# A Practical Dynamic Programming Approach to Datalog Provenance Computation

## Yann Ramusat ✉ 🔗
DI ENS, ENS, CNRS, PSL University & Inria, France

## Silviu Maniu ✉ 🔗
Université Paris-Saclay, LISN, CNRS, France

## Pierre Senellart ✉ 🔗
DI ENS, ENS, CNRS, PSL University & Inria & IUF, France

## ─── Abstract ───

We establish a translation between a formalism for dynamic programming over hypergraphs and the computation of semiring-based provenance for Datalog programs. The benefit of this translation is a new method for computing provenance for a specific class of semirings. Theoretical and practical optimizations lead to an efficient implementation using SOUFFLÉ, a state-of-the-art Datalog interpreter. Experimental results on real-world data suggest this approach to be efficient in practical contexts, even competing with our previous dedicated solutions for computing provenance in annotated graph databases. The cost overhead compared to plain Datalog evaluation is fairly moderate in many cases of interest.

## 1 Introduction

A notion of *provenance* for Datalog queries was introduced by Green, Karvounarakis, and Tannen [13]. It is based on the algebraic structure of *semirings* to encode additional meta-information about query results and extends the notion of semiring provenance of the positive fragment of the relational algebra, also introduced in [13]. The *full provenance* of a Datalog program (i.e., the provenance associated to each *derived tuple*) is expressed as a *system of equations* over an $\omega$-*continuous* semiring.

Beyond this initial definition, some research has proposed other representation frameworks to improve the computation of Datalog provenance. Circuits [6] have been proposed to obtain a compact representation of the provenance. Derivation-tree analysis [9] focuses on the natural mapping between systems of equations over $\omega$-continuous semirings and context-free grammars. These techniques rely on additional properties of the semiring, thus being limited to some specific applications; in particular, some semiring properties such as *absorptivity* serve as a basis for optimization in both scenarios.

The present work addresses two major challenges related to the computation of semiring-based provenance of Datalog programs.

On the first hand, much of the knowledge around the theory and applications of using semirings to capture some meta-information about a computation (i.e., *provenance* in database speak) is spread across different fields of computer sciences. Among other areas, semirings

have been successfully used in linguistic structure prediction [28], dynamic programming [14], constraint programming [3], and variations of the shortest-distance problem in graphs [23]. The structure of semirings is also at the heart of the study of formal languages [19, 20], weighted automata [8], graph theory and graph algorithms [22], and computational algebra [1]. The *algebraic path problem* [26] provides a common viewpoint for a large variety of problems involving semirings. Because of the widespread use of semirings in different areas, and rather unfortunately, it is common to find different names used for similar properties (e.g., 0-closed semirings are also called absorptive, bounded, simple, 0-stable, etc.). Even more concerning, the term "semiring" itself is replaced by similar names in some works (e.g. by "dioids" or "algèbre de chemin" [12]); in some cases, the semiring structure could have been used to model the problem, but was not acknowledged at the time [17]. We believe it is necessary to spend some effort establishing links between multiple representation frameworks of what is essentially semiring provenance to benefit from a wide range of techniques proposed in various domains of computer science. This would obviate the pitfall of designing algorithms and developing techniques in the jargon of a given field, whereas some key ideas may have already been investigated under similar algebraic frameworks, but by other research communities.

On the other hand, to the best of the knowledge of the authors, the bulk of the scientific literature about provenance for Datalog programs is mostly concerned with, either, formalization and extensions [16] to the theoretical framework behind semiring-based provenance, with the objective to reach a massive number of application domains, or elegantly leveraging quite abstract tools such as circuits [6] or Newton convergence [9] to optimize computations. Those works rarely go all the way to the implementation level, apart from some "proofs of concept" [10], raising the question of the applicability to real-world data.

A work of note that does come with an implementation is [4, 5], concerned about the computation of a *relevant* subset of provenance information through selection criteria based on tree patterns and the ranking of rules and facts, when the computation of the *full* provenance (how-provenance) is unfeasible in practice.

Our contributions can be organized into four parts.

**1.** We first establish a correspondence between dynamic programming over hypergraphs (as introduced in [14] under the name of *weighted hypergraphs*) and the *proof-theoretic* definition of the provenance for Datalog programs. We provide both-way translations and characterize for which class of semirings the *best-weight derivation* in the hypergraph corresponds to the provenance of the initial Datalog program.

**2.** The translation we thus introduced permits us to obtain a version of Knuth's generalization of Dijkstra's algorithm [17] to the *grammar problem*, adapted to the case of Datalog provenance computation. In the special setting where all hyperedges are of arity 1, we obtain the classical notion of semiring-based provenance for graph databases [24]. In the general setting, the algorithm steadily generalizes to Datalog the adapted Dijkstra's algorithm we presented in [24], under the same assumptions on the properties of the underlying semiring.

**3.** Such algorithm is unlikely to be efficient as-is in practical contexts. The main issue is closely related to the inefficiency of basic Datalog evaluation: many computations of facts (provenance values) have already been deduced, leading to redundant computations. We show that the *semi-naïve* evaluation strategy for Datalog is also applicable in our setting. An added advantage is that it greatly facilitates the extension of existing Datalog solvers to compute provenance annotations. We implement our strategy extending SOUFFLÉ [15], a state-of-the-art Datalog solver. We apply our solution to process rich graph queries (translated into Datalog programs) on several real-world and synthetic graph datasets,

as well as to Datalog programs used in previous works [5]. The performance of the implementation competes with dedicated solutions specific to graph databases [25].

4. The link we established with the framework of weighted hypergraphs, and hence, with the grammar problem, opens up new opportunities and challenges for semiring-based provenance for Datalog programs. We thus carefully discuss selected topics and problems that can now understood as closely linked to the provenance framework [13].

We start by introducing in Section 2 basic concepts on semirings and we recall the definition of provenance for Datalog programs. We formulate and prove in Section 3 the correspondence between weighted hypergraphs and semiring-based provenance for Datalog programs. In Section 4, we present the adapted version of Knuth's algorithm for the grammar problem and discuss theoretical aspects of its optimization. We then dive into the practical aspects of its implementation using SOUFFLÉ, and present in Section 5 experimental results witnessing the efficiency of our approach for practical scenarios. We finally discuss, in Section 6 and Section 7, respectively opportunities raised by our translation and related work. For space reasons, proofs and some auxiliary material are deported to an appendix.

## 2 Background

In the following, we recall basic concepts of semiring theory underlying the optimization techniques we provide in this paper. For more background on the theory and applications of semirings, examples of relevant semirings, as well as references to the literature on advanced notions of semiring theory, see our previous work [25]. We follow the definitions in [25] and highlight notions that occur under different names depending on the application domain.

▶ **Definition 1** (Semiring). *A* semiring *is an algebraic structure* $(S, \oplus, \otimes, \bar{0}, \bar{1})$ *where* $S$ *is some set,* $\oplus$ *and* $\otimes$ *are binary operators over* $S$, *and* $\bar{0}$ *and* $\bar{1}$ *are elements of* $S$, *satisfying the following axioms:*

- $(S, \oplus, \bar{0})$ *is a* commutative monoid*:* $(a \oplus b) \oplus c = a \oplus (b \oplus c)$, $a \oplus b = b \oplus a$, $a \oplus \bar{0} = \bar{0} \oplus a = a$;
- $(S, \otimes, \bar{1})$ *is a* monoid*:* $(a \otimes b) \otimes c = a \otimes (b \otimes c)$, $\bar{1} \otimes a = a \otimes \bar{1} = a$;
- $\otimes$ *distributes over* $\oplus$*:* $a \otimes (b \oplus c) = (a \otimes b) \oplus (a \otimes c)$;
- $\bar{0}$ *is an annihilator for* $\otimes$*:* $\bar{0} \otimes a = a \otimes \bar{0} = \bar{0}$.

A semiring is *commutative* if for all $a, b \in S$, $a \otimes b = b \otimes a$. A semiring is *idempotent* if for all $a \in S$, $a \oplus a = a$. For an idempotent semiring we can introduce the *natural order* defined by $a \leqslant b$ iff $a \oplus b = a$.[1] Note that this order is compatible with the two binary operations of the semiring: for all $a, b, c \in S$, $a \leqslant b$ implies $a \oplus c \leqslant b \oplus c$ and $a \otimes c \leqslant b \otimes c$. This is also called the *monotonicity* property.

An important property is that of *k-closedness* [23], i.e., a semiring is *k*-closed if: $\forall a \in S$, $\bigoplus_{i=0}^{k+1} a^i = \bigoplus_{i=0}^{k} a^i$. Here, by $a^i$ we denote the repeated application of the $\otimes$ operation $i$ times, i.e., $a^i = \underbrace{a \otimes a \otimes \cdots \otimes a}_{i}$. 0-closed semirings (i.e., those in which $\forall a \in S, \bar{1} \oplus a = \bar{1}$) have also been called *absorptive*, *bounded*, or *simple* depending on the literature. Note that any 0-closed semiring is idempotent (indeed, $a \oplus a = a \otimes (\bar{1} \oplus \bar{1}) = a \otimes \bar{1} = a$) and therefore admits a natural order.

---

[1] There are unfortunately two definitions of natural order commonly found in the literature; we use here that found in [23, 14] which matches the standard order on the tropical semiring; other works [19, 13, 25] define it as the reverse order. Our choice obviously has some impacts: in particular, when defining Datalog provenance, we need greatest fixpoints in lieu of the least fixpoints used in [19, 13].

Huang [14] introduces the notion of *superiority* of a semiring $S$ with respect to a partial order $\leqslant$, defined by: $\forall a, b \in S\ a \leqslant a \otimes b,\ b \leqslant a \otimes b$. The natural order satisfies this notion for 0-closed semirings:

▶ **Lemma 2.** *Let $S$ be an idempotent semiring and $\leqslant$ the natural order over $S$. Then $S$ is superior with respect to $\leqslant$ if and only if $S$ is 0-closed.*

An easier way of understanding natural order in 0-closed semirings is to note that for any idempotent semiring $S$, $\bar{0}$ is the greatest element ($\forall a \in S,\ a \oplus 0 = a \leqslant \bar{0}$) while, if the semiring is also 0-closed (i.e., *bounded*), $\bar{1}$ is the smallest ($\forall a \in S,\ \bar{1} \oplus a = \bar{1} \geqslant a$). Thus a bounded semiring $S$ verifies $\bar{1} \leqslant a \leqslant \bar{0}$ for all $a \in S$.

▶ **Definition 3** ($\omega$-Continuous semiring). *A semiring $(S, \oplus, \otimes, \bar{0}, \bar{1})$ is $\omega$-continuous if*

1. $(S, \geqslant)$ *is a $\omega$-complete partial order, i.e., the infimum $\inf\limits_{i \in \mathbb{N}} a_i$ of any infinite chain $a_0 \geqslant a_1 \geqslant \ldots$ exists in $S$.*

2. *both addition and multiplication are $\omega$-continuous in both arguments, i.e., for all $a \in S$ and infinite chain $a_0 \geqslant a_1 \geqslant \ldots$, $(a \oplus \inf\limits_{i \in \mathbb{N}} a_i) = \inf\limits_{i \in \mathbb{N}}(a \oplus a_i)$, $(a \otimes \inf\limits_{i \in \mathbb{N}} a_i) = \inf\limits_{i \in \mathbb{N}}(a \otimes a_i)$, $(\inf\limits_{i \in \mathbb{N}} a_i) \otimes a = \inf\limits_{i \in \mathbb{N}}(a_i \otimes a)$.*

In such semirings we can define countable sums: $\bigoplus\limits_{n \in \mathbb{N}} a_n = \inf\limits_{m \in \mathbb{N}} \bigoplus\limits_{i=0}^{m} a_i$.
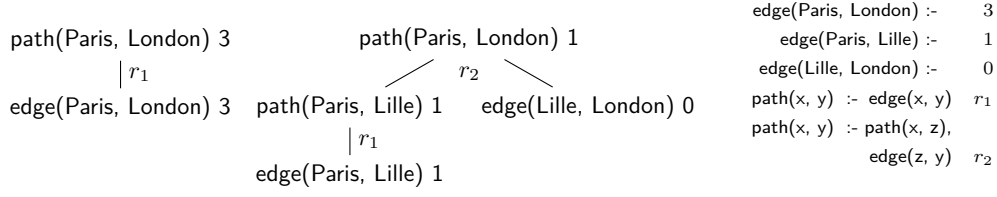
A *system of fixpoint equations* over an $\omega$-continuous semiring $S$ is a finite set of equations: $X_1 = f_1(X_1, X_2, \ldots, X_n), \ldots, X_n = f_n(X_1, X_2, \ldots, X_n)$, where $X_1, \ldots, X_n$ are variables and $f_1, \ldots f_n$ are polynomials with coefficients in $S$. We extend the notion of natural order from semiring elements to tuples of semiring elements by simply considering the product order. We then have the following on solutions of a system of equations over an $\omega$-continuous semiring:

▶ **Theorem 4** (Theorem 3.1 of [19]). *Every system of fixpoint equations $\mathbf{X} = f(\mathbf{X})$ over a commutative $\omega$-continuous semiring has a greatest solution $\mathsf{gfp}(f)$ w.r.t. $\leqslant$, and $\mathsf{gfp}(f)$ is equal to the infimum of the Kleene sequence: $\mathsf{gfp}(f) = \inf\limits_{m \in \mathbb{N}} f^m(\bar{\mathbf{0}})$.*

We now recall some basics about the Datalog query language and refer to [2] for more details. A *Datalog rule* is of the form $R(\vec{x}) \coloneq R_1(\vec{x_1}), \ldots, R_n(\vec{x_n})$ with $R$'s representing relations of a given arity and the $\vec{x}$'s tuples of variables of corresponding arities. Variables occurring on the left-hand side, the *head* of the rule, are required to occur in at least one of the atoms on the right-hand side, the *body* of the rule. A *Datalog program* is a finite set of Datalog rules. We call *fact* a rule with an empty body and variables replaced by constants. We divide relations into *extensional* ones (which can only occur as head of a fact, or in rule bodies) and *intensional* ones (which may occur as heads of a non-fact rule). The set of extensional facts is called the *extensional database* (EDB). We distinguish one particular relation occurring in the head of a rule, the *output predicate* of the Datalog program. We refer to [2] for the semantics of such a program and the notion of *derivation tree*.

There are two ways of defining the *provenance* of a Datalog program $q$ with output predicate $G$ over an $\omega$-continuous semiring. We can first base this definition on the proof-theoretic definition of standard Datalog:

▶ **Definition 5** (Proof-theoretic definition for Datalog provenance [13]). *Let $(S, \oplus, \otimes, \bar{0}, \bar{1})$ be a commutative $\omega$-continuous semiring and $q$ a Datalog program with output predicate $S$ and such that all extensional facts $R(t')$ are annotated with an element of $S$, denoted as $\mathsf{prov}_R^q(t')$.*

path(Paris, London) 3          path(Paris, London) 1          edge(Paris, London) :-        3
         $| r_1$                              $r_2$                    edge(Paris, Lille) :-        1
edge(Paris, London) 3   path(Paris, Lille) 1   edge(Lille, London) 0   edge(Lille, London) :-        0
                                  $| r_1$                    path(x, y)  :- edge(x, y)    $r_1$
                         edge(Paris, Lille) 1                 path(x, y)  :- path(x, z),
                                                                          edge(z, y)   $r_2$

**Figure 1** Derivation trees along their weights for the fact path(Paris, London) using the transitive closure Datalog program over the tropical semiring with an EDB containing 3 facts

*Then the* provenance *of $G(t)$ for $q$, where $G(t)$ is in the output of $q$, is defined as:*

$$\mathsf{prov}_G^q(t) = \bigoplus_{\tau \text{ yields } t} \left( \bigotimes_{t' \in \text{ leaves}(\tau)} \mathsf{prov}_R^q(t') \right).$$

The first sum ranges over all the derivation trees of the fact $t$ (see Figure 1 for examples of derivation trees), the second sum ranges over all leaves of the tree (extensional facts). This definition describes how the provenance propagates across the deduction process given an initial assignment of provenance weights to the extensional relations of $q$, $\mathsf{prov}_R^q$.

▶ **Example 6.** The *tropical semiring* is $(\mathbb{R}^+ \cup \{\infty\}, \min, +, \infty, 0)$. We show in Figure 1 an example Datalog program (right) with tropical semiring annotations on extensional facts, as well as the (only) two derivation trees of the fact path(Paris, London) along their weight. This witnesses that the provenance of path(Paris, London) is $\min(1, 3) = 1$.

Since some tuples can have infinitely many derivations, the Datalog semantics given above cannot be used as an algorithm. As pointed out in [13] it is possible instead to use a fixpoint-theoretic definition of the provenance of a Datalog query $q$: introduce a fresh variable for every possible *intensional* tuple (i.e., every possible *ground atom*), and produce for this variable an equation that reflects the *immediate consequence operator $T_q$ – extensional* facts appearing as their semiring annotations in these equations. This yields a system of fixpoint equation $f_q$. The provenance of $G(t)$ for $q$ is now simply the value of the variable corresponding to $G(t)$ in $\mathsf{gfp}(f_q)$.

The fixpoint-theoretic definition directly yields an algorithm, albeit a very inefficient one because of the need of generating a rule for every intensional tuple. In this paper, we investigate more efficient algorithms, for specific types of semirings.

## 3   Datalog provenance and dynamic programming over hypergraphs

We now show how to convert a Datalog program into a weighted hypergraph (as introduced in [14]) and characterize the semirings where the best-weight derivation in the hypergraph corresponds to the provenance for the initial Datalog program, mimicking the proof-theoretic definition. We first recall basic definitions and notation related to hypergraphs.

▶ **Definition 7** (Weighted hypergraph [14]). *Given a semiring $S$, a* weighted hypergraph *on $S$ is a pair $H = \langle V, E \rangle$, where $V$ is a finite set of vertices and $E$ is a finite set of hyperedges, where each element $e \in E$ is a triple $e = \langle h(e), \mathrm{T}(e), f_e \rangle$ with $h(e) \in V$ its head vertex, $\mathrm{T}(e) \in V$ an ordered list of tail vertices and $f_e$ a weight function from $S^{|\mathrm{T}(e)|}$ to $S$.*

We note $|e| = |\mathrm{T}(e)|$ the *arity* of a hyperedge. If $|e| = 0$, we say $e$ is nullary and then $f_e()$ is a constant, an element of the semiring; we assume there exists at most one nullary edge for a given vertex. In that case, $v = h(e)$ is called a *source vertex* and we note $f_e()$ as $f_v$. The *arity of a hypergraph* is the maximum arity of any hyperedge.

The *backward-star* $\mathrm{BS}(v)$ of a vertex $v$ is the set of incoming hyperedges $\{e \in E \mid h(e) = v\}$. The *graph projection* of a hypergraph $H = \langle V, E \rangle$ is a directed graph $G = (V, E')$ where $E' = \{(u, v) \mid \exists e \in \mathrm{BS}(v), u \in \mathrm{T}(e)\}$. A hypergraph $H$ is acyclic if its graph projection $G$ is acyclic; then a topological ordering of $H$ is an ordering of $V$ that is a topological ordering of $G$.

With these definitions in place, we can encode a Datalog program with semiring annotations as a weighted hypergraph in a straightforward manner:

▶ **Definition 8** (Hypergraph representation of a Datalog program). *Given a Datalog program $q$ as a set of rules $\{q_1, \cdots, q_n\}$ and the semiring $S$ used for annotations, we define the* weighted hypergraph representation *of $q$ as $H_q = \langle V_q, E_q \rangle$ with $V_q$ being all ground atoms and, for each instantiation of a rule $t(\vec{x}) \leftarrow r_1(\vec{x_1}), \ldots, r_n(\vec{x_n})$, a corresponding edge $\langle t(\vec{x}), (r_1(\vec{x_1}), \ldots, r_n(\vec{x_n})), \otimes \rangle$. For a fact $R(\vec{x}) \in \mathsf{EDB}(q)$ we add a nullary edge $e$ with $h(e) = R(\vec{x})$ and $f_e = \mathsf{prov}_R^q(\vec{x})$.*

The notion of *derivations* is the hypergraph counterpart to paths in graph. We recall the definition of derivations and we define it in a way that is reminiscent of Datalog-related notions.

▶ **Definition 9** (Derivation in hypergraph [14]). *We recursively define a* derivation $D$ *of a vertex $v$ in a hypergraph $H$ (as a pair formed of a hyperedge and a list of derivations), its size $|D|$ (a natural integer) and its weight $w(D)$ (a semiring element) as follows:*

- *If $e \in \mathrm{BS}(v)$ with $|e| = 0$, then $D = \langle e, \langle \rangle \rangle$ is a derivation of $v$, $|D| = 1$, and $w(D) = f_e()$.*
- *If $e \in \mathrm{BS}(v)$ with $|e| \geqslant 0$, $D_i$ is a derivation of $T_i(e)$ for $i = 1 \ldots |e|$, then $D = \langle e, \langle D_1 \cdots D_{|e|} \rangle \rangle$ is a derivation of $v$, $|D| = 1 + \sum_{i=1}^{|e|} |D_i|$, $w(D) = f_e(w(D_1), \ldots, w(D_{|e|}))$.*

*We note $\mathcal{D}_H(v)$ the set of derivations of $v$ in $H$.*

When modeling Datalog provenance in a semiring $S$ as weighted hypergraphs on $S$, all non-source weight functions are bound to the $\otimes$ operation of the semiring. Note that, if $S$ is idempotent, the natural order on $S$ induces an ordering on derivations: $D \leqslant D'$ if $w(D) \leqslant w(D')$.

We now show that in this formalism, the Datalog provenance of an output predicate can be understood as the best-weight for the corresponding vertex in the hypergraph.

▶ **Definition 10** (Best-weight [14]). *The* best-weight $\delta_H(v)$ *of a vertex $v$ of a hypergraph $H$ on a semiring $(S, \oplus, \otimes, \bar{0}, \bar{1})$ is the weight of the best derivation of $v$:*

$$
\delta_H(v) = \begin{cases} f_v & \text{if } v \text{ is a source vertex;} \\ \bigoplus_{D \in \mathcal{D}_H(v)} w(D) & \text{otherwise.} \end{cases}
$$

The best-weight generally requires additional properties of either the hypergraph or the semiring to be well-defined. Acyclicity for the hypergraph is a sufficient condition. Existence of an infinitary summation operator in the semiring extending $\oplus$, guaranteed in $\omega$-continuous semirings, is also a sufficient condition. To guarantee semantics compatible with the intuitive meaning of provenance, a more restrictive sufficient condition is for the semiring to be a *c-complete star-semiring* [18], see [25] for details.

We can now show that Datalog provenance can be computed through the formalism of weighted hypergraphs. Let us start with a lemma exhibiting a one-to-one mapping between derivations in the hypergraph and proofs in Datalog.

▶ **Lemma 11.** *For any Datalog query $q$ and grounding of an atom $v$ of $q$, there is a bijection between $\mathcal{D}_{H_q}(v)$ and $\{\tau \mid \tau \text{ yields } v\}$.*

We then show that the weight of each derivation of a tuple is equal to the corresponding proof tree weight in Datalog.

▶ **Lemma 12.** *For any Datalog query $q$ and grounding of an atom $v$ of $q$, for any derivation $D$ of $v$ in $H_q$ $w(D) = \bigotimes\limits_{t' \in \text{leaves}(\tau_D)} \text{prov}_R^q(t')$ where $\tau_D$ is the proof tree corresponding to $D$ in the bijection given by Lemma 11.*

Finally, we obtain:

▶ **Theorem 13.** *Let $t$ be a tuple of a Datalog program $q$ with output predicate $G$ and $H_q$ its hypergraph representation, then $\text{prov}_G^q(t) = \delta_{H_q}(G(t))$.*

## 4    Best-first method

Knuth [17] generalized the Dijkstra algorithm to what he calls the *grammar problem* (i.e., finding the *best-weight derivation* from a given non-terminal, where each terminal has a specific weight and each rule comes with an associated weight function). This has been identified as corresponding to the search problem in a monotonic superior hypergraph – for each $e \in H$, $f_e$ is monotone and superior in each argument (see Table 3 in [14]). We showed in Lemma 2 that superiority corresponds to 0-closedness in semirings with natural orders. The definition of the grammar problem assumes a total order on weights as the weights are real numbers. In the special case where all hyperedges are of arity 1 (and all weight functions bound to $\otimes$), we obtain the classical notion of semiring-based provenance for graph databases [24]. Thus, Knuth's algorithm can be seen as a generalization to hypergraphs (and therefore, by the results of the previous section, to Datalog provenance computation) of the modified Dijkstra algorithm we presented in [24], working on 0-*closed totally-ordered semirings*, which are generalizations of the tropical semiring.

### Optimized version of Best-first method

In the original paper of Knuth [17], the question of efficient construction of the set of candidate facts for the extraction of the minimal-valued fact is not dealt with. A lot of redundant work may be carried out if the implementation is not carefully designed.

In the following, we show how to obtain a ready-to-be-implemented version incorporating ideas from the *semi-naïve* evaluation of Datalog programs. Semi-naïve evaluation of Datalog, as described in [2, Chapter 13] introduces a number of ideas aiming at improving the efficiency of the *naïve* Datalog evaluation method; we show how to leverage these tricks in our setting.

The *naïve* evaluation of a Datalog program $q$ processes iteratively, applying at each step the *consequence operator* $T_q$. Many redundant derivations are computed, leading to practical inefficiency. The *semi-naïve* evaluation addresses this problem by considering only facts derived using at least one new fact found at the previous step. Note, however, whereas many new facts can be found at one step of the semi-naïve evaluation, only one is to be added by the Best-first method to respect the $\leqslant$-minimality ordering of added facts.

■ **Algorithm 1** Basic semi-naïve version of Best-first method for Datalog provenance

---

**Input:** Datalog query $q$, EDB $D$ with provenance indications over a 0-closed totally-ordered semiring $S$.

**Output:** full Datalog provenance for the IDB of $q$.

1: **function** RELAX($r_0(\vec{x_0})$, $S$)
2:     **for** each instantiation of a rule $r(\vec{x}) \leftarrow r_1(\vec{x_1}), \ldots, r_m(\vec{x_m}), \cdots, r_n(\vec{x_n})$ where $r_i(\vec{x_i}) \in D \cup S \cup \{r_0(\vec{x_0})\}$, $1 \leqslant i < m$, $r_m(\vec{x_m}) = r_0(\vec{x_0})$ and $r_i(\vec{x_i}) \in D \cup S$, $m < i \leqslant n$ **do**
3:         $\nu(r(\vec{x})) \quad \oplus= \bigotimes\limits_{1 \leqslant i \leqslant n} r_i(\vec{x_i})$

4:
5: $I \leftarrow \emptyset$
6: Let $\nu$ be the function that maps all facts of $D$ to their annotation in $S$ and all potential facts of the intensional schema of $q$ to $\bar{0}$
7: **for** each intensional atom $r(\vec{x}) \notin I$ **do**
8:     $\nu(r(\vec{x})) \quad \oplus= \bigotimes\limits_{1 \leqslant i \leqslant n} r_i(\vec{x_i})$ for each instantiation of a rule
        $r(\vec{x}) \leftarrow r_1(\vec{x_1}), \ldots, r_n(\vec{x_n})$ with $r_i(\vec{x_i}) \in D$
9: **while** $\min_{\nu \setminus I} r(\vec{x}) \neq \bar{0}$ **do**
10:     Add such minimal $r(\vec{x})$ to $I$
11:     RELAX($r(\vec{x})$, $I \setminus r(\vec{x})$)
12: **return** $\nu$

---

This algorithm starts by initializing the priority queue with IDB facts that are derivable from EDB facts. Then, at each step, the minimum valued-fact is added, and only derivations using this new fact are computed to update the value of the facts in $I$. This algorithm stops whether: 1. the maximal value is reached for a candidate fact, 2. the list is empty - the minimal value of the list is by default the maximal value of the semiring.

▶ **Theorem 14.** *Algorithm 1 computes the full Datalog provenance for* 0*-closed totally-ordered semirings.*

**Proof.** We show the algorithm verifies the following invariant: whenever a tuple is added to $I$ in Line 10, it has optimal value. This implies that $I$ is populated in increasing order: each new derivation computed in the RELAX() procedure only updates the priority queues with values greater than the value of the tuple relaxed (by superiority of $\otimes$).

Suppose by contradiction that some output tuples are not correctly labeled and take such a minimal tuple $\nu = r(\vec{x})$. At the moment where $\nu$ is extracted with value $n$ let us consider an optimal derivation path of $\nu$ that leads to the optimum value $opt < n$. By superiority each tuple occurring in the tail of the rule has value less than $opt$. Thus a tuple occurring in the tail is either wrong-valued or not present in $I$ at the moment where $\nu$ is found. In both cases and because tuples are added to $I$ in increasing order we obtain a new minimal tuple incorrectly labeled by the algorithm, contradicting the hypothesis. ◀

The structure of the Datalog program can be analysed to provide clues about the predicates to focus on. Following [2], we introduce the notion of *precedence graph* $G_P$ of a Datalog program $P$. The nodes are the IDB predicates and the edges are pairs of IDB predicates $(R, R')$ where $R'$ occurs at the head of a rule of $P$ with $R$ belonging to the tail. $P$ is a *recursive* program if $G_P$ has a directed cycle. Two predicates $R$ and $R'$ are mutually recursive if $R = R'$ or $R$ and $R'$ participate in the same cycle of $G_P$. This defines equivalence classes. Putting it together, we obtain as Algorithm 2 out final algorithm.

◼ **Algorithm 2** Semi-naïve version of Best-first method for Datalog provenance

---

**Input:** $q$ a Datalog query with provenance indication over a 0-closed totally-ordered semiring $S$.
**Output:** full Datalog provenance for the IDB of $q$.
 1: Compute the equivalence classes of $q$
 2: **for** each equivalence class in a topological order **do**
 3:     Apply Algorithm 1 over IDB predicates in the equivalence class
 4:     considering previous equivalence classes as EDB predicates
 5: **return** $\nu$

---

#### Generalization to distributive lattices

In [25], we outlined a new algorithm, based on the Dijkstra algorithm and solving the single-source provenance in graph databases with provenance indications over 0-closed multiplicatively idempotent semirings (equivalents to distributive lattices). This new method stems from a tentative to bridge the strong complexity gap for computing the provenance in the case of a semiring not 0-closed and totally ordered. A similar gap also appears when we consider provenance over Datalog queries (see Section 7). Thus, we show how to apply this method for computing provenance for Datalog queries over distributive lattices.

We provide a brief review of the key ideas presented in [25]. Any element of a distributive lattice is decomposable into a product of *join-irreducible* elements of the lattice, and there exists an embedding of the distributive lattice into a chain decomposition of its join-irreducible elements. This ensures we can 1) work on a totally ordered environment and apply algorithms that require total ordering over the elements, 2) independently compute partial provenance annotations for each dimension to form the final provenance annotation. Given $m$ the number of dimensions in the decomposition, our solution (described in Algorithm 3) performs $m$ launches of the Best-first method and thus, roughly has a cost increased by a factor $m$.

## 5 Implementation and experiments

In numerous application domains, Datalog is used as a *domain specific language* (DSL) to express logical specifications for static program analysis. A formal specification, written as a *declarative* Datalog program is usually translated into an efficient *imperative* implementation by a *synthesizer*. This process simplifies the development of program analysis compared to hand-crafted solutions (highly optimized C++ applications specialized in enforcing a fixed set of specifications). SOUFFLÉ [15, 27] has been introduced to provide efficient synthesis of Datalog specifications to executable C++ programs, competing with state-of-the-art handcrafted code for program analysis. The inner workings of SOUFFLÉ were of interest to our work; the algorithm implementations are similar to the evaluation strategy followed by the Best-first method we introduced here. We present a brief overview of the architecture of SOUFFLÉ and discuss how we extended it.

#### Architecture and implementation

Following what is described in [15], an input datalog program $q$ goes through a staged specialization hierarchy. After parsing, the first stage of SOUFFLÉ specializes the semi-naïve evaluation strategy applied to $q$, yielding a relational algebra machine program (RAM). Such a program consists in basic relational algebra operations enriched with I/O operators and

---

■ **Algorithm 3** Generalized Best-first method for Datalog provenance

---

**Input:** $q$ a Datalog query with provenance indication over a 0-closed multiplicatively idempotent semiring $S$.

**Output:** full Datalog provenance for the IDB of $q$.

1: **for** each EDB fact $R(\vec{x})$ **do**
2:     DECOMPOSE($R(\vec{x})$)
3: **for** each dimension $i$ **do**
4:     $\nu_i \leftarrow$ BEST-FIRST($q, i$)
5: **return** RECOMPOSE($\nu_1, \ldots, \nu_n$)

---

■ **Algorithm 4** Input Datalog program computing the transitive closure (SOUFFLÉ syntax)

---

1: **.decl** edge(s:number, t:number[, @prov:semiring value])
2: **.input** edge
3: **.decl** path(s:number, t:number[, @prov:semiring value])
4: **.output** path
5: path(x, y) :- edge(x, y).
6: path(x, y) :- path(x, z), edge(z, y).

---

■ **Algorithm 5** Corresponding SOUFFLÉ RAM program for Algorithm 4

---

1: **if** $\neg(\text{edge} = \emptyset)$ **then**
2:     **for** t0 **in** edge **do**
3:         **add** (t0.0, t0.1) **in** path
4:         **add** (t0.0, t0.1) **in** $\delta$path
5: **loop**
6:     **if** $\neg(\delta\text{path} = \emptyset) \wedge \neg(\text{edge} = \emptyset)$ **then**
7:         **for** t0 **in** $\delta$path **do**
8:             **for** t1 **in** edge **on index** t1.0 = t0.1 **do**
9:                 **if** $\neg$(t0.0, t0.1) $\in$ path **then**
10:                     **add** (t0.0, t0.1) **in** path'
11:     **if** path' $= \emptyset$ **then**
12:         **exit**
13:     **for** t0 **in** path' **do**
14:         **add** (t0.0, t0.1) **in** path
15:     **swap** $\delta$path **with** path'
16:     **clear** path

---

■ **Algorithm 6** Modification of RAM program of Algorithm 5 to implement Best-first strategy

---

1: **if** $\neg(\text{edge} = \emptyset)$ **then**
2:     **for** t0 **in** edge **do**
3:         **update** (t0.0, t0.1, t0.prov) **in** path
4:     **for** t0 **in** path **do**
5:         **add** (t0.0, t0.1, t0.prov) **in** $\delta$path
6: **loop**
7:     **if** $\neg(\delta\text{path} = \emptyset) \wedge \neg(\text{edge} = \emptyset)$ **then**
8:         **for** t0 **in** $\delta$path **do**
9:             **for** t1 **in** edge **on index** t1.0 = t0.1 **do**
10:                 **if** $\neg$(t0.0, t1.1, $\bot$) $\in$ path **then**
11:                     **update** (t0.0, t0.1, t0.prov $\otimes$ t1.prov) **in** pq
12:     **clear** $\delta$path
13:     **if** pq is empty **then**
14:         **exit**
15:     **add** pq.top() **in** pq.top().relation **and in** pq.top().$\delta$relation

---

fix-point computations. As a final step, the RAM program is finally either interpreted or compiled into an executable. For this work, we have only used the interpreter. Our code was inserted in two different stages of SOUFFLÉ: a new *translation strategy* from the parsed program to the RAM program, a *priority queue*, replacing the code in charge of adding at run-time the tuples to the relations.

We showcase the result of our translation strategy in Algorithms 4, 5, and 6 for a Datalog query computing the transitive closure of a graph; this program is given in Algorithm 4 in its SOUFFLÉ syntax. Algorithm 5 presents the corresponding SOUFFLÉ RAM program resulting from applying the semi-naïve evaluation strategy and Algorithm 6 our modification to the RAM program to provide provenance annotation via the Best-first strategy and use the priority queue pq for provenance computation. The ⊥ notation corresponds to a wildcard. Importantly, modifying directly at the RAM level of SOUFFLÉ allows us to benefit of all implemented optimizations.

### Experiments

Our implementation was tested on an Intel Xeon E5-2650 computer with 176 GB of RAM. The source code is freely available on GITHUB[2].

The initial motivation for this work stems from a key observation we outlined in the conclusion of [25], where we pointed out the similarity between Datalog and the classes of semirings and their optimized provenance algorithms discussed in that work, focused on graph provenance algorithms. To translate this graph setting into Datalog, the graph structure has been encoded into an EDB with one binary predicate *edge* encoding the edges, and with edge notations depending on the provenance semiring we chose. We run the *transitive closure* Datalog program outlined in Algorithm 4. Full information on the graph datasets used can be found in [25]. We provide, in Figure 2, a comparison between the best-first method introduced here (SOUFFLÉ-PROV), the plain SOUFFLÉ without provenance computation, and a previous provenance-based algorithm [25] computing all-pairs shortest-distances over graph databases (the NODEELIMINATION algorithm, with a choice of node to eliminate based on its id or its degree), in the tropical semiring. Similarly, in Figure 3, we compare with previous solutions for single-source shortest-distances, in the same semiring, in particular the adaptation of DIJKSTRA algorithm of [25], and, for comparison purposes, a BFS algorithm that does not compute provenance. The main focus of this work was to provide an effective Datalog based solution for all-pairs provenance in graph databases. For the all-pairs problem, depending on the dataset, (see, e.g., YEAST), SOUFFLÉ-PROV is significantly faster than the previous best known algorithm, NODEELIMINATION. Unsurprisingly, BFS and DIJKSTRA perform respectively better than SOUFFLÉ and SOUFFLÉ-PROV in the single-source context. What favors both graph algorithms strongly is the fact that they reduce redundant computation: the algorithms abort whenever the target vertex has been reached. SOUFFLÉ-PROV performs between 1 and 2 orders of magnitude faster than MOHRI [23] – an algorithm designed for single-source provenance on $k$-closed semirings. This fact highlights the potential of adapting the best-first method to also handle $k$-closed semirings.

We now turn to evaluating the overhead induced by adding provenance computations to Datalog via SOUFFLÉ. This appears fairly modest in the experiments of Figures 2, 3. We further consider two datasets that go beyond the setting of graph databases: a non-recursive query over synthetic data (IRIS [4], obtained from the authors of that paper) in Figure 4 and
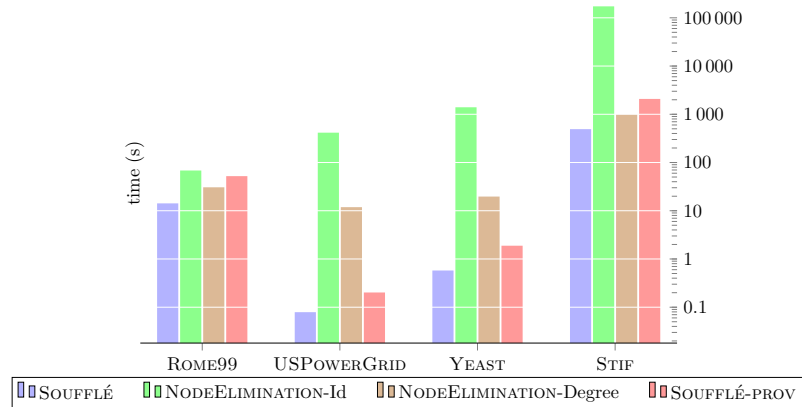
---

[2] `https://github.com/yannramusat/souffle-prov`

**Figure 2** Comparison between algorithms for all-pairs shortest-distances (Tropical). Values greater than 100 000 s are timeouts.
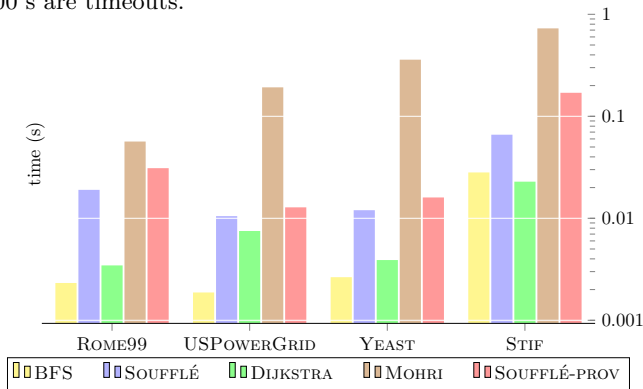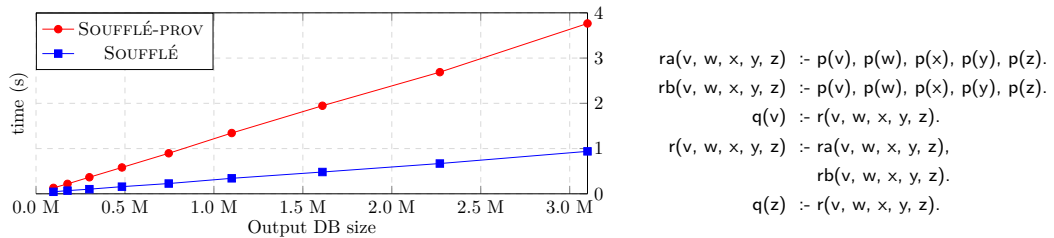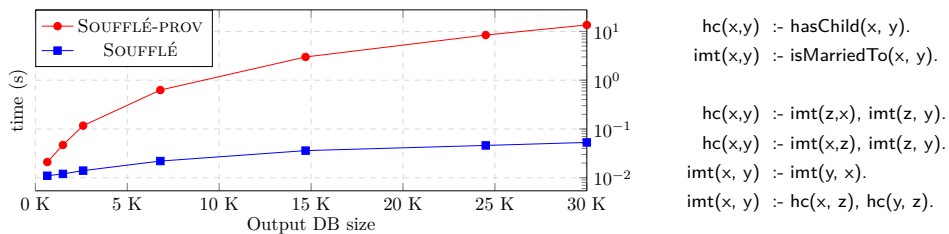


**Figure 3** Comparison between algorithms for single-source shortest-distances (Tropical).



ra(v, w, x, y, z)  :- p(v), p(w), p(x), p(y), p(z).
rb(v, w, x, y, z)  :- p(v), p(w), p(x), p(y), p(z).
         q(v)  :- r(v, w, x, y, z).
  r(v, w, x, y, z)  :- ra(v, w, x, y, z),
                       rb(v, w, x, y, z).
         q(z)  :- r(v, w, x, y, z).

**Figure 4** Comparison between SOUFFLÉ and SOUFFLÉ-PROV (IRIS)



hc(x,y)  :- hasChild(x, y).
imt(x,y)  :- isMarriedTo(x, y).

hc(x,y)  :- imt(z,x), imt(z, y).
hc(x,y)  :- imt(x,z), imt(z, y).
imt(x, y)  :- imt(y, x).
imt(x, y)  :- hc(x, z), hc(y, z).

**Figure 5** Comparison between SOUFFLÉ and SOUFFLÉ-PROV (AMIE)

a program over a knowledge base (AMIE [11], available online[3]) populated with real-world data and automatic rules in Figure 5. The evaluation was made by varying the output database size. In IRIS, the overhead induced is a modest constant factor of approximately 4 times the original cost. This is not the case for AMIE, where the overhead tends to increase with the size of the answer. We conjecture this is due to the lack of native support for tuple updates in SOUFFLÉ (in Datalog, the fixed-point operator is inflationary); this leads to inefficiency in updating provenance values associated to each tuple.

## 6 Reverse translations and opportunities

Our focus so far in this study was to express Datalog programs into hypergraphs. We now consider the opposite direction: given a weighted hypergraph $H$ with all its weight functions derived from a semiring $S$ via its $\otimes$ operator we aim to design an equivalent Datalog program which computes the provenance on the given semiring. We propose two translations – a straightforward one and one with a fixed program, to benefit of the low data complexity of Datalog.

▶ **Definition 15** (Weighted hypergraph to Datalog program). *Given an hypergraph $H = \langle V, E \rangle$ of maximal hyperedge arity $n$, we define a Datalog query $q_H$ with extensional predicates $E_n$ of arity $n + 1$ for $0 \leqslant i \leqslant n$ and an unary intensional predicate $R$. For each $1 \leqslant i \leqslant n$, we add a single rule of the form: $R(x) \leftarrow E_{i+1}(x, x_1, \ldots, x_n), R(x_1), \ldots, R(x_n)$, degenerating into $R(x) \leftarrow E_1(x)$ for $i = 0$. For each $e \in E$ of arity $n \geqslant 1$, we populate $E_{n+1}$ with $E_{n+1}(h(e), \mathrm{T}(e))$ having provenance value $\bar{1}$. Source vertices (i.e., heads of nullary edges) may have an initial constant value $s \in S$. For each $e \in E$ of arity $0$ we tag $E_1(h(e))$ with the provenance value $f_e$.*

We then have:

▶ **Theorem 16.** *Let $v$ be a vertex of a weighted hypergraph $H$ with all weight functions derived from the $\otimes$ operation of a semiring $S$ and $q_H$ its translation into a Datalog query, then $\delta_H(v) = \mathsf{prov}_R^{q_H}(v)$.*

Let us know consider the case of programs where we have a fixed schema and arity:

▶ **Definition 17** (Weighted hypergraph to Datalog program with fixed schema). *Given an hypergraph $H = \langle V, E \rangle$ let us consider the following Datalog program $q_{H_f}$ over unary predicates $R$ and Nullary, binary predicates $E$, $N$ and First and ternary predicate Next:*

$$R(x) \leftarrow E(x, e), H(e)$$
$$H(e) \leftarrow First(e, x), R(x), N(e, x) \qquad N(e, x) \leftarrow Next(e, x, y), R(y), N(e, y)$$
$$\leftarrow Nullary(e) \qquad\qquad\qquad \leftarrow End(e, x)$$

*For each $e \in E$ of arity $n \geqslant 1$, we populate $E$ with $E(h(e), e)$ with provenance value $\bar{1}$. Let $x_1, \ldots, x_n$ be the elements of $\mathrm{T}(e)$), we populate First with $First(e, x_1)$, Next with $Next(e, x_i, y_{i+1})$ for $1 \leqslant i \leqslant n - 1$ and End with $End(e, x_n)$; all have provenance value $\bar{1}$.*

*Source vertices (heads of nullary edges) may have an initial constant value $s \in S$, then for each $e \in E$ of arity $0$ we tag $E(h(e), e)$ with the provenance value $f_e$ and Nullary(e) with provenance value $\bar{1}$.*

---

[3] https://bitbucket.org/amirgilad/selp/src/Journal/

The reasoning follows a similar flow:

▶ **Theorem 18.** *Let $v$ be a vertex of a weighted hypergraph $H$ with all weight functions derived from the $\otimes$ operation of a semiring $S$ and $q_{H_f}$ its translation into a Datalog query with fixed schema, then $\delta_{H_f}(v) = \mathsf{prov}_R^{q_{H_f}}(v)$.*

### Case study: AND/OR graphs

Section 4.3 of [14] showcases some formalisms in which equivalent hypergraphs can be constructed in the dynamic programming framework proposed in their work. One advantage is that we can directly benefit from these translations in our Datalog provenance framework. However, there are some issues with these translations, as we now discuss.

One example formalism that can be easily translated to our setting are the AND/OR graphs, a special case of graphs – e.g., used in scheduling tasks – in which nodes can express two types of restrictions: AND restrictions (in which tasks have to executed sequentially) and OR restrictions (in which tasks can be overlapping). These two operations, not unlike the $\otimes$ and $\oplus$ operations on semirings, suggest our translation strategy is readily usable. Indeed, by our translation, vertices of the hypergraph would represent OR vertices and the hyperedges correspond to the AND vertices. This translation does not add any other restriction to the shape of the AND/OR graph. AND and OR vertices can have multiples outgoing edges and the AND/OR graph is not necessarily bipartite. However, not all applications in which AND/OR graphs are useful can be used with semiring operations. For instance, scheduling using AND/OR graphs [21, 7] requires three operations (`min`, `max` and $+$) instead of the two allowed using semirings.

Another crucial issue, not obvious at first glance, is that when expressing these graphs into our Datalog framework the semantics can change. As explained in [7], zero-edge cycles express mutual events (occurring in groups, with no causal relations): either they all occur together or none of them do. In a deductive system, where semantics is given by a least fixed-point, they are believed to not occur at all. The semantics of Datalog provenance or derivation trees cannot express this situation: in Datalog there would be no finite proof of a set of given facts. However, adding all of these facts to the solution would not break the deductive rules; simply they do not belong to the least fixed-point.

## 7 Related work

There are two major notions of provenance for information systems currently in use in the literature, each focusing on different usages. The first notion can be broadly categorized as "informational": the provenance encapsulates information about the deductive process leading to a specific result. Declarative debugging is one application of this concept. Analysis of provenance information in an interactive manner simplifies the user's burden of identifying which part of the rule specification is responsible for a *faulty* derivation. SOUFFLÉ contains its own integrated debugging system, based on a provenance-based evaluation strategy [29].

In the current study we consider a "computational" notion of provenance, where operations (and queries) over provenance values are permitted. In this case, the underlying properties of the semiring are important for optimized algorithms. For instance, *absorptivity* is an ubiquitous property in the literature allowing efficient algorithms for provenance evaluation. This applies not only for Datalog programs, but also for e.g., regular path queries over (semiring)-annotated graph databases [25].

With respect to Datalog provenance, it has been shown in [6] that, for a Datalog program having $n$ candidate IDB tuples, a circuit for representing Datalog provenance in the semiring **Sorp**$(X)$ (the most general absorptive semiring) only needs $n+1$ layers. For binary relations, e.g., representing the edge relation of a graph, this construction is at least quadratic in the number of vertices, thus not practically applicable for the graphs we analyzed in our experiments. Similarly, in [9], absorptive semirings (i.e., 0-closed semirings) have the property that derivation trees of size $\geqslant n$ are "pumpable" (they do not contribute to the final result). A concrete implementation [10] computes the provenance for commutative and idempotent semirings using $n$ Newton iterations.

Fairly recently, [16] introduced POPS (*Partially Ordered, Pre-Semiring*), a structure decoupling the order on which the fixed-point is computed from the semiring structure. Complex and recursive computations over vectors, matrices, tensors are now expressible using this framework. The study also generalized the semi-naïve method from plain Datalog evaluation to idempotent semirings (aka dioids). In comparison, our method is restricted to semirings that are totally ordered (a subclass of distributive dioids[4]), leveraging the invariant that once a fact is first labeled with a provenance value, we are certain it is the correct one.

In cases where keeping the *full* provenance of a program (*how*-provenance) is still prohibitively large, [5, 4] propose to select only a relevant subset of such trees using *selection criteria* based on tree patterns and ranking over rules and facts occurring in the derivation. First, given a Datalog program $P$ and a pattern $q$, an *offline* instrumentation is performed, leading to an *instrumented program $P_q$*. Then, given any database $D$, an efficient algorithm can be used to retrieve only the top-$k$ best derivation trees for $P_q(D)$. The top-1 algorithm of the study is closely related to our solution, but does not mention the use of a priority queue nor does it take into account the optimization provided by the semi-naïve evaluation strategy we describe in Section 4.

Our solution can be seen as a hybrid of the ideas introduced in [16] and [5]. We generalize the semi-naïve evaluation to a specific class of semirings in order to achieve a more efficient algorithm, one that can be used in practical real-world scenarios.

## 8 Conclusion

In this work, we developed a novel method for Datalog provenance computation based on the link between dynamic programming over hypergraphs and the proof structure of provenance of Datalog programs. We introduced Knuth's algorithm for computing the provenance, and optimized it for practical use. We showed its feasibility by providing an implementation on top of SOUFFLÉ and tested it on several graph databases and Datalog programs.

As a continuation of this work, we aim to establish a taxonomy of recent algorithms for provenance computation of Datalog programs (e.g. [16]), i.e., classify them by their data complexity and the class of semiring usable in each algorithm – with a focus on $k$-closed semirings. The objective is to provide a meaningful comparison with the already established taxonomy for semiring provenance in graph databases [25]. Another direction of interest would be to investigate a "fix" to the problems raised by Datalog provenance in AND/OR graphs (see Section 6). Finally, in a more practical direction, it might be feasible to implement the semi-naïve evaluation for Datalog° over distributive dioids (Algorithm 1 in [16]) using SOUFFLÉ. Most of the inner workings can easily be adapted to do so, but extending the current data structures storing facts (B-Trie, B-Tree) to allow updates may be a challenge.

---

[4] Distributive dioids are POPS structures over a distributive lattice being the natural order of the dioid.

─── **References** ───

**1**   Kamal Abdali. Parallel computations in *-semirings. In Klaus G. Fischer, Philippe Loustaunau, Jay Shapiro, Edward L.Green, and Daniel Farkas, editors, *Computational Algebra*, chapter 1. Marcel Dekker, Inc., 1994.

**2**   Serge Abiteboul, Richard Hull, and Victor Vianu. *Foundations of Databases*. Addison-Wesley, 1995.

**3**   Stefano Bistarelli, Ugo Montanari, and Francesca Rossi. Semiring-based constraint satisfaction and optimization. *Journal of the ACM*, 44(2):201–236, March 1997. `doi:10.1145/256303.256306`.

**4**   Daniel Deutch, Amir Gilad, and Yuval Moskovitch. selP: selective tracking and presentation of data provenance. In *2015 IEEE 31st International Conference on Data Engineering*, pages 1484–1487, 2015. `doi:10.1109/ICDE.2015.7113407`.

**5**   Daniel Deutch, Amir Gilad, and Yuval Moskovitch. Efficient provenance tracking for datalog using top-k queries. *The VLDB Journal*, 27:245–269, 2018.

**6**   Daniel Deutch, Tova Milo, Sudeepa Roy, and Val Tannen. Circuits for Datalog Provenance. In *ICDT*, pages 201–212, 2014.

**7**   Yefim Dinitz, Paz Carmi, Shahar Golan, Omri Liba, and Guy Rozenwald. An $O(|V||E|)$ algorithm for scheduling with and/or precedence constraints. In *Haifa Workshop on Interdisciplinary Applications of Graphs, Combinatorics and Algorithms*, 2011.

**8**   Manfred Droste, Werner Kuich, and Heiko Vogler. *Handbook of Weighted Automata*. Springer Publishing Company, Incorporated, 1st edition, 2009.

**9**   Javier Esparza and Michael Luttenberger. Solving fixed-point equations by derivation tree analysis. In *International Conference on Algebra and Coalgebra in Computer Science*, pages 19–35. Springer, 2011.

**10**   Javier Esparza, Michael Luttenberger, and Maximilian Schlund. Fpsolve: A generic solver for fixpoint equations over semirings. In Markus Holzer and Martin Kutrib, editors, *Implementation and Application of Automata*, pages 1–15, Cham, 2014. Springer International Publishing.

**11**   Luis Galárraga, Christina Teflioudi, Katja Hose, and Fabian Suchanek. Amie: Association rule mining under incomplete evidence in ontological knowledge bases. In *WWW*, pages 413–422, 2013. `doi:10.1145/2488388.2488425`.

**12**   M. Gondran. Algèbre linéaire et cheminement dans un graphe. *RAIRO - Operations Research - Recherche Opérationnelle*, 9(V1):77–99, 1975.

**13**   Todd J. Green, Grigoris Karvounarakis, and Val Tannen. Provenance semirings. In *PODS*, 2007. `doi:10.1145/1265530.1265535`.

**14**   Liang Huang. Advanced dynamic programming in semiring and hypergraph frameworks. In *Coling 2008: Advanced Dynamic Programming in Computational Linguistics: Theory, Algorithms and Applications*, pages 1–18, 2008.

**15**   Herbert Jordan, Bernhard Scholz, and Pavle Subotić. Soufflé: On synthesis of program analyzers. In Swarat Chaudhuri and Azadeh Farzan, editors, *Computer Aided Verification*, pages 422–430, Cham, 2016. Springer International Publishing.

**16**   Mahmoud Abo Khamis, Hung Q. Ngo, Reinhard Pichler, Dan Suciu, and Yisu Remy Wang. Convergence of datalog over (pre-) semirings. *CoRR*, abs/2105.14435, 2021. `arXiv:2105.14435`.

**17**   Donald E. Knuth. A generalization of Dijkstra's algorithm. *Information Processing Letters*, 6(1), 1977. `doi:10.1016/0020-0190(77)90002-3`.

**18**   Daniel Krob. Monoides et semi-anneaux complets. In *Semigroup Forum*, volume 36, pages 323–339. Springer, 1987.

**19**   Werner Kuich. Semirings and formal power series: Their relevance to formal languages and automata. In *Handbook of Formal Languages*, volume 1, chapter 9, pages 609–677. Springer, 1997.

**20**   Werner Kuich and Arto Salomaa. *Semirings, Automata and Languages*. Springer-Verlag, Berlin, Heidelberg, 1985.

**21** Adelson-Velsky M, Alexander Gelbukh, and Eugene Levner. On fast path-finding algorithms in and-or graphs. *Mathematical Problems in Engineering*, 8, 2002. `doi:10.1080/10241230306728`.

**22** Michel Minoux and Michel Gondran. *Graphs, Dioids and Semirings. New Models and Algorithms*, volume 41 of *Operations Research/Computer Science Interfaces Series*. Springer, 2008. `doi:10.1007/978-0-387-75450-5`.

**23** Mehryar Mohri. Semiring frameworks and algorithms for shortest-distance problems. *J. Autom. Lang. Comb.*, 7(3):321–350, January 2002.

**24** Yann Ramusat, Silviu Maniu, and Pierre Senellart. Semiring provenance over graph databases. In *10th USENIX Workshop on the Theory and Practice of Provenance (TaPP 2018)*, London, United Kingdom, July 2018. URL: `https://hal.inria.fr/hal-01850510`.

**25** Yann Ramusat, Silviu Maniu, and Pierre Senellart. Provenance-based algorithms for rich queries over graph databases. In *EDBT 2021 - 24th International Conference on Extending Database Technology*, Nicosia / Virtual, Cyprus, March 2021.

**26** Günter Rote. Path problems in graphs. *Computing Supplementum*, 7, 1999. `doi:10.1007/978-3-7091-9076-0_9`.

**27** Bernhard Scholz, Herbert Jordan, Pavle Subotić, and Till Westmann. On fast large-scale program analysis in datalog. In *Proceedings of the 25th International Conference on Compiler Construction*, CC 2016, page 196–206, New York, NY, USA, 2016. Association for Computing Machinery. `doi:10.1145/2892208.2892226`.

**28** Noah A. Smith. Linguistic Structure Prediction. *Synthesis Lectures on Human Language Technologies*, 4(2):1–274, May 2011. `doi:10.2200/S00361ED1V01Y201105HLT013`.

**29** David Zhao, Pavle Subotić, and Bernhard Scholz. Debugging large-scale datalog: A scalable provenance evaluation strategy. *ACM Trans. Program. Lang. Syst.*, 42(2), April 2020. `doi:10.1145/3379446`.

## A    Proofs for Section 2 (Background)

▶ **Lemma 2.** *Let $S$ be an idempotent semiring and $\leqslant$ the natural order over $S$. Then $S$ is* superior *with respect to $\leqslant$ if and only if $S$ is 0-closed.*

**Proof.** First assume $S$ superior with respect to $\leqslant$. Then for any $a$, $\bar{1} \leqslant \bar{1} \otimes a = a$, which means that $\bar{1} + a = \bar{1}$, i.e., $S$ is 0-closed.

Now assume $S$ 0-closed. Since $a \oplus a \otimes b = a \otimes (\bar{1} \oplus b) = a$, we have: $a \leqslant a \otimes b$, and similarly for $b \leqslant a \otimes b$. Thus $S$ is superior with respect to $\leqslant$.    ◀

## B    Proofs for Section 3 (Datalog provenance and dynamic programming over hypergraphs)

▶ **Lemma 11.** *For any Datalog query $q$ and grounding of an atom $v$ of $q$, there is a bijection between $\mathcal{D}_{H_q}(v)$ and $\{\tau \mid \tau$ yields $v\}$.*

**Proof.** By definition of $H_q$ each instantiation of a rule corresponds to a unique hyperedge. Then, we can inductively construct for a given derivation $D$ its associated (unique) Datalog proof tree $\tau_D$:

- If $|D| = 1$, then $v$ is a source vertex and thus an extensional tuple, we get the empty proof.
- If $|D| \geqslant 1$, then there exists $e \in \mathrm{BS}(v)$ where $|e| \geqslant 0$ and $D_i$ a derivation of $T_i(e)$ for $1 \leqslant i \leqslant |e|$, where $D = \langle e, D_1 \cdots D_{|e|} \rangle$. By definition, this hyperedge corresponds to the grounding of a rule $t(\vec{x}) \leftarrow r_1(\vec{x_1}), \ldots, r_n(\vec{x_n})$. By induction, for $1 \leqslant i \leqslant |e|$, $\tau_{D_i}$ is the corresponding proof of the derivation $D_i$. Then by composition we obtain $\tau_D$ the proof for $D$.

◀

▶ **Lemma 12.** *For any Datalog query $q$ and grounding of an atom $v$ of $q$, for any derivation $D$ of $v$ in $H_q$ $w(D) = \bigotimes\limits_{t' \in leaves(\tau_D)} \mathsf{prov}_R^q(t')$ where $\tau_D$ is the proof tree corresponding to $D$ in the bijection given by Lemma 11.*

**Proof.** By induction on the size of the derivation $D$:

- If $|D| = 1$ then, there exists a nullary edge $e \in E_q$ with $h(e) = v$ and $w(D) = f_v = \mathsf{prov}_R^q(r(\vec{x})) = \prod\limits_{t' \in leaves(\tau_D)} \mathsf{prov}_R^q(t')$.
- If $|D| \geqslant 1$ then there exists $e \in E_q$ and $D$ is of the form $\langle e, D_1 \cdots D_{|e|} \rangle$ with $D_i$ a derivation of $T_i(e)$ for $1 \leqslant i \leqslant |e|$. We have $w(D) = f_e(w(D_1), \ldots, w(D_{|e|}))$ and by definition of $f_e = \otimes$ and by IHP $w(D) = \bigotimes\limits_{t' \in leaves(\tau_D)} \mathsf{prov}_R^q(t')$.

◀

▶ **Theorem 13.** *Let $t$ be a tuple of a Datalog program $q$ with output predicate $G$ and $H_q$ its hypergraph representation, then $\mathsf{prov}_G^q(t) = \delta_{H_q}(G(t))$.*

**Proof.**

$$
\begin{aligned}
\delta_{H_q}(t) \;\; &= \bigoplus_{D \in \mathcal{D}_{H_q}(G(t))} w(D) && \text{and by Lemma 12,} \\[2mm]
&= \bigoplus_{D \in \mathcal{D}_{H_q}(G(t))} \left( \bigotimes_{t' \in leaves(\tau_D)} \mathsf{prov}_R^q(t') \right) && \text{and by Lemma 11,} \\[2mm]
&= \bigoplus_{\tau \ yields \ t} \left( \bigotimes_{t' \in leaves(\tau)} \mathsf{prov}_R^q(t') \right) &&= \;\; \mathsf{prov}_T^q(t)
\end{aligned}
$$

◀

## C  Proofs for Section 4 (Best-first method)

We present as Algorithm 7 the Best-first method, the reformulation of the Knuth's algorithm in terms of semiring-based Datalog provenance, the basis of the optimizations introduced in the main text.

■ **Algorithm 7** Naïve version of Best-first method for Datalog provenance

**Input:** Datalog query $q$, EDB $D$ with provenance indications over a 0-closed totally-ordered semiring $S$.
**Output:** full Datalog provenance for the IDB of $q$ over $D$.
 1: $I \leftarrow \emptyset$
 2: Let $\nu$ be the function that maps all facts of $D$ to their annotation in $S$ and all potential facts of the intensional schema of $q$ to $\bar{0}$
 3: **repeat**
 4:     **for** each intensional fact $r(\vec{x}) \notin I$ **do**
 5:         $\nu(r(\vec{x})) \quad \oplus= \bigotimes_{1 \leqslant i \leqslant n} \nu(r_i(\vec{x_i}))$ for each instantiation of a rule
             $r(\vec{x}) \leftarrow r_1(\vec{x_1}), \dots, r_n(\vec{x_n})$ with $r_i(\vec{x_i}) \in D \cup I$
 6:         Add to $I$ the tuple $r(\vec{x_{\min}})$ such that $\nu(r(\vec{x_{\min}}))$ is $\leqslant$-minimal among all potential intensional facts $r(\vec{x})$ not in $I$
 7: **until** $\nu(r(\vec{x_{\min}})) = \bar{0}$ or $I$ contains all potential intensional facts
 8: **return** $\nu_{|I}$

## D  Proofs for Section 6 (Reverse translations and opportunities)

Let us start with a lemma showing the one-to-one correspondence between derivations in the hypergraph and proofs in Datalog.

▶ **Lemma 19.** *Given a vertex $v \in V$ of the hypergraph, $\mathcal{D}_H(v) \simeq \{\tau \mid \tau \text{ yields } v\}$.*

**Proof.** By definition of $q_H$ each rule corresponds to a unique hyperedge. Then, we can inductively construct for a given derivation $D$ its associated (unique) Datalog proof $\tau_D$:
- If $|D| = 1$, then $v$ is a source vertex (head of a nullary hyperedge). Thus, the corresponding proof is the single instantiation of the rule $R(v) \leftarrow E_1(v)$.
- If $|D| \geqslant 1$, then it exists $e \in BS(v)$ where $|e| \geqslant 0$ and $D_i$ is a derivation of $T_i(e)$ for $1 \leqslant i \leqslant |e|$, then $D = \langle e, D_1 \cdots D_{|e|} \rangle$ is a derivation of $v$. By definition, this edge corresponds to a single rule $R(v) \leftarrow E_{i+1}(v, T_1(e), \dots, T_{|e|}(e)), R(T_1(e)), \dots, R(T_{|e|}(e))$. By IHP, for $1 \leqslant i \leqslant |e|$, $\tau_{D_i}$ is the corresponding proof of the derivation $D_i$. Then by composition we obtain $\tau_D$ the proof for $D$.

◀

We then show the weight of each derivation of a tuple is equal to the corresponding proof weight in Datalog.

▶ **Lemma 20.** $w(D) = \bigotimes_{t' \in leaves(\tau_D)} \mathsf{prov}^{q_H}_{E_1, \dots, E_{n+1}}(t')$

**Proof.** By induction on the size of the derivation $D$:
- If $|D| = 1$ then, there exists a nullary edge $e \in E$ with $h(e) = v$ and $w(D) = f_e = \mathsf{prov}^{q_H}_{E_1}(e) = \bigotimes_{t' \in leaves(\tau_D)} \mathsf{prov}^{q_H}_{E_1, \dots, E_{n+1}}(t')$.

- If $|D| \geqslant 1$ then it exists $e \in E$ and $D$ is of the form $\langle e, D_1 \cdots D_{|e|} \rangle$ with $D_i$ a derivation of $T_i(e)$ for $1 \leqslant i \leqslant |e|$. We have $w(D) = f_e(w(D_1), \ldots, w(D_{|e|}))$ and by definition of $f_e = \otimes$ and by IHP $w(D) = \bigotimes\limits_{t' \in leaves(\tau_D)} \mathsf{prov}^{q_H}_{E_1, \ldots, E_{n+1}}(t')$.

◀

▶ **Theorem 16.** *Let $v$ be a vertex of a weighted hypergraph $H$ with all weight functions derived from the $\otimes$ operation of a semiring $S$ and $q_H$ its translation into a Datalog query, then $\delta_H(v) = \mathsf{prov}^{q_H}_R(v)$.*

**Proof.**

$$
\begin{aligned}
\delta_H(t) &= \bigoplus_{D \in \mathcal{D}_H(v)} w(D) && \text{and by Lemma 20,} \\[2mm]
&= \bigoplus_{D \in \mathcal{D}_H(v)} \left( \bigotimes_{t' \in leaves(\tau_D)} \mathsf{prov}^{q_H}_{E_1, \ldots, E_{n+1}}(t') \right) && \text{and by Lemma 19,} \\[2mm]
&= \bigoplus_{\tau \ yields \ t} \left( \bigotimes_{t' \in leaves(\tau)} \mathsf{prov}^{q_H}_{E_1, \ldots, E_{n+1}}(t') \right) &&= \ \mathsf{prov}^{q_H}_R(t)
\end{aligned}
$$

◀

▶ **Lemma 21.** *Given a vertex $v \in V$ of the hypergraph, $\mathcal{D}_{H_f}(v) \simeq \{\tau \mid \tau \ yields \ v\}$.*

**Proof.** We inductively construct for a given derivation $D$ its associated (unique) Datalog proof $\tau_D$:
- If $|D| = 1$, then $v$ is a source vertex (head of a nullary hyperedge $e$). Thus, the corresponding proof is composed of $R(v) \leftarrow E(v, e), H(e)$ and $H(e) \leftarrow Nullary(e)$.
- If $|D| \geqslant 1$, then there exists $e \in \mathrm{BS}(v)$ where $|e| \geqslant 0$ and $D_i$ is a derivation of $T_i(e)$ for $1 \leqslant i \leqslant |e|$, then $D = \langle e, D_1 \cdots D_{|e|} \rangle$ is a derivation of $v$. Let $v_1, \ldots, v_n$ be the elements of $\mathrm{T}(e))$, this hyperedge corresponds to the following proof tree:

$$
\begin{aligned}
R(v) &\quad\leftarrow\quad E(v, e), H(e) \\
H(e) &\quad\leftarrow\quad First(e, v_1), R(v_1), N(e, v_1) \\
&\quad For \ 1 \leqslant i \leqslant n - 1 : \\
N(e, v_i) &\quad\leftarrow\quad Next(e, v_i, v_{i+1}), R(v_{i+1}), N(e, v_{i+1}) \\
N(e, v_n) &\quad\leftarrow\quad End(e, v_n)
\end{aligned}
$$

By IHP, for $1 \leqslant i \leqslant |e|$, $\tau_{D_i}$ is the corresponding proof of the derivation $D_i$ (derivation of the fact $R(v_i)$). Then, by composition, we obtain $\tau_D$ the proof for $D$.

◀

▶ **Lemma 22.** $w(D) = \bigotimes\limits_{t' \in leaves(\tau_D)} \mathsf{prov}^{q_{H_f}}_{E, Nullary, First, Next, End}(t')$

**Proof.** By induction on the size of the derivation $D$:
- If $|D| = 1$ then there exists a nullary edge $e \in E$ with $h(e) = v$ and its associated (unique) derivation is $R(v) \leftarrow E(v, e), H(e)$ and $H(e) \leftarrow Nullary(e)$. Thus $w(D) = f_e = \mathsf{prov}^{q_{H_f}}_E(v, e) \otimes \mathsf{prov}^{q_{H_f}}_{Nullary}(e)$.
- If $|D| \geqslant 1$ then there exists $e \in E$ and $D$ is of the form $\langle e, D_1 \cdots D_{|e|} \rangle$ with $D_i$ a derivation of $T_i(e)$ for $1 \leqslant i \leqslant |e|$. The derivation tree given in Lemma 21 corresponding to the hyperedge has provenance:

$$
\mathsf{prov}^{q_{H_f}}_E(v, e) \ \otimes \ \mathsf{prov}^{q_{H_f}}_{First}(e) \ \otimes \ \mathsf{prov}^{q_{H_f}}_R(v_1) \otimes
$$
$$
\bigotimes_{1 \leqslant i \leqslant n-1} \left( \mathsf{prov}^{q_{H_f}}_{Next}(e, v_i, v_{i+1}) \otimes \mathsf{prov}^{q_{H_f}}_R(v_{i+1}) \right) \otimes \mathsf{prov}^{q_{H_f}}_{End}(e, v_n)
$$

By IHP, for $1 \leqslant i \leqslant |e|$, $w(D_i) = \mathsf{prov}_R^{q_{H_f}}(v_i)$, thus, we obtain

$$w(D) = \bigotimes_{t' \in \, leaves(\tau_D)} \mathsf{prov}_{E, Nullary, First, Next, End}^{q_{H_f}}(t').$$

◀

▶ **Theorem 18.** *Let $v$ be a vertex of a weighted hypergraph $H$ with all weight functions derived from the $\otimes$ operation of a semiring $S$ and $q_{H_f}$ its translation into a Datalog query with fixed schema, then $\delta_{H_f}(v) = \mathsf{prov}_R^{q_{H_f}}(v)$.*

**Proof.** We note EDB $= \{E, Nullary, First, Next, End\}$.

$$
\begin{aligned}
\delta_{H_f}(t) \quad &= \quad \bigoplus_{D \in \mathcal{D}_{H_f}(v)} w(D) &&\text{and by Lemma 22,} \\[2mm]
&= \quad \bigoplus_{D \in \mathcal{D}_{H_f}(v)} \left( \bigotimes_{t' \in \, leaves(\tau_D)} \mathsf{prov}_{\mathrm{EDB}}^{q_{H_f}}(t') \right) &&\text{and by Lemma 21,} \\[2mm]
&= \quad \bigoplus_{\tau \; yields \; t} \left( \bigotimes_{t' \in \, leaves(\tau)} \mathsf{prov}_{\mathrm{EDB}}^{q_{H_f}}(t') \right) \quad = \quad \mathsf{prov}_R^{q_{H_f}}(t)
\end{aligned}
$$

◀