(t)\}) + 1$.

Then, the Skolem depth $SkD(E)$ of a set E of terms is $\max(\{SkD(y) \mid y \in E\})$. The Skolem depth of an atom or a trigger (or a set thereof) is then the Skolem depth of the set of variables appearing in it.

The rough idea behind this definition is that existential rules can be “Skolemized” by replacing their existential quantifiers with Skolem functions, and the Skolem depth of a term is the maximal depth of nested Skolem functions. This notion is interesting here because of Theorem 7 from [4], that we restate here.

Definition 19. A couple rule set-query has bounded depth if there is k such that for all databases D , the following are equivalent:

- $(\mathcal{R}, D) \models q$.
- There is a finite derivation \mathcal{D} such that $\text{res}(\mathcal{D}) \models q$ and all the terms in $\text{res}(\mathcal{D})$ are of Skolem depth less than k .

Theorem 20. If (\mathcal{R}, q) has bounded depth, then it is Datalog-expressible.

Thus, in the following, we show that $(\mathcal{R}_M, \text{Goal})$ has a bounded depth, and that the bound is one. To do so, we distinguish two cases: the databases that entail **Fail** and those that do not. For the databases that entail **Fail**, we show that they do it fast: we only need to close under Datalog to see if a database entails **Fail**. For the other databases, though, we show that (if **Goal** is not in them), then the body of rules that generate **Goal** cannot contain terms of Skolem depth greater than one. However, this is not enough to conclude: it could be the case that, for instance, there is a derivation that can entail $\mathbf{q}_f(x, y)$ with x and y constants, but that this derivation needs terms of Skolem depth greater than one to generate this atom. Thus, we introduce the notion of Skolem-increasing derivation, a restriction of Datalog-first derivations in which after closing under Datalog, we apply the existential trigger that create the variable with the smallest Skolem depth.

Definition 21. A Skolem-increasing derivation is a Datalog-first derivation $\mathcal{D} = \mathcal{D}_{init}, t_0, \mathcal{D}_0, \dots$ such that for all n , there is no existential trigger t applicable on $\text{res}(\mathcal{D}_n)$ such that $SkD(t) < SkD(t_{n+1})$.

Note that Skolem-increasing derivations always exist, as the property that characterizes them is only a criterion for choosing which existential trigger to apply. We do not require for the sequence $(SkD(t_i))_{i \in \mathbb{N}}$ to be increasing, even though we will see that it is indeed the case for our rule set. We also show another important property, which is that to generate an atom featuring terms of Skolem depth less than n , then the sequence of trigger applications that leads to it only features terms of depth less than n .

Proposition 22. *Let D be a database and $\mathcal{D} = \mathcal{D}_0, t_0, \mathcal{D}_1, \dots$ a Datalog-first derivation from (\mathcal{R}_M, D) . Then, the following holds for all terms x and y in $\text{res}(\mathcal{D})$:*

1. *If $\mathbf{N}(x, y) \in \text{res}(\mathcal{D})$, then either $\text{SkD}(x) = \text{SkD}(y) = 0$ or $\text{SkD}(y) = \text{SkD}(x) + 1$.*
2. *If $\mathbf{N}^+(x, y) \in \text{res}(\mathcal{D})$, then $\text{SkD}(x) \leq \text{SkD}(y)$.*
3. *For all states q , if $\mathbf{q}(x, y) \in \text{res}(\mathcal{D})$, then $\text{SkD}(x) \leq \text{SkD}(y)$.*
4. *For all triggers t in \mathcal{D} and atoms $A \in \text{output}(t)$ not using \mathbf{Eq} , we have $\text{SkD}(\text{support}(t)) \leq \text{SkD}(A)$.*

To show this, we first prove a few results about the \mathbf{Eq} predicate. Let us first describe the role of the rules in which \mathbf{Eq} appears.

3.3.2 An axiomatization of equality

First, rules $R_{\mathbf{Eq}}^{\text{trans}}$, $R_{\mathbf{Eq}}^{\text{sym}}$ and $R_{P,i}^{\text{repl}}$ constitute an axiomatization of non-reflexive equality. Then, rules R_{func}^+ and R_{func}^- ensure that \mathbf{N} is both functional and anti-functional, and rule R_{q_0} ensures that there is only one \mathbf{q}_0 -atom. We have already seen in Lemma 14 that making equality reflexive would not yield any new fact (and it is complicated to do properly with existential rules). Thus, we define an equivalence relation \approx which is the reflexive closure of the \mathbf{Eq} predicate, and we show that this relation is indeed an axiomatization of equality.

Lemma 23. *Let D be a Datalog closed database w.r.t \mathcal{R}_M . The relation \approx defined as:*

$$x \approx y \quad \text{iff} \quad \mathbf{Eq}(x, y) \in D \text{ or } x = y$$

is an equivalence relation. In addition, for all Datalog closed instances D , predicates P of arity n and variables $x_1, \dots, x_n, y_1, \dots, y_n$, if $x_i \approx y_i$ for all $i \leq n$, then $P(x_1, \dots, x_n) \in D$ if and only if $P(y_1, \dots, y_n)$.

Proof. Let us first prove that \approx is an equivalence relation.

Symmetry: Let x and y be terms in D such that $x \approx y$. If $x = y$, then $y = x$ so $y \approx x$. Otherwise, it means $\mathbf{Eq}(x, y) \in D$. In this case, since D satisfies rule $R_{\mathbf{Eq}}^{\text{sym}}$, $\mathbf{Eq}(y, x) \in D$. Thus, $y \approx x$.

Reflexivity: If x is a term in D , then $x = x$ so $x \approx x$.

Transitivity: Let x, y and z be terms in D such that $x \approx y$ and $y \approx z$. Then:

- If $x = y$, then since $y \approx z$, we have $x \approx z$.
- If $y = z$, then since $x \approx y$, we have $x \approx z$.
- Otherwise, we have both $\mathbf{Eq}(x, y) \in D$ and $\mathbf{Eq}(y, z) \in D$. Since D is a model of rule $R_{\mathbf{Eq}}^{\text{trans}}$, it means $\mathbf{Eq}(x, z) \in D$, so $x \approx z$.

Thus, \approx is indeed an equivalence relation. Now, for the second part of the lemma, assume that $P(x_1, \dots, x_n) \in D$. We then show by induction on i that $P(y_1, \dots, y_i, x_{i+1}, \dots, x_n) \in D$.

If $i = 0$, by hypothesis we have $P(x_1, \dots, x_n) \in D$, so there is nothing to prove. Now, assume that $i > 0$ and $P(y_1, \dots, y_{i-1}, x_i, \dots, x_n) \in D$. If $x_i = y_i$, then the result follows immediately from the induction hypothesis. Otherwise, we have $\text{Eq}(x_i, y_i) \in D$. Thus, since D is a model of rule $R_{P,i}^{\text{repl}}$, we have $P(y_1, \dots, y_i, x_{i+1}, \dots, x_n) \in D$.

In particular, this result is true for $i = n$, meaning that $P(y_1, \dots, y_n) \in D$. \square

We then consider the quotient of Datalog closed instances, and show that this quotient is also Datalog closed.

Definition 24. For D a Datalog closed instance w.r.t. \mathcal{R}_M , we denote by $[x]$ the class of a term x in D for the equivalence relation \approx , and define D/\approx the quotient of D by \approx as the database such that:

- Its terms are the classes $[x]$ for all $x \in D$.
- We have $P([x_1], \dots, [x_n]) \in D/\approx$ if and only if $P(x_1, \dots, x_n) \in D$.

Lemma 25. If D is a Datalog closed database w.r.t. \mathcal{R}_M , then D/\approx satisfies all the Datalog rules in \mathcal{R}_M that do not feature an Eq-atom in their heads.

Proof. Let $R = B(\vec{x}, \vec{y}) \rightarrow H(\vec{x})$ be a rule in \mathcal{R}_M not featuring an Eq-atom in its head, and π be a homomorphism from $B(\vec{x}, \vec{y})$ to D_{SC} . As $\pi(B(\vec{x}, \vec{y})) \in D/\approx$, if $\pi(\vec{x}) = \vec{c}$ and $\pi(\vec{y}) = \vec{d}$, then for all $\vec{u} \in \vec{c}$ and $\vec{v} \in \vec{d}$, we have $B(\vec{u}, \vec{v}) \in D$. Thus, since D models R , we also have $H(\vec{u}) \in D$, so $H(\vec{c}) = \pi(H(\vec{x})) \in D_{SC}$, and D_{SC} satisfies R . \square

Note that Lemma 23 and Lemma 25 are independent of the rules in \mathcal{R}_M other than our axiomatization, and could be used in any rule set in which the same axiomatization of equality is. Though, since rule $R_{P,i}^{\text{repl}}$ can create any atom, if we want to show that **Fail** cannot appear too late in our context, we need to show that the rules we just described do not equalize too many individuals. This informal intuition calls for the following result, and especially point (4).

Lemma 26. Let D be a database, and $\mathcal{D} = D_0, t_0, D_1, t_1, D_2, \dots$ a Datalog-first restricted derivation from (\mathcal{R}_M, D) (with $D_0 = D$), and \mathcal{D}_0 the prefix of \mathcal{D} composed only of Datalog triggers. Then, the following claims hold for all n :

1. All the atoms $\mathbf{N}(v, w)$ in D_n with w constant are in $\text{res}(\mathcal{D}_0)$.
2. If w is the variable created by an application of rule R_0 , there is no v such that $\mathbf{N}(v, w) \in D_n$.
3. All $\mathbf{q_0}$ -atoms in D_n either feature only constant and are in $\text{res}(\mathcal{D}_0)$, or feature a variable and are created by R_0 .
4. If $\text{Eq}(x, y) \in D_n \setminus \text{res}(\mathcal{D}_0)$, then $x = y$.

Proof. We show the claim by induction on n . The base case is trivial. Indeed, in D there is no variable so (2) holds, and everything that is in D is in $\text{res}(\mathcal{D}_0)$, so the other points hold too. We then prove the induction step. Assume that all the points are true at step $n - 1$, and that $t_n = (R_n, \pi_n)$.

First, regarding point (1), assume by contradiction that t_n creates an atom $\mathbf{N}(v, w)$ with v any term and w a constant. Then, the rule R_n must be either $R_{\mathbf{N},1}^{repl}$ or $R_{\mathbf{N},2}^{repl}$, or w would be a variable. If R_n is $R_{\mathbf{N},1}^{repl}$, then the atoms $\mathbf{N}(u, w)$ and $\text{Eq}(u, v)$ must be in D_{n-1} , with u a term such that $u \neq v$. Then, by point (4) of the induction hypothesis, we have $\text{Eq}(u, v) \in \text{res}(\mathcal{D}_0)$. In addition, $\mathbf{N}(u, w)$ is an \mathbf{N} -atom with a constant on second position in D_{n-1} , so by point (1) of the induction hypothesis, we also have $\mathbf{N}(u, w) \in \text{res}(\mathcal{D}_0)$. Finally, since $\text{res}(\mathcal{D}_0)$ models rule $R_{\mathbf{N},1}^{repl}$, we have $\mathbf{N}(v, w) \in \text{res}(\mathcal{D}_0)$, which contradicts our hypothesis, so rule $R_{\mathbf{N},1}^{repl}$ cannot generate $\mathbf{N}(v, w)$. An analogous argument proves that rule $R_{\mathbf{N},2}^{repl}$ cannot do so either, concluding the proof of this point.

Regarding point (2), assume that w is the variable introduced by an application of rule R_0 , and that t_n creates the atom $\mathbf{N}(v, w)$. First, rules R_1 and R_{+1} both create an \mathbf{N} -atom with an existentially quantified variable in its second position. Since w has already been introduced, it cannot be introduced again by these two rules. Rule $R_{\mathbf{N},1}^{repl}$ requires an atom $\mathbf{N}(u, w)$ to be applicable, which is impossible as $\mathbf{N}(v, x)$ is the first with w in second position. Finally, rule $R_{\mathbf{N},2}^{repl}$ requires an Eq -atom involving w to be applicable, which by induction hypothesis does not exist, which concludes this point.

To show point (3), we first treat the constant case. Assume that t_n creates $\mathbf{q}_0(v, w)$ with v and w constants. Then, as it only features constants, it must have been generated by rule $R_{\mathbf{q}_0,1}^{repl}$ or $R_{\mathbf{q}_0,2}^{repl}$. If $R_n = R_{\mathbf{q}_0,1}^{repl}$, then there are $\mathbf{q}_0(u, w)$ and $\text{Eq}(u, v)$ in D_n , with a term u such that $u \neq v$. Then, by point (4) of the induction hypothesis, we have $\text{Eq}(u, v) \in \text{res}(\mathcal{D}_0)$. In addition, since $\text{Eq}(u, v) \in \text{res}(\mathcal{D}_0)$, u is a constant, so $\mathbf{q}_0(u, w)$ is a \mathbf{q}_0 -atom over constants. Thus, by point (3) of the induction hypothesis, $\mathbf{q}_0(u, w) \in \text{res}(\mathcal{D}_0)$. Then, since $\text{res}(\mathcal{D}_0)$ is closed under rule $R_{\mathbf{q}_0,1}^{repl}$, then $\mathbf{q}_0(v, w) \in \text{res}(\mathcal{D}_0)$. An analogous argument holds if $R_n = R_{\mathbf{q}_0,2}^{repl}$, which concludes the constant case.

Then, for the variable case of point (3), assume that t_n creates the atom $\mathbf{q}_0(v, w)$ with v or w a variable, and that $R_n \neq R_0$. If $R_n = R_{\mathbf{q}_0,1}^{repl}$, then there must be $\mathbf{q}_0(u, w)$ and $\text{Eq}(u, v)$ in D_{n-1} , with u a term such that $u \neq v$. Thus, u and v must be constants. Since we assume v or w is a variable, it must be that w is a variable. Thus, by induction hypothesis, the atom $\mathbf{q}_0(u, w)$ must have been created by rule R_0 , so $u = w$, and then u is both a constant and a variable, which is absurd. An analogous argument concludes this point if $R_n = R_{\mathbf{q}_0,2}^{repl}$.

Then, to prove point (4), assume by contradiction that t_n is not in $\text{res}(\mathcal{D}_0)$, and creates an atom $\text{Eq}(x, y)$ with $x \neq y$ two terms. We then do a case distinction on R_n .

If $R_n = R_{func}^+$: In this case, there are atoms $\mathbf{N}(z, x)$ and $\mathbf{N}(z, y)$ in D_{n-1} with z a term, so that t is applicable. If x and y are constant, then both $\mathbf{N}(z, x)$ and $\mathbf{N}(z, y)$ are in $\text{res}(\mathcal{D}_0)$ by point (1) of the induction hypothesis, so by rule R_{func}^+ , $\text{Eq}(x, y) \in \text{res}(\mathcal{D}_0)$, contradicting our hypothesis. Then, for the case where x or y is a variable, assume that x is a variable (x and y are symmetric in this case), and let $t_k = (R_k, \pi_k)$ be the trigger that creates x .

If $R_k = R_0$, then by point (2), the atom $\mathbf{N}(z, x)$ cannot be in D_{n-1} . Thus, R_k is either R_1 or R_{+1} . In either case, t_k creates the atom $\mathbf{N}(v, x)$ along with x . We show that there is no atom $\mathbf{N}(w, y)$ in D_{n-1} with $w \in [v]$.

Assume by contradiction that the atom $N(w, y)$ is the first atom of this form created. We then study the rule that introduces $N(w, y)$.

If the rule that creates $N(w, y)$ is R_1 or R_{+1} , then there is l such that t_l introduces $N(w, y)$, and $l \neq k$. If $l < k$, then $\text{res}(\mathcal{D}_l) \subseteq \text{res}(\mathcal{D}_{k-1})$. As $\text{res}(\mathcal{D}_l)$ is Datalog closed, by Lemma 23, $N(v, y) \in \text{res}(\mathcal{D}_l)$, so t_k is not applicable, which contradicts our hypothesis. If $l > k$, as $\text{res}(\mathcal{D}_k)$ is Datalog closed, by Lemma 23, $N(w, x) \in \text{res}(\mathcal{D}_k)$, so t_l is not applicable. Thus, these rules cannot create $N(w, y)$.

Rule $R_{N,2}^{\text{repl}}$ requires an Eq-atom involving y to produce $N(w, y)$, so y has to be a constant. Then, by point (1), the atom $N(v, y)$ is in $\text{res}(\mathcal{D}_0)$. Since $\text{res}(\mathcal{D}_0)$ is Datalog closed, by Lemma 23, we have $N(v, y) \in \text{res}(\mathcal{D}_0)$. Thus, t_k is not applicable, as one can map x to y .

For rule $R_{N,1}^{\text{repl}}$ to produce $N(w, y)$, there must be $N(u, y)$ and $\text{Eq}(u, w)$ in D_{n-1} before it. Thus, $u \in [v]$, so since $N(w, y)$ is the first one, this rule cannot create it either.

Then, there is no atom $N(w, y)$ in D_{n-1} with $w \in [v]$. As $v \in [v]$, the atom $N(v, y)$ is not in D_{n-1} either, so t is not applicable in this case.

We can now show that for any z , the atoms $N(z, x)$ and $N(z, y)$ cannot exist at the same time. If the rule that generates $N(z, x)$ is R_1 or R_{+1} , then $z = v$ and we have shown that for all $w \in [v]$, including v itself, $N(w, y) \notin D_{n-1}$. Then, $N(z, y) \notin D_{n-1}$, so these rules cannot create $N(z, x)$. If this rule is $R_{N,1}^{\text{repl}}$, then $z \approx v$, and by the same result $N(z, y) \notin D_{n-1}$. Finally, rule $R_{N,2}^{\text{repl}}$ requires an Eq-atom over x , a variable, so it cannot be used either. Thus, the atoms $N(z, x)$ and $N(z, y)$ cannot exist at the same time.

If $R_n = R_{\text{func}}^-$: In this case, there are atoms $N(x, z)$ and $N(y, z)$ in D_{n-1} with z a term, so that t is applicable. First notice that if z is a constant, by point (1), $N(x, z)$ and $N(y, z)$ are in $\text{res}(\mathcal{D}_0)$. Thus, the closure under rule R_{func}^- entails that $\text{Eq}(x, y) \in \text{res}(\mathcal{D}_0)$, which contradicts our hypothesis. Thus, z must be a variable, so the atoms $N(x, z)$ and $N(y, z)$ must have been created during the chase. The term z being a variable also means that $R_{N,2}^{\text{repl}}$ cannot create any atom $N(v, z)$, as in this case z would be a constant. Let $t_k = (R_k, \pi_k)$ be the trigger that created z .

If $R_k = R_0$, then by point (2), the atom $N(x, z)$ cannot be in D_{n-1} . Thus, R_k must be either R_1 or R_{+1} . In either case, t_k creates the atom $N(v, z)$. We show that for all $w \notin [v]$, the atom $N(w, z)$ is not in D_{n-1} . Assume by contradiction that $N(w, z)$ is the first of these atoms to appear in D_{n-1} . As z is a variable created along the atom $N(v, z)$, only rules $R_{N,1}^{\text{repl}}$ and $R_{N,2}^{\text{repl}}$ can generate the atom $N(w, z)$. As we saw earlier, rule $R_{N,2}^{\text{repl}}$ cannot be used, so there is only rule $R_{N,1}^{\text{repl}}$. This rule requires the atoms $N(u, z)$ and $\text{Eq}(u, w)$ to be in D_{n-1} for it to be applicable (for some u). As by hypothesis $N(w, z)$ is the first one with $w \notin [v]$, we have $u \in [v]$ so $u \approx v$, but since we have $\text{Eq}(u, w)$, we also have $u \approx w$. Thus, by transitivity, we have $w \approx v$, so $w \in [v]$, which is a contradiction. Thus, the only w such that $N(w, z) \in D_{n-1}$ are in $[v]$. Then, if we have $N(x, z)$ and $N(y, z)$

in D_{n-1} , then $x \approx v \approx y$ so as $x \neq y$, we have $\text{Eq}(x, y) \in D_{n-1}$, which contradicts our hypothesis.

If $R_n = R_{q_0}$: In this case, in D_{n-1} there are either the atoms $\mathbf{q_0}(x, y)$ and $\mathbf{q_0}(v, w)$, or the atoms $\mathbf{q_0}(v, x)$ and $\mathbf{q_0}(y, w)$ (with v and w terms). In the first case, as $x \neq y$, the terms x and y must both be constants. Then, $\mathbf{q_0}(x, y) \in \text{res}(\mathcal{D}_0)$ by point (3). Thus, since $\text{res}(\mathcal{D}_0)$ is closed under rule R_{q_0} , we have $\text{Eq}(x, y) \in \text{res}(\mathcal{D}_0)$, contradicting our hypothesis. Thus, there must be atoms $A_x = \mathbf{q_0}(v, x)$ and $A_y = \mathbf{q_0}(y, w)$ in D_{n-1} . By point (3), there are four cases possible.

- If A_x and A_y both only feature constants, they are both in $\text{res}(\mathcal{D}_0)$, so rule R_{q_0} ensures that $\text{Eq}(x, y)$ is in $\text{res}(\mathcal{D}_0)$ too, which contradicts our hypothesis.
- If A_x features constants and A_y a variable, then by rule R_{q_0} , we have $\text{Eq}(v, x) \in \text{res}(\mathcal{D}_0)$, and by rule $R_{\mathbf{q_0}, 1}^{repl}$, the atom $\mathbf{q_0}(x, x)$ is in $\text{res}(\mathcal{D}_0)$. Thus, rule R_0 is not applicable (one can map y to x), and A_y cannot appear.
- The case where A_y features constants and A_x a variable is symmetric to the previous one.
- Finally, if A_x and A_y both feature variables, then they are both created by rule R_0 . As it has an empty body, it can be applied at most once, so as $x \neq y$, this case is not possible either. \square

3.3.3 Skolem-increasing derivations

We are now equipped with the notions we need to show Proposition 22, which we recall here.

Proposition 22. *Let D be a database and $\mathcal{D} = \mathcal{D}_0, t_0, \mathcal{D}_1, \dots$ a Datalog-first derivation from (\mathcal{R}_M, D) . Then, the following holds for all terms x and y in $\text{res}(\mathcal{D})$:*

1. If $\mathbf{N}(x, y) \in \text{res}(\mathcal{D})$, then either $\text{SkD}(x) = \text{SkD}(y) = 0$ or $\text{SkD}(y) = \text{SkD}(x) + 1$.
2. If $\mathbf{N}^+(x, y) \in \text{res}(\mathcal{D})$, then $\text{SkD}(x) \leq \text{SkD}(y)$.
3. For all states q , if $\mathbf{q}(x, y) \in \text{res}(\mathcal{D})$, then $\text{SkD}(x) \leq \text{SkD}(y)$.
4. For all triggers t in \mathcal{D} and atoms $A \in \text{output}(t)$ not using Eq , we have $\text{SkD}(\text{support}(t)) \leq \text{SkD}(A)$.

Proof. For point (1), assume for a contradiction that there is a first atom $\mathbf{N}(x, y) \in \text{res}(\mathcal{D})$ such that neither $\text{SkD}(x) = \text{SkD}(y) = 0$ nor $\text{SkD}(y) = \text{SkD}(x) + 1$. If y is a constant, then by point (1) of Lemma 26, x must also be a constant so $\text{SkD}(x) = \text{SkD}(y) = 0$. Otherwise, if y is a variable, then $\mathbf{N}(x, y)$ has been generated during the chase. Let $t = (R, \pi)$ be the trigger that generates this atom.

If $R \in \{R_1, R_{+1}\}$: Then t also creates y from x , so $\text{SkD}(y) = \text{SkD}(x) + 1$.

If $R \in \{R_{N,1}^{repl}, R_{N,2}^{repl}\}$: Then there is $N(z, t)$ with $z \approx x$ and $t \approx y$. Since y is a variable, it must be that $t = y$ by Lemma 26. As $N(x, y)$ is the first N -atom for which the statement does not hold, and $SkD(y) \neq 0$, it must be that $SkD(y) = SkD(z) + 1$. If x is a constant, then z is too so $SkD(z) = 0$ so $SkD(y) = 1$. Thus, as $SkD(x) = 0$, it is indeed the case that $SkD(y) = SkD(x) + 1$. Now, if x is a variable, then $z = x$ so the atom $N(x, y)$ must be used to generate itself, which is a contradiction.

Thus, point (1) holds.

For point (2), assume for a contradiction that there is a first atom $N^+(x, y) \in res(\mathcal{D})$ such that $SkD(x) > SkD(y)$ for some x, y . Then, x must be a variable so $N^+(x, y)$ must have been generated during the chase; let $t = (R, \pi)$ be the corresponding trigger.

If $R = R_{N^+}$: Then the atom $N(x, y)$ is in $res(\mathcal{D})$, so by point (1), $SkD(x) \leq SkD(y)$.

If $R = R_{N^+}^{trans}$: Then there are atoms $N(x, z)$ and $N^+(z, y)$ in $res(\mathcal{D})$. By point (1), we then have $SkD(x) \leq SkD(z)$. In addition, as $N^+(x, y)$ is the first N^+ -atom such that $SkD(x) > SkD(y)$, we also have $SkD(z) \leq SkD(y)$. Thus, by transitivity, we get $SkD(x) \leq SkD(y)$.

If $R \in \{R_{N^+,1}^{repl}, R_{N^+,2}^{repl}\}$: Then there is an atom $N^+(z, t)$ with $z \approx x$ and $t \approx y$. As x is a variable, $z = t$. Since $N^+(x, y)$ is the first N^+ -atom such that $SkD(x) > SkD(y)$, $SkD(z) \leq SkD(t)$. Thus, as $x = z$ and x is a variable, t must also be a variable, so $t = y$, which means $N^+(x, y)$ is necessary to generate itself, a contradiction.

Thus, point (2) holds too.

Then, for point (3), assume for a contradiction that there is a first atom $q(x, y) \in res(\mathcal{D})$ such that $SkD(x) > SkD(y)$ for some x, y, q . Then, x must be a variable so $q(x, y)$ must have been generated during the chase; let $t = (R, \pi)$ be the corresponding trigger.

If $R = R_0$: In this case $x = y$ which is not compatible with them having different Skolem depth.

If $R = R_{q,c}^{\leftarrow}$: Then there must be atoms $N(x, z)$ and $N^+(z, y)$ for some term z (among other atoms) in $res(\mathcal{D})$. Thus, by points (1) and (2), $SkD(x) \leq SkD(z) \leq SkD(y)$.

If $R = R_{q,c}^{\rightarrow}$: Then there must be an atom $N^+(x, y)$ in $res(\mathcal{D})$, so by point (2), $SkD(x) \leq SkD(y)$.

If $R = R_{q,c}^{max}$: Then $x = y$ which is not compatible with them having different Skolem depth.

If $R \in \{R_{q_0,1}^{repl}, R_{q_0,2}^{repl}\}$: Then there must be an atom $q(z, t)$ with $z \approx x$ and $t \approx y$. Note that by Lemma 26, for all terms a and $b \in [a]$, $SkD(a) = SkD(b)$. Thus, as $q(z, t)$ must be in $res(\mathcal{D})$ before $q(x, y)$, $SkD(z) \leq SkD(t)$, so $SkD(x) \leq SkD(y)$.

Thus, such an atom cannot exist, so point (3) holds.

Finally, for point (4), let $t = (R, \pi)$ be a trigger. As rules R_{Goal} , $R_{q,c}^0$, $R_{N^+}^{\text{cycle}}$, R_{-1} , $R_{q_1,q_2}^{\text{step}}$, $R_{q_1,q_2}^{\text{cell}}$, $R_{c_1,c_2}^{\text{cell}}$, R_c^{OoB} and R_{Fail} do not feature any variable in their heads, if R is among them, point (4) is vacuously true. The same holds for rules R_{func}^+ , R_{func}^- , R_{q_0} , $R_{\text{Eq}}^{\text{trans}}$ and $R_{\text{Eq}}^{\text{sym}}$, which only introduce Eq-atoms. We then show that for all the other rules, point (4) holds by showing that for all the variables in their bodies, there is a variable in each atom of the head that has a greater Skolem depth. We establish the order between the variables' Skolem depths using the previous points.

If $R = R_0$: This rule features no atoms in its body, so the condition is vacuously true.

If $R = R_1$: We have $\text{SkD}(w_0) \leq \text{SkD}(w_1)$ and w_1 is in the only atom of the head.

If $R = R_{+1}$: We have $\text{SkD}(w_0) \leq \text{SkD}(w_i) \leq \text{SkD}(w_{i+1})$ and w_{i+1} is in the head.

If $R = R_{N^+}$: The same variables appear in the only atoms in the body and the head of this rule.

If $R = R_{N^+}^{\text{trans}}$: We have $\text{SkD}(w_i) \leq \text{SkD}(w_{i+1}) \leq \text{SkD}(w_k)$ and w_k is in the head.

If $R = R_B^0$: The same variable appear in the only atoms in the body and the head of this rule.

If $R = R_B$: We have $\text{SkD}(w_0) \leq \text{SkD}(w_i)$ and w_i is in the head.

If $R = R_{q,c}^{\leftarrow}$: We have $\text{SkD}(p_1) \leq \text{SkD}(p_2) \leq \text{SkD}(t_1) \leq \text{SkD}(t_2)$, and t_2 is in both the atoms in the head.

If $R = R_{q,c}^{\rightarrow}$: We have $\text{SkD}(p_1) \leq \text{SkD}(p_2) \leq \text{SkD}(t_2)$ and $\text{SkD}(p_1) \leq \text{SkD}(t_1) \leq \text{SkD}(t_2)$, so t_2 has the greatest Skolem depth of all these variables, and appears in both atoms in the head.

If $R = R_{q,c}^{\text{max}}$: We have $\text{SkD}(w_1) \leq \text{SkD}(w_2)$, and w_2 appears in the head.

If $R = R_{q,c}^L$: We have $\text{SkD}(p_i) \leq \text{SkD}(p_k) \leq \text{SkD}(t_1) \leq \text{SkD}(t_2)$, and t_2 appears in the head.

If $R = R_{q,c}^R$: We have $\text{SkD}(p_i) \leq \text{SkD}(p_k)$ and $\text{SkD}(p_i) \leq \text{SkD}(t_1) \leq \text{SkD}(t_2)$, and p_k and t_2 both appear in the head.

If $R = R_{P,i}^{\text{repl}}$: The variables x_j for $j \neq i$ all appear in the head. Regarding x_i , by point 4, either it is a constant, and y_i is a constant too so they both have a Skolem depth of zero, or it is a variable and $x_i = y_i$. Thus, either way, the head has a greater Skolem depth than the body. \square

We are now equipped to show that in Skolem-increasing derivations, all the facts entailed by the knowledge base about terms of depth k are in the database before a term of depth $k + 1$ is produced.

Proposition 27. *Consider a database D , a Skolem-increasing derivation $\mathcal{D} = \mathcal{D}_0, t_1, \mathcal{D}_1$ from (\mathcal{R}_M, D) , and an atom A not using Eq such that $(\mathcal{R}_M, D) \models A$. Then, for all n , if $SkD(A) < SkD(t_{n+1})$, then $A \in \mathcal{D}_n$.*

Proof. Assume for a contradiction that there is an atom A that contradicts the proposition. Then, there must be a sequence of trigger t^1, \dots, t^k appearing in \mathcal{D} such that $t^1 = t_n$, $A \in output(t^k)$, and for all $i < k$, $support(t^{i+1}) \cap output(t^i) \neq \emptyset$. As we saw in the proof of Proposition 22, all the rules either have only one atom in the head or both their atoms feature the variable with the greatest Skolem depth. Thus, for all $i < k$, $SkD(output(t^i)) \leq SkD(support(t^{i+1}))$ by point (4) of Proposition 22, for all $i \leq k$, $SkD(support(t^i)) \leq SkD(output(t^i))$. Thus, both these results together give by transitivity that $SkD(t^1) \leq SkD(t^n)$, which entails that $SkD(t_{n+1}) \leq SkD(A)$, and contradicts our hypothesis. \square

Note that, on a side note, this also entails that the sequence $SkD(t_n)$ is increasing. Now that we have this fundamental result, we can go back to our dichotomy about the Fail predicate.

3.3.4 Some structure — The Fail predicate

First, we describe the rules that generate Fail, through a description of what a Datalog closed database without Fail looks like. First, rule $R_{N^+}^{cycle}$ ensures that N , which is both functional and inverse functional, is without cycles. Thus, the graph formed by the equivalence classes by Eq as nodes and N -atoms as edges is a disjoint union of what we call chains.

Definition 28. *In a graph $G = (V, E)$, a chain C of size n is a sequence of vertices v_1, \dots, v_n such that for all $(v, w) \in E$, either there is i such that $v = v_i$ and $w = v_{i+1}$, or v and w are both not in C .*

Then, rule R_{-1} ensure that the only q_0 -atom is at the start of a chain. Finally, rules R_{q_1, q_2}^{step} , R_{q_1, q_2}^{cell} , R_{c_1, c_2}^{cell} and R_c^{OoB} ensure that on each chain, there is only one Turing machine simulation at a step. We formalize this into the following proposition.

Definition 29. *Consider an instance D and a binary predicate P . Then, the P -graph associated with D , denoted with G_D^P , is the graph having as vertices the terms of D , and an edge between x and y if and only if $P(x, y) \in D$.*

Proposition 30. *Let D be a Datalog closed instance that does not contain Fail. Then:*

- *The graph $G_{D/\approx}^N$ is a union of chains.*
- *There is at most one q_0 -atom in D/\approx , and if there is one, it is of the form $q_0([z_0], [z_0])$ with $[z_0]$ being the first element of a chain. We call this chain the starting chain.*

Proof. We first show that $G_{D/\approx}^N$ is a union of chain. First, we show that N is functional, that is that for all $c \in D/\approx$, there is at most one $d \in D/\approx$ such that $N(c, d) \in D/\approx$. Indeed, if d_1 and d_2 are such that $N(c, d_1)$ and $N(c, d_2)$ are in D/\approx , let $x \in c$, $y_1 \in d_1$ and $y_2 \in d_2$. Then, we have $N(x, y_1)$ and $N(x, y_2)$ in D . Thus, since D satisfies rule R_{func}^+ , we have $Eq(y_1, y_2) \in D$, so $y_1 \approx y_2$,

and $d_1 = d_2$. The same reasoning using rule R_{func}^- shows that \mathbf{N} is also inverse functional in D/\approx .

Then, to show that $G_{D/\approx}^{\mathbf{N}}$ is a union of chains, we must show that it is without cycles. Assume by the absurd that it contains a cycle c_1, \dots, c_n . Let x_i be a representative of c_i for all i . Then, we have $\mathbf{N}(x_n, x_0)$ and $\mathbf{N}(x_i, x_{i+1})$ in D for all $i < n$. We then show by induction on i that for all i , $\mathbf{N}^+(x_{n-i}, x_0) \in D$. For the base case, since D is a model of rule $R_{\mathbf{N}^+}$, we have $\mathbf{N}^+(x_n, x_0)$ in D . For the induction step, assume that $\mathbf{N}^+(x_{n-i+1}, x_0) \in D$. Then, since $\mathbf{N}(x_{n-i}, x_{n-i+1}) \in D$ and D models $R_{\mathbf{N}^+}^{trans}$, we have $\mathbf{N}^+(x_{n-i}, x_0) \in D$, which concludes the inclusion. Thus, in particular for $i = n$, we have $\mathbf{N}^+(x_0, x_0) \in D$. Since D models rule $R_{\mathbf{N}^+}^{cycle}$, the predicate **Fail** is in D , which contradicts our initial hypothesis. Thus, $G_{D/\approx}^{\mathbf{N}}$ contains no cycle, and is indeed a union of chains.

Now for the second point, if we have $\mathbf{q}_0(c_1, c_2) \in D/\approx$ and $\mathbf{q}_0(c_3, c_4) \in D/\approx$, let x_i be a representative of c_i for all i , and then we have $\mathbf{q}_0(x_1, x_2) \in D$ and $\mathbf{q}_0(x_3, x_4) \in D$. As D is a model of $R_{\mathbf{q}_0}$, the atoms $\mathbf{Eq}(x_1, x_2)$, $\mathbf{Eq}(x_2, x_3)$ and $\mathbf{Eq}(x_3, x_4)$ are all in D . Thus, $x_1 \approx x_2 \approx x_3 \approx x_4$, so $c_1 = c_2 = c_3 = c_4$. Thus, there is at most one \mathbf{q}_0 -atom in D/\approx , which is of the form $\mathbf{q}_0(c_0, c_0)$. We then show that if it exists, there is no c such that $\mathbf{N}(c, c_0) \in D/\approx$. Indeed, if by contradiction we assume there is one, let $x \in c$ and $z_0 \in c_0$. We then have $\mathbf{N}(x, z_0) \in D$. Since D models rule R_{-1} , the predicate **Fail** is entailed, which contradicts our initial hypothesis, and concludes the proof. \square

We now are equipped to show that if the Datalog closure of a database D does not contain **Fail**, and $\mathcal{D} = \mathcal{D}_0, t_0, \mathcal{D}_1, \dots$ is a Datalog-first restricted derivation from (\mathcal{R}_M, D) , then $res(\mathcal{D}_n)$ does not contain **Fail** either. Thus, Proposition 30 applies on $res(\mathcal{D}_n)$ too, for all n . Therefore, this structure is maintained through all derivations from databases that do not entail **Fail**.

Proposition 31. *Consider a database D and a Datalog-first restricted derivation $\mathcal{D} = \mathcal{D}_0, t_1, \mathcal{D}_1, \dots$. If **Fail** $\notin \mathcal{D}_0$, then for all n , **Fail** $\notin \mathcal{D}_n$.*

Proof. We proceed by induction on n . The base case is trivial, so we focus on the induction step. Assume that **Fail** $\notin \mathcal{D}_{n-1}$. First, as $t_n = (R_n, \pi_n)$ is an existential trigger, R_n must be among R_0, R_1 and R_{+1} . Thus, it follows from its \mathbb{R} -applicability that t_n can only either create a starting chain of length one if there is none, or add one element at the end of a chain. Then, by point (4) of Lemma 26 and the fact that no other rule can create \mathbf{q}_0 -atoms or \mathbf{N} -atoms, this structure cannot change during \mathcal{D}_n . Thus, rules $R_{\mathbf{N}^+}^{cycle}$ and R_{-1} cannot fire and entail **Fail**.

Regarding the other four rules, they follow the same pattern: if they were applicable on $res(\mathcal{D}_n)$, then their supports must have been introduced at a previous step, which is not possible by induction hypothesis. \square

3.3.5 Controlling the chains

We then come back to our result through several lemmas. The main idea is that since Datalog closed databases are composed of chains, we then show that each chain is well-behaved, and that nothing inconvenient happens between chains. We start with the starting chain, by showing that it behaves exactly like the chase on $(\mathcal{R}_M, \emptyset)$.

Lemma 32. *Consider D be a Datalog closed instance that neither contain **Fail** nor **Goal**, and such that D/\approx has a starting chain. Then, let $\mathcal{D} = t_0, \mathcal{D}_0, t_1, \dots$ be a Datalog-first derivation from $(\mathcal{R}_M, \emptyset)$, and D_{SC} be the instance obtained by removing all the classes in D/\approx that are not on the starting chain, and all the atoms in which such a term appears. Finally, let n be the number of classes on the starting chain. Then, D_{SC} and $\text{res}(\mathcal{D}_{n-1})$ are isomorphic.*

Proof. First, by Proposition 30, $G_{D_{SC}}^N$ is a chain of size n , and by Proposition 16, $G_{\text{res}(\mathcal{D}_{n-1})}^N$ is also a chain of size n . We thus denote the elements of $\text{res}(\mathcal{D}_{n-1})$ and D_{SC} as z_i and d_i , respectively, where i is z_i and d_i 's position in the chain (and ranges from 0 to $n-1$). We then show that the map h that sends z_i to d_i , which is a bijection, is an isomorphism.

Let us start by showing that $h(\text{res}(\mathcal{D}_{n-1})) \subseteq D_{SC}$. More precisely, we show by induction on i that for all atoms $A \in \text{res}(\mathcal{D}_{n-1})$ featuring only variables among $\{z_j \mid 0 \leq j \leq i\}$, we have $h(A) \in D_{SC}$.

For the base case, the atoms in $\text{res}(\mathcal{D}_{n-1})$ featuring only z_0 are $\mathbf{q}_0(z_0, z_0)$ and $\mathbf{B}(z_0, z_0)$. Then, $\mathbf{q}_0(d_0, d_0) \in D_{SC}$ since D_{SC} is the starting chain in D/\approx . Finally, as $\mathbf{q}_0(d_0, d_0) \in D_{SC}$, since D_{SC} models rule $R_{\mathbf{B}}^0$ by Lemma 25, we have $\mathbf{B}(d_0, d_0) \in D_{SC}$, which concludes the base case.

For the inductive step, assume that for all atoms $A \in \text{res}(\mathcal{D}_{n-1})$ featuring only variables among $\{z_j \mid 0 \leq j \leq i\}$, we have $h(A) \in D_{SC}$, and let A be an atom in $\text{res}(\mathcal{D}_{n-1})$ featuring only variables among $\{z_j \mid 0 \leq j \leq i+1\}$. If it only features variables among $\{z_j \mid 0 \leq j \leq i\}$, then the induction hypothesis suffices, so we assume that it features z_{i+1} . As by Proposition 16 we have $\text{res}(\mathcal{D}_{n-1})$ characterized, we now exactly which atoms fulfil this specification.

$N(z_i, z_{i+1})$: As we said earlier, the graph $G_{D_{SC}}^N$ is a chain, so our naming convention ensures that $N(d_i, d_{i+1}) \in D_{SC}$.

$N^+(z_j, z_{i+1})$ for $j \leq i$: By the previous point, $N(d_i, d_{i+1}) \in D_{SC}$, so by rule R_{N^+} , we have $N^+(d_i, d_{i+1}) \in D_{SC}$. Then, if by induction we assume that $N^+(c_{j+1}, d_{i+1}) \in D_{SC}$, as by induction hypothesis $N(c_j, c_{j+1}) \in D_{SC}$, then $N^+(c_j, d_{i+1}) \in D_{SC}$.

$\mathbf{q}_{i+1}(z_{h_{i+1}}, z_{i+1})$: As we saw earlier, the atom $N(d_i, d_{i+1})$ is in D_{SC} . Then, as for all t we have $h_t \leq t$, in particular we have $h_i \leq i$, so by induction hypothesis the atoms $\mathbf{q}_i(c_{h_i}, d_i)$, $N^+(c_{h_i}, d_{i+1})$ and $\mathbf{c}(c_{h_i}, z_i)$ with $c = M(h_i, i)$ are all in D_{SC} . We also have $h_{i+1} \leq i+1$, so the atom $N(c_{h_{i+1}}, c_{h_i})$ or $N(c_{h_i}, c_{h_{i+1}})$ is in D_{SC} , depending on whether $\delta(q_i, c) = (q_{i+1}, c', \leftarrow)$ or $\delta(q_i, c) = (q_{i+1}, c', \rightarrow)$. Thus, by rules $R_{q,c}^{\leftarrow}$ and $R_{q,c}^{\rightarrow}$, we have $\mathbf{q}_{i+1}(c_{h_{i+1}}, d_{i+1})$ and $\mathbf{c}'(z_{h_i}, z_{i+1})$ in D_{SC} , with $c' = M(h_i, i+1)$.

$\mathbf{c}_j(z_j, z_{i+1})$ for $j \leq i+1$, with $c_j = M(j, i+1)$: By the previous case, if $j = h_i$, this atom is in D_{SC} , so we assume $j \neq h_i$ in the following. The atom $N(z_i, z_{i+1})$ is in D_{SC} , and by induction hypothesis, the atoms $\mathbf{q}_i(d_{h_i}, d_i)$ and $\mathbf{c}_j(d_j, z_i)$ are in D_{SC} , since for all $j \neq h_i$, $M(j, i) = M(j, i+1)$. Then, if $j < h_i$, we have $N^+(d_j, d_{h_i}) \in D_{SC}$ so by rule $R_{q,c}^L$, the atom $\mathbf{c}_j(z_j, z_{i+1})$ is in D_{SC} . If $j > h_i$, we reach the same conclusion using rule $R_{q,c}^R$ instead.

Thus, this concludes the proof that for all i , and all atoms $A \in \text{res}(\mathcal{D}_{n-1})$ featuring only variables among $\{z_j \mid 0 \leq j \leq i\}$, we have $h(A) \in D_{SC}$. Thus, in particular, $h(\text{res}(\mathcal{D}_{n-1})) \subseteq D_{SC}$.

We then prove that h^{-1} is a homomorphism. Assume by contradiction that A is an atom in D_{SC} such that $h^{-1}(A) \notin \text{res}(\mathcal{D}_{n-1})$. First note that we assumed that **Fail** and **Goal** were not in D , so they are not in D_{SC} , and since $G_{D_{SC}}^N$ and $G_{\text{res}(\mathcal{D}_{n-1})}^N$ are both chains, A cannot use the predicate **N**. Then, we show that D_{SC} must contain **Fail**, contradicting our initial hypothesis on D .

If $A = \mathbf{N}^+(d_i, d_j)$: As $h^{-1}(A) = \mathbf{N}^+(z_i, z_j) \notin \text{res}(\mathcal{D}_{n-1})$, it must then be that $j \leq i$. We then show by induction on $i - j$ that if $\mathbf{N}^+(z_i, z_j) \in D_{SC}$ with $j \leq i$, then **Fail** $\in D_{SC}$. For the base case, if $i - j = 0$, then $j = i$. Since D_{SC} models rule $R_{\mathbf{N}^+}^{\text{cycle}}$, **Fail** must then be in D_{SC} . Otherwise, if $i - j > 0$, then $j < i$, so $\mathbf{N}(d_{i-1}, d_i) \in D_{SC}$. Thus, by rule $R_{\mathbf{N}^+}^{\text{trans}}$, $\mathbf{N}(d_j, d_{i-1}) \in D_{SC}$. Then, $(i - 1) - j = (i - j) - 1$, so by induction hypothesis, we have **Fail** $\in D_{SC}$.

If $A = \mathbf{q}(d_i, d_j)$: As $\text{time}_\varepsilon(M) = \infty$, for all j , we have $\mathbf{q}_j(z_{h_j}, z_j) \in \text{res}(\mathcal{D}_{n-1})$, so by the previous point, $\mathbf{q}_j(d_{h_j}, d_j) \in D_{SC}$. Since $h^{-1}(A) \notin \text{res}(\mathcal{D}_{n-1})$, the atom A must be different from $\mathbf{q}_j(d_{h_j}, d_j)$. If $h_j \neq i$, then either $\mathbf{N}^+(d_{h_j}, d_i)$ or $\mathbf{N}^+(d_i, d_{h_j})$ is in D_{SC} . In either case, since D_{SC} models rule $R_{q_1, q_2}^{\text{step}}$, **Fail** is in D_{SC} . Otherwise, if $i = h_j$, it must be that $q \neq q_j$. Thus, as D_{SC} models rule $R_{q_1, q_2}^{\text{cell}}$, **Fail** is again in D_{SC} .

If $A = \mathbf{c}(d_i, d_j)$: First, if $j < i$, then $\mathbf{N}^+(z_j, z_i) \in \text{res}(\mathcal{D}_{n-1})$, so $\mathbf{N}^+(d_j, d_i) \in D_{SC}$. Thus, as D_{SC} models rule R_c^{OoB} , **Fail** is in D_{SC} . We then assume that $i \leq j$. As by hypothesis, $\text{time}_\varepsilon(M) = \infty$, for all $i \leq j$, the atom $\mathbf{c}_i(z_i, z_j)$ is in $\text{res}(\mathcal{D}_{n-1})$, with $c_j = M(i, j)$. Thus, by the previous point, $\mathbf{c}_i(d_i, d_j) \in D_{SC}$. Since $h^{-1}(A) \notin \text{res}(\mathcal{D}_{n-1})$, the atom A must be different from $\mathbf{c}_i(d_i, d_j)$, so $c \neq c_i$. Thus, since D_{SC} models rule $R_{c_1, c_2}^{\text{cell}}$, **Fail** is in D_{SC} .

Thus, $\text{res}(\mathcal{D}_{n-1})$ and D_{SC} are indeed isomorphic. \square

Now, since we assume that M does not halt on the empty word, by Theorem 13, we have $(\mathcal{R}_M, \emptyset) \not\models \text{Goal}$, so what happens on the starting chain alone cannot entail **Goal**. We then show that the non-starting chains cannot grow much.

Lemma 33. *Let D be a database, and $\mathcal{D} = \mathcal{D}_0, t_0, \mathcal{D}_1, \dots$ a Datalog-first restricted derivation from (\mathcal{R}_M, D) . Then, for all $n \in \mathbb{N}$, if $\text{res}(\mathcal{D}_0)$ does not contain **Fail**, then:*

1. *If x or y is a variable and $\mathbf{N}^+([x], [y]) \in \text{res}(\mathcal{D}_n) / \approx$, then $[x]$ and $[y]$ are on the same chain.*
2. *All non-starting chains only feature classes of Skolem depth at most one.*

Proof. For point (1), assume for a contradiction that $\mathbf{N}^+(x, y)$ is the first \mathbf{N}^+ -atom in $\text{res}(\mathcal{D}_n)$ featuring a variable, and such that x and y are two terms on different chains. Then, there is a trigger $t = (R, \pi)$ that created $\mathbf{N}^+(x, y)$.

If $R = R_{\mathbf{N}^+}$: Then the atom $\mathbf{N}(x, y)$ is in $\text{res}(\mathcal{D}_n)$, so x and y are on the same chain.

If $R = R_{N^+}^{trans}$: Then the atoms $\mathbf{N}(x, z)$ and $\mathbf{N}^+(z, y)$ must be in $res(\mathcal{D}_n)$ for some term z . First, since $\mathbf{N}(x, z)$ is in $res(\mathcal{D}_n)$, then x and z are on the same chain. It remains to show that z and y are on the same chain.

If z is a variable, then $\mathbf{N}^+(z, y)$ is an \mathbf{N}^+ -atom containing a variable created before $\mathbf{N}^+(x, y)$, so z and y are on the same chain by hypothesis.

If z is a constant, then by point (1) of Lemma 26, we have $\mathbf{N}(x, z) \in res(\mathcal{D}_0)$, so x is a constant. Thus, y has to be a variable, and $\mathbf{N}^+(z, y)$ is again an \mathbf{N}^+ -atom containing a variable created before $\mathbf{N}^+(x, y)$, so z and y are on the same chain by induction hypothesis.

Thus, in both cases, the result holds.

If $R \in \{R_{N^+,1}^{repl}, R_{N^+,2}^{repl}\}$: Then there must be $\mathbf{N}^+(z, t)$ in $res(\mathcal{D}_n)$ with $z \approx x$ and $t \approx y$. Then, if x is a variable, by point (4) of Lemma 26, $[x] = \{x\}$, so $z = x$ and $\mathbf{N}^+(z, t)$ is an \mathbf{N}^+ -atom featuring a variable introduced before $\mathbf{N}^+(x, y)$, so z and t are on the same chain. The same reasoning applies if y is a variable. Thus, in both cases, z and t are on the same chain, so x and y are on the same chain too.

Thus, no rule can generate such an atom, which then cannot exist, so point (1) holds.

We then prove point (2). Let y be a variable on a non-starting chain, and let $t = (R, \pi)$ be the rule that created it. We then prove that its Skolem depth must be one. As y is on a non-starting chain, we cannot have $\mathbf{q}_0(y, y) \in res(\mathcal{D}_n)$, so $R \neq R_0$. In addition, for the same reason, $R \neq R_1$, as otherwise there would be a term x such that $\mathbf{q}_0(x, x)$ and $\mathbf{N}(x, y)$ are in $res(\mathcal{D}_n)$, so y would be on the starting chain. Thus, it must be that $R = R_{+1}$. Then, there are some terms z_0 and x such that $\mathbf{q}_0(z_0, z_0)$, $\mathbf{N}^+(z_0, x)$ and $\mathbf{N}(x, y)$ are in $res(\mathcal{D}_n)$. As y is on a non-starting chain, z_0 is not on the same chain as y , and thus as x . Thus, by point (1), z_0 and x must both be constants, so $SkD(y) = 1$, concluding the proof. \square

Finally, we need to control what happens between chains. What happens between non-starting chains is not an issue as all the terms involved are of Skolem depth at most one by Lemma 33, but if state and character atoms are between a non-starting and a starting chain, there could be a false start that is allowed to grow enough to entail **Goal** using terms of unbounded depth. To prevent this, we show that there can be no atom using a state predicate and terms of Skolem depth greater than one on different chains.

Lemma 34. *Let D be a database, and $\mathcal{D} = \mathcal{D}_0, t_0, \mathcal{D}_1, \dots$ a Datalog-first restricted derivation from (\mathcal{R}_M, D) . Then, for all $n \in \mathbb{N}$, if $res(\mathcal{D}_0)$ does not contain **Fail**, and $\mathbf{q}(x, y) \in res(\mathcal{D}_n)$ for some state q and terms x, y where x or y is a variable, then x and y are on the same chain.*

Proof. Assume for a contradiction that there is a first atom $\mathbf{q}(x, y) \in res(\mathcal{D}_n)$ for some state q and terms x and y on different chains with x or y a variable. This atom must have been generated by some trigger $t = (R, \pi)$, as it features a variable.

If $R = R_0$: Then $x = y$ so x and y are on the same chain.

If $R = R_{q,c}^{\leftarrow}$: Then there are atoms $N(x, z)$ and $N^+(z, y)$ in $res(\mathcal{D}_n)$ for some term z . By point (1) of Lemma 26, if z is a constant, then x has to be a constant too, so y must be a variable. Thus, in all cases, the atom $N^+(z, y)$ features a variable, so z and y must be on the same chain by point (1) of Lemma 33. Thus, as x and z are on the same chain, x and y must be too.

If $R = R_{q,c}^{\rightarrow}$: Then there must be an atom $N^+(x, y)$ in $res(\mathcal{D}_n)$, so x and y must be on the same chain by point (1) of Lemma 33.

If $R = R_{q,c}^{max}$: Then again $x = y$ so x and y are on the same chain.

If $R \in \{R_{q,1}^{repl}, R_{q,2}^{repl}\}$: Then there must be an atom $\mathbf{q}(z, t)$ with $z \approx x$ and $t \approx y$. Since this atom has to come before $\mathbf{q}(x, y)$, and x or y is a variable, so either $z = x$ or $t = y$ by (26), z and t must be on the same chain. As being on the same chain is a notion invariant through \approx , x and y are on the same chain too.

□

3.3.6 Skolem depth and Datalog-expressibility

We now are fully equipped to show that $(\mathcal{R}_M, \text{Goal})$ has bounded depth.

Theorem 35. *If M does not halt on the empty word, then $(\mathcal{R}_M, \text{Goal})$ has bounded depth.*

Proof. Let D be a database and $\mathcal{D} = \mathcal{D}_0, t_1, \mathcal{D}_1, \dots$ any Skolem-increasing derivation. Then, assume that $(\mathcal{R}_M, D) \models \text{Goal}$, and let n be the smallest integer such that $\text{Goal} \in res(\mathcal{D}_n)$. We then show that all the terms appearing in $res(\mathcal{D}_n)$ are of Skolem depth at most one. Notice that this is enough to show that $(\mathcal{R}_M, \text{Goal})$ has bounded depth, as the other direction is trivial.

First, if $n = 0$, as $res(\mathcal{D}_0)$ only contains constants (of Skolem depth zero), the statement holds. Thus, in the following, we assume $\text{Goal} \notin res(\mathcal{D}_0)$. As $res(\mathcal{D}_0)$ models rule R_{Fail} , this also means that $\text{Fail} \notin res(\mathcal{D}_0)$, which by Proposition 31 entails that $(\mathcal{R}_M, D) \not\models \text{Fail}$. Let n be the smallest integer such that $\text{Goal} \in res(\mathcal{D}_n)$. Then, Goal must have been introduced by some trigger $t = (R, \pi)$ during $res(\mathcal{D}_n)$. As $(\mathcal{R}_M, D) \not\models \text{Fail}$, it must be that $R \neq R_{\text{Fail}}$.

If $R = R_{\text{Goal}}$, then an atom $\mathbf{q}_f(x, y)$ must be in $res(\mathcal{D}_n)$. As $\text{Goal} \notin res(\mathcal{D}_0)$, by Proposition 27, x or y must be a variable. Thus, by Lemma 34, x and y are on the same chain. In addition, as M does not halt on the empty word, by Lemma 32, x and y cannot be on the starting chain, so they are both on a non-starting chain. By Lemma 33, $SkD(x)$ and $SkD(y)$ are both less than one. Thus, by the contrapositive of Proposition 27, $res(\mathcal{D}_n)$ only contains terms of Skolem depth less than one.

If $R = R_{q,c}^0$, then atoms $\mathbf{q}_0(x, x)$, $\mathbf{q}(x, y)$ and $\mathbf{c}(x, y)$ must be in $res(\mathcal{D}_n)$ (with $\delta(q, c) = (q', c', \leftarrow)$ for some q', c'). As $\text{Goal} \notin res(\mathcal{D}_0)$, by Proposition 27, x or y must be a variable. Thus, by Lemma 34, x and y are on the same chain. In addition, as M does not halt on the empty word, by Lemma 32, x and y cannot be on the starting chain, so they are both on a non-starting chain. By Lemma 33, $SkD(x)$ and $SkD(y)$ are both less than one. Thus, by the contrapositive of Proposition 27, $res(\mathcal{D}_n)$ only contains terms of Skolem depth less than one.

□

Theorem 17. *The pair $(\mathcal{R}_M, \text{Goal})$ is Datalog-expressible.*

Proof. If M halts on the empty word, then by Theorem 13, $(\mathcal{R}_M, \emptyset) \models \text{Goal}$. Thus, by monotonicity of first-order entailment, for all databases D , $(\mathcal{R}_M, D) \models \text{Goal}$, so $(\{\rightarrow \text{Goal}\}, \text{Goal})$ is a suitable Datalog rewriting.

If M does not halt on the empty word, then by Theorem 35, it has bounded depth, so by Theorem 20, it is Datalog expressible. \square

We can finally prove the main result.

Theorem 3. *Query answering over Datalog-expressible rule sets is undecidable.*

Proof. Assume for a contradiction that query answering over Datalog-expressible rule sets is decidable, and consider a Turing machine M . Then, by Theorem 17, the couple $(\mathcal{R}_M, \text{Goal})$ is Datalog-expressible. Thus, we can decide if $(\mathcal{R}_M, \emptyset) \models \text{Goal}$. By Theorem 13, $(\mathcal{R}_M, \emptyset) \models \text{Goal}$ if and only if M halts on the empty word. Thus, we are able to decide the halting problem, which is absurd. \square

Bibliography

- [1] C. Beeri and M. Vardi. The implication problem for data dependencies. In *Proc. of Automata, Languages and Programming, Eighth Colloquium (ICALP 1981)*, volume 115 of *Lecture Notes in Computer Science*, pages 73–85, 1981.
- [2] Alin Deutsch, Alan Nash, and Jeff Remmel. The chase revisited. In *Proceedings of the Twenty-Seventh ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems, PODS '08*, pages 149–158, New York, NY, USA, 2008. Association for Computing Machinery.
- [3] Mélanie König, Michel Leclere, Marie-Laure Mugnier, and Michaël Thomazo. Sound, complete and minimal ucq-rewriting for existential rules. *Semantic Web*, 6, 11 2013.
- [4] Bruno Marnette. Resolution and Datalog rewriting under value invention and equality constraints. *CoRR*, abs/1212.0254, 2012.