

Modular Materialisation of Datalog Programs

Pan Hu and Boris Motik and Ian Horrocks

Department of Computer Science, University of Oxford
Oxford, United Kingdom
firstname.lastname@cs.ox.ac.uk

Abstract

The seminaïve algorithm can be used to materialise all consequences of a datalog program, and it also forms the basis for algorithms that incrementally update a materialisation as the input facts change. Certain (combinations of) rules, however, can be handled much more efficiently using custom algorithms. To integrate such algorithms into a general reasoning approach that can handle arbitrary rules, we propose a modular framework for computing and maintaining a materialisation. We split a datalog program into modules that can be handled using specialised algorithms, and we handle the remaining rules using the seminaïve algorithm. We also present two algorithms for computing the transitive and the symmetric-transitive closure of a relation that can be used within our framework. Finally, we show empirically that our framework can handle arbitrary datalog programs while outperforming existing approaches, often by orders of magnitude.

1 Introduction

Datalog (Abiteboul, Hull, and Vianu 1995) is a prominent rule language whose popularity is mainly due to its ability to express recursive definitions such as transitive closure. Datalog captures OWL 2 RL (Motik et al. 2009) ontologies with SWRL rules (Horrocks et al. 2004), so it supports query answering on the Semantic Web. It has been implemented in many systems, including but not limited to WebPIE (Urbani et al. 2012), VLog (Urbani, Jacobs, and Krötzsch 2016), Oracle’s RDF Store (Wu et al. 2008), OWLIM (Bishop et al. 2011a), and RDFox (Nenov et al. 2015).

Datalog reasoning is often realised by precomputing and storing all consequences of a datalog program and a set of facts; this process and its output are both called *materialisation*. A materialisation must be updated when the input facts change, but doing so ‘from scratch’ can be inefficient if changes are small. To minimise the overall work, *incremental maintenance algorithms* have been developed. These include the well-known *Delete/Rederive* (DRed) (Gupta, Mumick, and Subrahmanian 1993; Staudt and Jarke 1996) and *Counting* (Gupta, Mumick, and Subrahmanian 1993) algorithms, and the more recent *Backward/Forward* (B/F) (Motik et al. 2015), DRed^c, and B/F^c (Hu, Motik, and Horrocks 2018) algorithms.

Materialisation and all aforementioned incremental algorithms compute rule consequences using *seminaïve evaluation* (Abiteboul, Hull, and Vianu 1995). The main benefit of this approach is that each applicable inference is performed exactly once. However, all consequences of certain rules or rule combinations can actually be computed without considering every applicable inference. For example, consider applying a program that axiomatises a relation R as symmetric and transitive to input facts that describe a connected graph consisting of n vertices. In Section 3 we show that computing all consequences using seminaïve evaluation involves $O(n^3)$ rule applications, whereas a custom algorithm can achieve the same goal using only $O(n^2)$ steps. Since incremental maintenance algorithms are based on the seminaïve algorithm, they can suffer from similar deficiencies.

Approaches that can maintain the closure of specific datalog programs have already been considered in the literature. For example, maintaining transitive closure of a graph has been studied extensively (Ibaraki and Katoh 1983; La Poutre and van Leeuwen 1987; King 1999; Demetrescu and Italiano 2000). Subercaze et al. (2016) presented an algorithm for the materialisation of the transitive and symmetric properties in RDFS-Plus. Dong, Su, and Topor (1995) showed that insertions into a transitively closed relation can be maintained by evaluating four nonrecursive first-order queries. However, these approaches can only handle datalog programs for which they have been specifically developed—that is, the programs are not allowed to contain any additional rules. The presence of other rules introduces additional complexity since updates computed by specialised algorithms must be propagated to the remaining rules and vice versa. Moreover, many of these approaches cannot handle deletion of input facts, which is a key problem in incremental reasoning. Thus, it is currently not clear whether and how customised algorithms can be used in general-purpose datalog systems that must handle arbitrary datalog rules and support incremental additions and deletions.

To address these issues, in this paper we present a modular framework for materialisation computation and incremental materialisation maintenance that can integrate specialised reasoning algorithms with the seminaïve evaluation. The framework partitions the rules of a datalog program into disjoint subsets called *modules*. For each module, four plugable functions are used to compute certain consequences

of the module’s rules; there are no restrictions on how these functions are implemented, as long as their outputs satisfy certain conditions. Moreover, if no specialised algorithm for a module is available, the four functions can be implemented using seminaïve evaluation. Thus, our framework can efficiently handle certain combinations of rules, but it can also handle arbitrary rules while avoiding repeated inferences.

We then examine a module that axiomatises the transitive closure, and a module that axiomatises the symmetric–transitive closure. These modules capture node reachability in directed and undirected graphs, respectively, both of which frequently occur in practice and are thus highly relevant. We present the functions necessary to integrate these modules into our framework and show that they satisfy the properties needed for correctness. We also discuss the kinds of input that are likely to benefit from modular reasoning.

We have implemented our algorithms and compared them on several real-life and synthetic datasets. Our experiments illustrate the potential benefits of the proposed solution: our approach often outperforms state-of-the-art algorithms, sometimes by orders of magnitude. Our system and test data are available online.¹ All proofs of our results are given in the appendix.

2 Preliminaries

We now introduce datalog with stratified negation. A *term* is a constant or a variable. An *atom* has the form $P(t_1, \dots, t_k)$, where P is a k -ary *predicate* with $k \geq 0$, and each t_i , $1 \leq i \leq k$, is a term. A *fact* is a variable-free atom, and a *dataset* is a finite set of facts. A rule r has the form

$$B_1 \wedge \dots \wedge B_m \wedge \text{not } B_{m+1} \wedge \dots \wedge \text{not } B_n \rightarrow H,$$

where $0 \leq m \leq n$, and B_i and H are atoms. For r a rule, $h(r) = H$ is the *head*, $b^+(r) = \{B_1, \dots, B_m\}$ is the set of *positive body atoms*, and $b^-(r) = \{B_{m+1}, \dots, B_n\}$ is the set of *negative body atoms*. Each rule r must be *safe*—that is, each variable occurring in r must occur in at least one positive body atom. A *program* is a finite set of rules.

A *stratification* λ of a program Π maps each predicate occurring in Π to a positive integer such that, for each rule $r \in \Pi$ with predicate P in its head, $\lambda(P) \geq \lambda(R)$ (resp. $\lambda(P) > \lambda(R)$) holds for each predicate R occurring in $b^+(r)$ (resp. $b^-(r)$). Such r is *recursive* w.r.t. λ if $\lambda(P) = \lambda(R)$ holds for some predicate R occurring in $b^+(r)$; otherwise, r is *nonrecursive* w.r.t. λ . Program Π is *stratifiable* if a stratification λ of Π exists. For s an integer, the *stratum* s of Π is the program Π^s containing each rule $r \in \Pi$ whose head predicate P satisfies $\lambda(P) = s$. Moreover, let Π_r^s and Π_{nr}^s be the recursive and the nonrecursive subsets, respectively, of Π^s . Finally, let $O^s = \{P(c_1, \dots, c_n) \mid \lambda(P) = s \text{ and } c_i \text{ are constants}\}$.

A *substitution* σ is a mapping of finitely many variables to constants. For α a term, an atom, a rule, or a set thereof, $\alpha\sigma$ is the result of replacing each occurrence of a variable x in α with $\sigma(x)$, provided that the latter is defined.

If r is a rule and σ is a substitution mapping all variables of r to constants, then rule $r\sigma$ is an *instance* of r . For I

Algorithm 1 $\text{MAT}(\Pi, \lambda, E)$

```

1:  $I := \emptyset$ 
2: for each stratum index  $s$  with  $1 \leq s \leq S$  do
3:    $\Delta := (E \cap O^s) \cup \Pi_{nr}^s[I]$ 
4:   while  $\Delta \neq \emptyset$  do
5:      $I := I \cup \Delta$ 
6:      $\Delta := \Pi_r^s[I; \Delta] \setminus I$ 

```

a dataset, we define the set $\Pi[I]$ of all facts obtained by applying a program Π to I as

$$\Pi[I] = \bigcup_{r \in \Pi} \{h(r\sigma) \mid b^+(r\sigma) \subseteq I \text{ and } b^-(r\sigma) \cap I = \emptyset\}.$$

Let E be a dataset (called *explicit facts*) and let λ be a stratification of Π with maximum stratum index S . Then, let $I_\infty^0 = E$; for each $s \geq 1$, let $I_\infty^s = I_\infty^{s-1}$, let

$$I_i^s = I_{i-1}^s \cup \Pi^s[I_{i-1}^s] \text{ for } i > 0, \text{ and let } I_\infty^s = \bigcup_{i \geq 0} I_i^s.$$

Set I_∞^S is the *materialisation* of Π w.r.t. E and λ . It is known that I_∞^S does not depend on λ , so we write it as $\text{mat}(\Pi, E)$.

3 Motivation

In this section we show how custom algorithms can handle certain rule combinations much more efficiently than seminaïve evaluation. We consider here only materialisation, but similar observations apply to incremental maintenance algorithms as most of them use variants of seminaïve evaluation.

3.1 Seminaïve Evaluation

The seminaïve algorithm (Abiteboul, Hull, and Vianu 1995) takes as input a set of explicit facts E , a program Π , and a stratification λ of Π , and it computes $\text{mat}(\Pi, E)$. To apply each rule instance at most once, in each round of rule application it identifies the ‘newly applicable’ rule instances (i.e., instances that depend on a fact derived in the previous round) as shown in Algorithm 1. For each stratum, the algorithm initialises Δ , the set of newly derived facts, by combining the explicit facts in the current stratum ($E \cap O^s$) with the facts derivable from previous strata via nonrecursive rules ($\Pi_{nr}^s[I]$). Then, in lines 4–6 it iteratively computes all consequences of Δ . To this end, in line 6 it uses operator $\Pi[I; \Delta]$, which extends $\Pi[I]$ to allow identifying ‘newly applicable’ rule instances. Specifically, given datasets I and $\Delta \subseteq I$, operator $\Pi[I; \Delta]$ returns a set containing $h(r\sigma)$ for each rule $r \in \Pi$ and substitution σ such that $b^+(r\sigma) \subseteq I$ and $b^-(r\sigma) \cap I = \emptyset$ hold (i.e., rule instance $r\sigma$ is applicable to I), but also $b^+(r\sigma) \cap \Delta \neq \emptyset$ holds (i.e., a positive body atom of $r\sigma$ occurs in the set of facts Δ derived in the previous round of rule application). It is not hard to see that the algorithm computes $I = \text{mat}(\Pi, E)$, and that it considers each rule instance $r\sigma$ at most once.

3.2 Problems with the Seminaïve Evaluation

Although seminaïve evaluation does not repeat derivations, it always considers each applicable rule instance. However,

¹<http://krr-nas.cs.ox.ac.uk/2018/modular/>

facts are often derived via multiple, distinct rule instances; this is particularly common with recursive rules, but it can also occur with nonrecursive rules only. We are unaware of a general technique that can prevent such derivations. We next present two programs for which materialisation can be computed without considering all applicable rule instances, thus showing how seminaïve evaluation can be suboptimal.

Example 1. Let Π be the program containing rule (1) and let $E = \{R(c_i, c_{i+1}) \mid 0 \leq i \leq n\}$.

$$R(x, y) \wedge R(y, z) \rightarrow R(x, z) \quad (1)$$

Clearly, $I = \text{mat}(\Pi, E) = \{R(c_i, c_j) \mid 0 \leq i < j \leq n\}$, so each rule instance of the form

$$R(c_i, c_j) \wedge R(c_j, c_k) \rightarrow R(c_i, c_k) \quad (2)$$

with $1 \leq i < j < k \leq n$ is applicable to I . Algorithm 1 considers all of these $O(n^3)$ rule instances.

We next present an outline of an approach that is still cubic in general, but on this specific input runs in $O(n^2)$ time. The key is to distinguish the set X of ‘external’ facts given to Π as input from the ‘internal’ facts derived by Π . We can transitively close R by iteratively considering pairs of facts $R(u, v) \in X$ and $R(v, w)$. That is, we require the first fact to be in X , but place no restriction on the second fact. (We could have equivalently required the second fact to be in X .) In our example, we have $X = E$, so the algorithm considers only rule instances of the form

$$R(c_i, c_{i+1}) \wedge R(c_{i+1}, c_k) \rightarrow R(c_i, c_k) \quad (3)$$

for $0 \leq i < k \leq n$, of which there are $O(n^2)$ many. Intuitively, this is analogous to replacing the predicate R in all explicit facts with X , and using a linear rule

$$X(x, y) \wedge R(y, z) \rightarrow R(x, z) \quad (4)$$

instead of rule (1). In our approach, however, other rules can derive R -facts so the set X is not fixed; thus, rule (1) cannot be simply replaced with (4). Our approach ‘simulates’ such linearisation, and it can be expected to perform well whenever the other rules derive fewer facts than rule (1).

Example 2. Let Π consist of rules (1) and (5), and let $E = \{R(c_i, c_{i+1}) \mid 1 \leq i < n\} \cup \{R(c_n, c_1)\}$.

$$R(x, y) \rightarrow R(y, x) \quad (5)$$

Now $I = \text{mat}(\Pi, E) = \{R(c_i, c_j) \mid 1 \leq i, j \leq n\}$, so each instance of the form (2) with $1 \leq i, j, k \leq n$ is applicable to I . Algorithm 1 considers all of these $O(n^3)$ rule instances.

However, we can view any relation R as an undirected graph with n vertices. To compute the symmetric–transitive closure of R , we first compute the connected components of R , and, for each connected component C , we enumerate all $u, v \in C$ and derive $R(u, v)$. The first step is linear in the size of R and the second step requires $O(n^2)$ time, so the algorithm runs in $O(n^2)$ time on any R .

4 Framework

In this section we present a general framework for materialisation and incremental reasoning that can avoid the deficiencies outlined in Section 3 for certain rule combinations. Our framework focuses on recursive rules only: non-recursive rules Π_{nr}^s are evaluated just once in each stratum,

which is usually efficient. In contrast, the recursive part Π_r^s of each stratum Π^s must be evaluated iteratively, which is a common source of inefficiency. Thus, our framework splits Π_r^s into $n(s)$ mutually disjoint, nonempty programs $\Pi_r^{s,i}$, $1 \leq i \leq n(s)$, called *modules*. (We let $n(s) = 0$ if $\Pi_r^s = \emptyset$.) Our notion of modules should not be confused with ontology modules: the latter are subsets of an ontology that are semantically independent from each other in a well-defined way, whereas our modules are just arbitrary program subsets. Each module is handled using ‘plugin’ functions that compute certain consequences of $\Pi_r^{s,i}$. These functions can be implemented as desired, as long as their results satisfy certain properties that guarantee correctness. We present our framework in two steps: in Section 4.1 we consider materialisation, and in Section 4.2 we focus on incremental reasoning. Then, in Sections 5 and 6 we discuss how to realise these ‘plugin’ functions for certain common modules.

Before proceeding, we generalise operator $\Pi[I : \Delta]$ as follows. Given datasets I^p, I^n, Δ^p , and Δ^n where $\Delta^p \subseteq I^p$ and $\Delta^n \cap I^n = \emptyset$, let

$$\begin{aligned} \Pi[I^p, I^n : \Delta^p, \Delta^n] = \bigcup_{r \in \Pi} \{ & \mathbf{h}(r\sigma) \mid \\ & \mathbf{b}^+(r\sigma) \subseteq I^p \text{ and } \mathbf{b}^-(r\sigma) \cap I^n = \emptyset, \text{ and} \\ & \mathbf{b}^+(r\sigma) \cap \Delta^p \neq \emptyset \text{ or } \mathbf{b}^-(r\sigma) \cap \Delta^n \neq \emptyset\}. \end{aligned}$$

When the condition in the last line is not required, we simply write $\Pi[I^p, I^n]$. Moreover, we omit I^n when $I^p = I^n$, and we omit Δ^n when $\Delta^n = \emptyset$. Intuitively, this operator computes the consequences of Π by evaluating the positive and the negative body atoms in I^p and I^n , respectively, while ensuring in each derivation that either a positive or a negative body atom is true in Δ^p or Δ^n , respectively. Our incremental algorithm uses this operator to identify the consequences of Π that are affected by the changes to the facts matching the positive and the negative body atoms of the rules in Π . For example, if the facts in Δ^p are added to (resp. removed from) the materialisation, then $\Pi[I^p : \Delta^p]$ contains the consequences of the rule instances that start (resp. cease) to be applicable because a positive body atom matches to a fact in Δ^p . The set Δ^n is used to analogously capture the consequences of the negative body atoms of the rules in Π .

4.1 Computing the Materialisation

Our modular materialisation algorithm uses a ‘plugin’ function $\text{Add}^{s,i}$ for each module $\Pi_r^{s,i}$. The function takes as arguments datasets I^p, I^n , and Δ such that $\Delta \subseteq I^p$, and it closes I^p with all consequences of $\Pi_r^{s,i}$ that depend on Δ . Each invocation of these functions must satisfy the following properties in order to guarantee correctness of our algorithm.

Definition 3. Function Add captures a datalog program Π on datasets I^p, I^n , and Δ with $\Delta \subseteq I^p$ if the result of $\text{Add}[I^p, I^n, \Delta]$ is the smallest dataset J that satisfies $\Pi[I^p \cup J, I^n : \Delta \cup J] \subseteq I^p \cup J$.

For brevity, in the rest of the paper we often say just ‘ Add captures Π ’ without specifying the datasets whenever the latter are clear from the context. In the absence of a customised algorithm, Add can always be realised using the seminaïve evaluation strategy as follows:

Algorithm 2 MAT-MOD(Π, λ, E)

```
7:  $I := \emptyset$ 
8: for each stratum index  $s$  with  $1 \leq s \leq S$  do
9:    $\Delta_1 := \dots := \Delta_{n(s)} := \emptyset$ 
10:   $\Delta := (E \cap \mathcal{O}^s) \cup \Pi_{nr}^s[I]$ 
11:  while  $\Delta \neq \emptyset$  do
12:     $I := I \cup \Delta$ 
13:    for each  $i$  with  $1 \leq i \leq n(s)$  do
14:       $\Delta_i := \text{Add}^{s,i}[I, I, \Delta \setminus \Delta_i]$ 
15:     $\Delta := \Delta_1 \cup \dots \cup \Delta_{n(s)}$ 
```

- let $\Delta_0 = \Delta$ and $J_0 = \emptyset$,
- for i starting with 0 onwards, if $\Delta_i = \emptyset$, stop and return J_i ; otherwise, let $\Delta_{i+1} = \Pi[I^p \cup J_i, I^n : \Delta_i] \setminus (I^p \cup J_i)$ and $J_{i+1} = J_i \cup \Delta_{i+1}$ and proceed to $i + 1$.

However, a custom implementation of Add will typically not examine all rule instances from the above computation in order to optimise reasoning with certain modules.

Algorithm 2 formalises our modular approach to datalog materialisation. It takes as input a program Π , a stratification λ of Π , and a set of explicit facts E , and it computes $\text{mat}(\Pi, E)$. The algorithm’s structure is similar to Algorithm 1. For each stratum of Π , both algorithms first apply the nonrecursive rules, and then they apply the recursive rules iteratively up to a fixpoint. The main difference is that, given a set of facts Δ derived from the previous iteration, Algorithm 2 computes the consequences of Δ for each module independently using $\text{Add}^{s,i}$ (line 14); note that each Δ_i is closed under $\Pi_r^{s,i}$, which is key to the performance of our approach. The algorithm then combines the consequences of all modules (line 15) before proceeding to the next iteration.

If each $\text{Add}^{s,i}$ function is implemented using seminaïve evaluation as described earlier, then the algorithm does not consider a rule instance more than once. This is achieved by passing $\Delta \setminus \Delta_i$ to $\text{Add}^{s,i}$ in line 14: only facts derived by other modules in the previous iteration are considered ‘new’ for $\text{Add}^{s,i}$, which is possible since the facts in Δ_i have been produced by the i th module in the previous iteration. Theorem 4 captures these properties formally.

Theorem 4. *Algorithm 2 computes I as $\text{mat}(\Pi, E)$ if function $\text{Add}^{s,i}$ captures $\Pi_r^{s,i}$ in each of its calls. Moreover, if all $\text{Add}^{s,i}$ use the seminaïve strategy, each applicable rule instance is considered at most once.*

4.2 Incremental Updates

Our modular incremental materialisation maintenance algorithm is based on the DRed^c algorithm by Hu, Motik, and Horrocks (2018), which is a variant of the well-known DRed algorithm (Gupta, Mumick, and Subrahmanian 1993). For each fact, DRed^c maintains two counters that track the number of nonrecursive and recursive derivations of the fact. The algorithm proceeds in three steps. During the deletion phase, DRed^c iteratively computes the consequences of the deleted facts, similar to DRed, while adjusting the counters accordingly. However, to optimise overdeletion, deletion propagation stops on facts with a nonzero nonrecursive counter. In

the one-step rederivation phase, DRed^c identifies the facts that were overdeleted but can be rederived from the remaining facts in one step by simply checking the recursive counters: if the counter is nonzero, then the corresponding fact is rederived. In the insertion phase, DRed^c computes the consequences of the rederived and the inserted facts using seminaïve evaluation, which we have already discussed.

Our modular incremental algorithm handles nonrecursive rules in the same way as DRed^c. Thus, the nonrecursive counters, which record the number of nonrecursive derivations of each fact, can be maintained globally just as in DRed^c. In contrast, as discussed in Section 3, custom algorithms for recursive modules will usually not consider all applicable rule instances, so counters of recursive derivations cannot be maintained globally. Nevertheless, certain modules can maintain recursive counters internally (e.g., the module based on the seminaïve evaluation can do so).

In addition to function $\text{Add}^{s,i}$ from Section 4.1, our modular incremental reasoning algorithm uses three further functions: $\text{Diff}^{s,i}$, $\text{Del}^{s,i}$, and $\text{Red}^{s,i}$. Definition 5 captures the requirements on $\text{Diff}^{s,i}$. Intuitively, $\text{Diff}^{s,i}[I^p, \Delta^p, \Delta^n]$ identifies the consequences of $\Pi_r^{s,i}$ affected by the addition of the facts in Δ^p and removal of the the facts Δ^n , respectively, with both sets containing facts from earlier strata.

Definition 5. *Function Diff captures a datalog program Π on datasets I^p , Δ^p , and Δ^n where $\Delta^p \subseteq I^p$, $\Delta^n \cap I^p = \emptyset$, and both Δ^p and Δ^n do not contain predicates occurring in rule heads in Π if $\text{Diff}[I^p, \Delta^p, \Delta^n] = \Pi[I^p : \Delta^p, \Delta^n]$.*

Function $\text{Del}^{s,i}$ captures overdeletion: if the facts in Δ are deleted, then $\text{Del}^{s,i}[I^p, I^n, \Delta, C_{nr}]$ returns the consequences of $\Pi_r^{s,i}$ that must be overdeleted as well. The function can use the nonrecursive counters C_{nr} in order to stop overdeletion as in DRed^c. We do not specify exactly what the functions must return: as we discuss in Section 6, computing the smallest set that needs to be overdeleted might require considering all rule instances as in DRed^c, which would miss the point of modular reasoning. Instead, we specify the required output in terms of a lower bound J_l and an upper bound J_u . Intuitively, J_l and J_u contain facts that would be overdeleted in DRed^c and DRed, respectively.

Definition 6. *Function Del captures a datalog program Π on datasets I^p , I^n , Δ with $\Delta \subseteq I^p$, and a mapping C_{nr} of facts to integers if $J_l \subseteq \text{Del}[I^p, I^n, \Delta, C_{nr}] \subseteq J_u$ where*

- the lower bound J_l is the smallest dataset such that, for each $F \in \Pi[I^p, I^n : \Delta \cup J_l]$, either $F \in \Delta \cup J_l$ or $C_{nr}(F) > 0$ holds, and
- the upper bound J_u is the smallest dataset that satisfies $\Pi[I^p, I^n : \Delta \cup J_u] \subseteq \Delta \cup J_u$.

Finally, function $\text{Red}^{s,i}$ captures rederivation: if facts in Δ are overdeleted, then $\text{Red}^{s,i}[I^p, I^n, \Delta]$ returns all facts in Δ that can be rederived from $I^p \setminus \Delta$ and $\Pi_r^{s,i}$ in one or more steps. This is different from DRed and DRed^c, which both perform only one-step rederivation. This change is important in our framework because, as we shall see in Section 6, $\text{Red}^{s,i}$ provides the opportunity for a module to adjust its internal data structures after deletion.

Definition 7. Function $\text{Red}^{s,i}$ captures a datalog program Π on datasets I^p, I^n, Δ with $\Delta \subseteq I^p$ if the result of $\text{Red}[I^p, I^n, \Delta]$ is the smallest dataset J that satisfies $\Pi[(I^p \setminus \Delta) \cup J, I^n] \cap \Delta \subseteq J$.

Algorithm 3 formalises our modular approach to incremental maintenance. The algorithm takes as input a program Π , a stratification λ of Π , a set of explicit facts E , the materialisation $I = \text{mat}(\Pi, E)$, the sets of facts E^- and E^+ to delete from and add to E , and a map C_{nr} that records the number of nonrecursive derivations of each fact. The algorithm updates I to $\text{mat}(\Pi, (E \setminus E^-) \cup E^+)$. We next describe the two main steps of the algorithm.

In the overdeletion phase, the algorithm first initialises the set of facts to delete Δ as the union of the explicitly deleted facts ($E^- \cap O^s$) and the facts affected by changes in previous strata (lines 22 and 23). Then, in lines 26–30 the algorithm computes all consequences of Δ . In each iteration, function $\text{Del}^{s,i}$ is called for each module to identify the consequences of $\Pi_r^{s,i}$ that must be overdeleted due to the deletion of Δ (line 28). As in Algorithm 2, the third argument of $\text{Del}^{s,i}$ is $\Delta \setminus \Delta_i$, which guarantees that the function will not be applied to its own consequences.

In the second step, the algorithm first identifies the rederivable facts by calling $\text{Red}^{s,i}$ for each module (lines 33–35). Then, the consequences of the rederived facts, the explicitly added facts ($E^+ \cap O^s$), and the facts added due to changes in previous strata are computed in the loop of lines 36–40 analogously to Algorithm 2. Although $\text{Red}^{s,i}$ rederives facts in one or more steps as opposed to the one-step rederivation in DRed and DRed^c , this extra effort is not repeated during insertion since $\text{Add}^{s,i}$ is not applied to the consequences of module i . Theorem 8 states that the algorithm is correct.

Theorem 8. Algorithm 3 updates I from $\text{mat}(\Pi, E)$ to $\text{mat}(\Pi, (E \setminus E^-) \cup E^+)$ if functions $\text{Add}^{s,i}, \text{Del}^{s,i}, \text{Diff}^{s,i}$, and $\text{Red}^{s,i}$ capture $\Pi_r^{s,i}$ in all of their calls.

5 Transitive Closure

We now consider a module consisting of a single rule (1) axiomatising a relation R as transitive. Following the ideas from Example 1, we distinguish the ‘internal’ facts produced by rule (1) from the ‘external’ facts produced by other rules. We keep track of the latter in a global set X_R that is initialised to the empty set. A key invariant of our approach is that each fact $R(a_0, a_n)$ is produced by a chain $\{R(a_0, a_1), \dots, R(a_{n-1}, a_n)\} \subseteq X_R$ of ‘external’ facts. Thus, we can transitively close R by considering pairs of R -facts where at least one of them is contained in X_R , which can greatly reduce the number of inferences. A similar effect could be achieved by rewriting the input program: we introduce a fresh predicate X_R , and we replace by X_R each occurrence of R in the head of a rule, as well as one of the two occurrences of R in the body of rule (1). Such an approach, however, introduces the facts containing the auxiliary predicate X_R into the materialisation and thus reveals implementation details to the users. Moreover, the rederivation step can be realised very efficiently in our approach.

Algorithm 3 $\text{DRED}^c\text{-MOD}(\Pi, \lambda, E, I, E^-, E^+, C_{nr})$

```

16:  $D := A := \emptyset, E^- = (E^- \cap E) \setminus E^+, E^+ = E^+ \setminus E$ 
17: for each stratum index  $s$  with  $1 \leq s \leq S$  do
18:   OVERDELETE
19:   REDERIVE-INSERT
20:  $E := (E \setminus E^-) \cup E^+, I := (I \setminus D) \cup A$ 

21: procedure OVERDELETE
22:    $\Delta_1 := \dots := \Delta_{n(s)} := \emptyset$ 
23:    $\Delta := (E^- \cap O^s) \cup \Pi_{nr}^s[I \setminus D \setminus A, A \setminus D]$  and update  $C_{nr}$ 
24:   for each  $i$  with  $1 \leq i \leq n(s)$  do
25:      $\Delta := \Delta \cup \text{Diff}^{s,i}[I, D \setminus A, A \setminus D]$ 
26:   while  $\Delta \neq \emptyset$  do
27:     for each  $i$  with  $1 \leq i \leq n(s)$  do
28:        $\Delta_i := \text{Del}^{s,i}[I \setminus (D \setminus A), I \cup A, \Delta \setminus \Delta_i, C_{nr}]$ 
29:        $D := D \cup \Delta$ 
30:        $\Delta := \Delta_1 \cup \dots \cup \Delta_{n(s)}$ 

31: procedure REDERIVE-INSERT
32:    $\Delta := (E^+ \cap O^s) \cup \Pi_{nr}^s[(I \setminus D) \cup A \setminus D, D \setminus A]$ 
   and update  $C_{nr}$ 
33:   for each  $i$  with  $1 \leq i \leq n(s)$  do
34:      $\Delta_i := \text{Red}^{s,i}[I, (I \setminus D) \cup A, D \setminus A]$ 
35:      $\Delta := \Delta \cup \Delta_i \cup \text{Diff}^{s,i}[(I \setminus D) \cup A, A \setminus D, D \setminus A]$ 
36:   while  $\Delta \neq \emptyset$  do
37:      $A := A \cup \Delta$ 
38:     for each  $i$  with  $1 \leq i \leq n(s)$  do
39:        $\Delta_i := \text{Add}^{s,i}[(I \setminus D) \cup A, (I \setminus D) \cup A, \Delta \setminus \Delta_i]$ 
40:      $\Delta := \Delta_1 \cup \dots \cup \Delta_{n(s)}$ 

```

Based on the above idea, function $\text{Add}^{\text{tc}(R)}$, shown in Algorithm 4, essentially implements seminaïve evaluation for rule (4): the loops in lines 42–43 and 44–47 handle the two delta rules derived from (4). For $\text{Diff}^{\text{tc}(R)}$, note that sets $A \setminus D$ and $D \setminus A$ in lines 25 and 35 of Algorithm 3 always contain facts with predicates that do not occur in $\Pi_r^{s,i}$ in rule heads; thus, since R occurs in the head of rule (1), these sets contain facts whose predicate is different from R , so $\text{Diff}^{\text{tc}(R)}$ can simply return the empty set. Function $\text{Del}^{\text{tc}(R)}$, shown in Algorithm 5, implements seminaïve evaluation for rule (4) analogously to $\text{Add}^{\text{tc}(R)}$. The main difference is that only facts whose nonrecursive counter is zero are overdeleted, which mimics overdeletion in DRed^c . As a result, not all facts processed in lines 50 and 53 are added to J so, to avoid repeatedly considering such facts, the algorithm maintains the set S of ‘seen’ facts. Finally, function $\text{Red}^{\text{tc}(R)}$, shown in Algorithm 6, identifies for each source vertex u all vertices reachable by the external facts in X_R .

Theorem 9. In each call in Algorithms 2 and 3, functions $\text{Add}^{\text{tc}(R)}, \text{Del}^{\text{tc}(R)}, \text{Diff}^{\text{tc}(R)}$, and $\text{Red}^{\text{tc}(R)}$ capture a datalog program that axiomatises relation R as transitive.

6 Symmetric–Transitive Closure

We now consider a module consisting of two rules, (1) and (5), axiomatising a relation R as transitive and symmetric. As in Example 2, we can view relation R as an undirected graph. To compute the materialisation, we extract the set C_R

Algorithm 4 $\text{Add}^{\text{tc}(R)}[I^p, I^n, \Delta]$

```
41:  $J := \emptyset, Q := \Delta, X_R := X_R \cup \Delta$ 
42: for each  $R(u, v) \in \Delta$  and each  $R(v, w) \in I^p \setminus \Delta$  do
43:   add  $R(u, w)$  to  $Q$  and  $J$ 
44: while  $Q \neq \emptyset$  do
45:   remove an arbitrarily chosen fact  $R(v, w)$  from  $Q$ 
46:   for each  $R(u, v) \in X_R$  such that  $R(u, w) \notin I^p \cup J$  do
47:     add  $R(u, w)$  to  $Q$  and  $J$ 
48: return  $J$ 
```

Algorithm 5 $\text{Del}^{\text{tc}(R)}[I^p, I^n, \Delta, C_{\text{nr}}]$

```
49:  $J := \emptyset, Q := S := \Delta, X_R := X_R \setminus \Delta$ 
50: for each  $R(u, v) \in \Delta$  and each  $R(v, w) \in I^p \setminus S$  do
51:   add  $R(u, w)$  to  $Q$  and  $S$ 
52:   if  $C_{\text{nr}}(R(u, w)) = 0$  then add  $R(u, w)$  to  $J$ 
53: while  $Q \neq \emptyset$  do
54:   remove an arbitrarily chosen fact  $R(v, w)$  from  $Q$ 
55:   for each  $R(u, v) \in X_R$  such that  $R(u, w) \in I^p \setminus S$  do
56:     add  $R(u, w)$  to  $Q$  and  $S$ 
57:     if  $C_{\text{nr}}(R(u, w)) = 0$  then add  $R(u, w)$  to  $J$ 
58: return  $J \setminus \Delta$ 
```

Algorithm 6 $\text{Red}^{\text{tc}(R)}[I^p, I^n, \Delta]$

```
59:  $J := \emptyset$ 
60: for each  $u$  such that there exist  $v$  with  $R(u, v) \in \Delta$  do
61:   for each  $w$  reachable from  $u$  via  $R$  facts in  $X_R$  do
62:     add  $R(u, w)$  to  $J$ 
63: return  $J \cap \Delta$ 
```

of connected components—that is, each $U \in C_R$ is a set of mutually connected vertices in the symmetric–transitive closure of R ; finally, we derive $R(u, v)$ for all u and v in each component $U \in C_R$. Set C_R is global and is initially empty.

Based on this idea, function $\text{Add}^{\text{stc}(R)}$, shown in Algorithm 7, uses an auxiliary function CLOSEEDGES to incrementally update the set C_R by processing each fact $R(u, v) \in \Delta$ in lines 67–75: if either u or v does not occur in a component in C_R , then the respective component is created in C_R (lines 69 and 71); and if u and v belong to distinct components U and V , then U and V are merged into a single component and all R -facts connecting U and V are added (lines 72–75). For the same reasons as in Section 5, function $\text{Diff}^{\text{tc}(R)}$ can simply return the empty set. Function $\text{Del}^{\text{stc}(R)}$, shown in Algorithm 8, simply overdeletes all facts $R(u', v')$ whose nonrecursive counter is zero and where both u' and v' belong to a component U containing both vertices of a fact $R(u, v)$ in Δ . Those facts $R(u', v')$ for which the nonrecursive counter is nonzero will hold after overdeletion, so they are kept in an initially empty global set Y_R so that they can be used for rederivation later. Finally, function $\text{Red}^{\text{stc}(R)}$, shown in Algorithm 9, simply closes the set Y_R in the same way as during addition, and it empties the set Y_R . While this creates a dependency between $\text{Del}^{\text{stc}(R)}$ and $\text{Red}^{\text{stc}(R)}$, the order in which these functions are called in Algorithm 3 ensures that the set Y_R is maintained correctly.

Algorithm 7 $\text{Add}^{\text{stc}(R)}[I^p, I^n, \Delta]$

```
64: return  $\text{CLOSEEDGES}(\Delta) \setminus I^p$ 
65: function  $\text{CLOSEEDGES}(\Delta)$ 
66:    $J := \emptyset$ 
67:   for each  $R(u, v) \in \Delta$  do
68:     if no  $U \in C_R$  exists such that  $u \in U$  then
69:       add  $\{u\}$  to  $C_R$ , and  $R(u, u)$  to  $J$ 
70:     if no  $V \in C_R$  exists such that  $v \in V$  then
71:       add  $\{v\}$  to  $C_R$ , and  $R(v, v)$  to  $J$ 
72:     if  $u$  and  $v$  belong to distinct  $U, V \in C_R$ , resp. then
73:       remove  $U$  and  $V$  from  $C_R$ , and add  $U \cup V$  to  $C_R$ 
74:       for each  $u' \in U$  and each  $v' \in V$  do
75:         add  $R(u', v')$  and  $R(v', u')$  to  $J$ 
76:   return  $J$ 
```

Algorithm 8 $\text{Del}^{\text{stc}(R)}[I^p, I^n, \Delta, C_{\text{nr}}]$

```
77:  $J := \emptyset$ 
78: for each  $U \in C_R$  where  $\exists R(u, v) \in \Delta$  s.t.  $\{u, v\} \subseteq U$  do
79:   for each  $u' \in U$  and each  $v' \in U$  do
80:     if  $C_{\text{nr}}(R(u', v')) = 0$  then add  $R(u', v')$  to  $J$ 
81:     else add  $R(u', v')$  to  $Y_R$ 
82:   remove  $U$  from  $C_R$ 
83: return  $J \setminus \Delta$ 
```

Algorithm 9 $\text{Red}^{\text{stc}(R)}[I^p, I^n, \Delta]$

```
84:  $J := \text{CLOSEEDGES}(Y_R) \cap \Delta$ 
85:  $Y_R := \emptyset$ 
86: return  $J$ 
```

Theorem 10. *In each call in Algorithms 2 and 3, functions $\text{Add}^{\text{stc}(R)}$, $\text{Del}^{\text{stc}(R)}$, $\text{Diff}^{\text{stc}(R)}$, and $\text{Red}^{\text{stc}(R)}$ capture a datalog program that axiomatises R as symmetric–transitive.*

7 Evaluation

We have implemented our modular materialisation and incremental maintenance algorithms, as well as the seminaïve materialisation and the DRed^c algorithms, and we have compared their performance empirically.

Test Benchmarks We used the following real-world and synthetic benchmarks in our tests. LUBM (Guo, Pan, and Heflin 2005) is a well-known benchmark that models individuals and organisations in a university domain. Claros describes archeological artefacts. We used the LUBM and Claros datasets with the *lower bound extended* (-LE) programs by Motik et al. (2014); roughly speaking, we converted a subset of the accompanying OWL ontologies into datalog and manually extended them with several ‘difficult rules’. DBpedia (Lehmann et al. 2015) contains structured information extracted from Wikipedia. DBpedia represents Wikipedia categories using the SKOS vocabulary (Miles and Bechhofer 2009), which defines several transitive properties. We used the datalog subset of the SKOS RDF schema. Moreover, the materialisation of DBpedia-SKOS is too large to fit into the memory of our test server, so we used a random sample of the DBpedia dataset consisting of five mil-

Benchmark	$ E $	$ I $	S	$ \Pi_{nr} $	$ \Pi_r $	$ TC $	$ STC $	Mat-Mod	Mat
Claros-LE	18.8 M	533.3 M	11	1031	306	27	2	733.55	3593.32
LUBM-LE	133.6 M	332.6 M	5	85	22	1	2	291.90	1100.22
DBpedia-SKOS	5.0 M	97.0 M	5	26	15	2	1	103.23	3623.37
DAG-R	0.1 M	22.9 M	1	1	1	1	0	29.60	3238.86

Table 1: Running times for materialisation computation (seconds)

Benchmark	Small Deletions		Small Insertions		Large Deletions	
	DRed ^c -Mod	DRed ^c	DRed ^c -Mod	DRed ^c	DRed ^c -Mod	DRed ^c
Claros-LE	0.93	1035.28	0.17	0.80	314.33	3616.93
LUBM-LE	0.32	3.87	0.01	0.01	182.93	1369.77
DBpedia-SKOS	21.77	691.32	0.20	2.78	111.28	3826.87
DAG-R	64.92	3005.11	14.56	116.78	62.48	4316.71

Table 2: Running times for incremental maintenance (seconds)

lion facts. Finally, DAG-R is a synthetic benchmark consisting of a randomly generated dataset containing a directed acyclic graph with 10k nodes and 100k edges, and a program that axiomatises the path relation as transitive. Table 1 shows the numbers of explicit facts ($|E|$), derived facts ($|I|$), strata (S), nonrecursive rules ($|\Pi_{nr}|$), recursive rules ($|\Pi_r|$), transitivity modules ($|TC|$), and symmetric-transitive modules ($|STC|$) for each benchmark.

Test Setup and Results We conducted all experiments on a Dell PowerEdge R730 server with 512 GB RAM and two Intel Xeon E5-2640 2.6 GHz processors running Fedora 27, kernel version 4.17.6. For each benchmark, we loaded the test data into our system and then compared the performance of our modular algorithms with the seminaïve and DRed^c algorithms using the following methodology.

We first computed the materialisation and measured the wall-clock time. The results are shown in Table 1. We then conducted two groups of incremental reasoning tests.

In the first group, we tested the performance of our incremental algorithms on small changes. To this end, we used uniform sampling to select ten subsets $E_i \subseteq E$, $1 \leq i \leq 10$, each consisting of 1000 facts from the input dataset. We deleted and then reinserted E_i for each i while measuring the wall-clock times, and then we computed the average times for deletion and insertion over the ten samples. The results are shown in the ‘Small Deletions’ and ‘Small Insertions’ columns of Table 2, respectively.

In the second group, we tested the performance of incremental algorithms on large deletions. To this end, we used uniform sampling to select a subset $E^- \subseteq E$ containing 25% of the explicit facts, and we measured the wall-clock time needed to delete E^- from the materialisation. The results are shown in the ‘Large Deletions’ column of Table 2. We did not consider large insertions because our algorithms handle insertion in the same way as materialisation, so the relative performance of our algorithms should be similar to the performance of materialisation shown in Table 1.

Discussion Mat-Mod significantly outperformed Mat on all test inputs. For example, Mat-Mod was several times faster than Mat on Claros-LE and LUBM-LE. The programs

of both benchmarks contain transitivity and symmetric-transitivity modules, which are efficiently handled by our custom algorithm. The performance improvement is even more significant for DBpedia-SKOS and DAG-R: Mat-Mod is more than 30 times faster than Mat on DBpedia-SKOS, and the difference reaches two orders of magnitude on DAG-R. In fact, DBpedia contains long chains/cycles over the *skos:broader* relation (Bishop et al. 2011b), which is axiomatised as transitive in SKOS. Mat-Mod outperforms Mat in this case since our custom algorithm for transitivity skips a large number of rule instances. The same observation explains the superior performance of DAG-R.

Similarly, DRed^c-Mod considerably outperformed DRed^c on small deletions: the performance speedup ranges from around ten times on LUBM-LE to three orders of magnitude on Claros-LE. The program of Claros-LE contains a symmetric-transitive closure module for the predicate *relatedPlaces*, and the materialisation contains large cliques of constants connected to each other via this predicate. Thus, when a *relatedPlaces(a,b)* fact is deleted, DRed^c can end up considering up to n^3 rule instances where n is the number of constants in the clique containing a and b . In contrast, our custom algorithm for this module maintains a connected component for the clique and requires only up to n^2 steps. It is worth noticing that, while DRed^c-Mod significantly outperforms DRed^c on DAG-R, the incremental update times for small deletion were larger than both the update times for large deletions and even for the initial materialisation. This is because deleting one thousand edges from the graph (‘Small Deletion’) caused a large part of the materialisation to be overdeleted and rederived again. In contrast, when 25% of the explicit facts are deleted (‘Large Deletion’), a larger proportion of the materialisation is overdeleted, but only a few facts are rederived. For DRed^c the situation is similar, but rederivation in DRed^c benefits from a global recursive counter (at the expense of considering each applicable rule instance), which makes small deletion still faster than large deletion and initial materialisation. Finally, as shown in Table 2, DRed^c-Mod scaled well and maintained its advantage over DRed^c on large deletions.

Incremental insertions are in general easier to handle than

deletions since during insertion the algorithms can rely on the whole materialisation to prune the propagation of facts whereas during deletion the algorithms can only rely on the nonrecursive counters of facts to do the same. This is clearly reflected in Table 2. Nevertheless, in our tests for small insertions, DRed^c-Mod was several times faster than DRed^c in all cases but LUBM-LE, for which both algorithms updated the materialisation instantaneously.

8 Conclusion

We have proposed a modular framework for the computation and maintenance of datalog materialisations. The framework allows integrating custom algorithms for specific types of rules with standard datalog reasoning methods. Moreover, we have presented such custom algorithms for programs axiomatising the transitive and the symmetric-transitive closure of a relation. Finally, we have shown empirically that our algorithms typically significantly outperform then existing ones, sometimes by orders of magnitude. In future, we plan to extend our framework also to the B/F^c algorithm, which eliminates overdeletion by eagerly checking alternative derivations. This could potentially be useful in cases such as DBpedia-SKOS and DAG-R, where overdeletion is a major source of inefficiency.

Acknowledgements

We thank David Tena Cucala for his help with obtaining the SKOS benchmark. This work was supported by the EPSRC projects AnaLOG, DBOnto, and ED³.

References

- Abiteboul, S.; Hull, R.; and Vianu, V. 1995. *Foundations of Databases*. Addison-Wesley.
- Bishop, B.; Kiryakov, A.; Ognyanoff, D.; Peikov, I.; Tashev, Z.; and Velkov, R. 2011a. OWLIM: A family of scalable semantic repositories. *SWJ* 2(1):33–42.
- Bishop, B.; Kiryakov, A.; Ognyanov, D.; Peikov, I.; Tashev, Z.; and Velkov, R. 2011b. Factforge: A fast track to the web of data. *Semantic Web* 2(2):157–166.
- Demetrescu, C., and Italiano, G. F. 2000. Fully dynamic transitive closure: breaking through the $\mathcal{O}(n/\sup 2/)$ barrier. In *Foundations of Computer Science, 2000. Proceedings. 41st Annual Symposium on*, 381–389. IEEE.
- Dong, G.; Su, J.; and Topor, R. W. 1995. Nonrecursive Incremental Evaluation of Datalog Queries. *Annals of Mathematics and Artificial Intelligence* 14(2–4):187–223.
- Guo, Y.; Pan, Z.; and Heflin, J. 2005. Lubm: A benchmark for owl knowledge base systems. *Web Semantics: Science, Services and Agents on the World Wide Web* 3(2-3):158–182.
- Gupta, A.; Mumick, I. S.; and Subrahmanian, V. S. 1993. Maintaining Views Incrementally. In *SIGMOD*. ACM.
- Horrocks, I.; Patel-Schneider, P. F.; Boley, H.; Tabet, S.; Grosz, B.; Dean, M.; et al. 2004. SWRL: A Semantic Web Rule Language Combining OWL and RuleML. W3C Member Submission.
- Hu, P.; Motik, B.; and Horrocks, I. 2018. Optimised maintenance of datalog materialisations. In *AAAI*, 1871–1879.
- Ibaraki, T., and Katoh, N. 1983. On-line computation of transitive closures of graphs. *Information Processing Letters* 16(2):95–97.
- King, V. 1999. Fully dynamic algorithms for maintaining all-pairs shortest paths and transitive closure in digraphs. In *Foundations of Computer Science, 1999. 40th Annual Symposium on*, 81–89. IEEE.
- La Poutre, J. A., and van Leeuwen, J. 1987. Maintenance of transitive closures and transitive reductions of graphs. In *International Workshop on Graph-Theoretic Concepts in Computer Science*, 106–120. Springer.
- Lehmann, J.; Isele, R.; Jakob, M.; Jentzsch, A.; Kontokostas, D.; Mendes, P. N.; Hellmann, S.; Morsey, M.; Van Kleef, P.; Auer, S.; et al. 2015. Dbpedia—a large-scale, multilingual knowledge base extracted from wikipedia. *Semantic Web* 6(2):167–195.
- Miles, A., and Bechhofer, S. 2009. Skos simple knowledge organization system reference. *W3C recommendation* 18:W3C.
- Motik, B.; Patel-Schneider, P.; Parsia, B.; Bock, C.; Fokoue, A.; Haase, P.; Hoekstra, R.; Horrocks, I.; Ruttenberg, A.; Sattler, U.; et al. 2009. OWL 2 Web Ontology Language: Structural Specification and Functional-Style Syntax. W3C.
- Motik, B.; Nenov, Y.; Piro, R.; Horrocks, I.; and Olteanu, D. 2014. Parallel materialisation of datalog programs in centralised, main-memory rdf systems. In *AAAI*, 129–137.
- Motik, B.; Nenov, Y.; Piro, R.; and Horrocks, I. 2015. Incremental Update of Datalog Materialisation: the Backward/Forward Algorithm. In *AAAI*, 1560–1568.
- Nenov, Y.; Piro, R.; Motik, B.; Horrocks, I.; Wu, Z.; and Banerjee, J. 2015. RDFox: A Highly-Scalable RDF Store. In *ISWC*, 3–20.
- Staudt, M., and Jarke, M. 1996. Incremental Maintenance of Externally Materialized Views. In *VLDB*, 75–86.
- Subercaze, J.; Gravier, C.; Chevalier, J.; and Laforest, F. 2016. Inferray: fast in-memory RDF inference. *PVLDB* 9(6):468–479.
- Urbani, J.; Kotoulas, S.; Maassen, J.; Van Harmelen, F.; and Bal, H. 2012. WebPIE: A Web-scale Parallel Inference Engine using MapReduce. *JWS* 10:59–75.
- Urbani, J.; Jacobs, C. J.; and Krötzsch, M. 2016. Column-Oriented Datalog Materialization for Large Knowledge Graphs. In *AAAI*, 258–264.
- Wu, Z.; Eadon, G.; Das, S.; Chong, E. I.; Kolovski, V.; Annamalai, M.; and Srinivasan, J. 2008. Implementing an Inference Engine for RDFS/OWL Constructs and User-Defined Rules in Oracle. In *ICDE*, 1239–1248. IEEE.

A Appendix

A.1 Proof of Theorem 4

Theorem 4. *Algorithm 2 computes I as $\text{mat}(\Pi, E)$ if function $\text{Add}^{s,i}$ captures $\Pi_r^{s,i}$ in each of its calls. Moreover, if all $\text{Add}^{s,i}$ use the seminaïve strategy, each applicable rule instance is considered at most once.*

Proof. We first prove a property about each $\text{Add}^{s,i}$ function, which will be used later to establish the correctness of the algorithm. More specifically, for datasets I^p, I^n , and Δ with $\Delta \subseteq I^p$, let $\Delta_0 = \Delta$ and $J_0 = \emptyset$; for $k > 0$, let

$$\Delta_{k+1} = \Pi_r^{s,i}[I^p \cup J_k, I^n : \Delta_k] \setminus (I^p \cup J_k) \quad \text{and} \quad J_{k+1} = J_k \cup \Delta_{k+1};$$

and let $\text{Semi}[\Pi_r^{s,i}, I^p, I^n, \Delta] = J_k$ for k such that $J_k = J_{k+1}$. We show that $\text{Semi}[\Pi_r^{s,i}, I^p, I^n, \Delta]$ is the smallest set of facts J such that $\Pi_r^{s,i}[I^p \cup J, I^n : \Delta \cup J] \subseteq I^p \cup J$ holds—that is, property (6) holds.

$$\text{Add}^{s,i}[I^p, I^n, \Delta] = \text{Semi}[\Pi_r^{s,i}, I^p, I^n, \Delta] \tag{6}$$

To simplify the notation, let $J = \text{Semi}[\Pi_r^{s,i}, I^p, I^n, \Delta]$; we first prove that $\Pi_r^{s,i}[I^p \cup J, I^n : \Delta \cup J] \subseteq I^p \cup J$ holds. To this end, consider an arbitrary fact $F \in \Pi_r^{s,i}[I^p \cup J, I^n : \Delta \cup J]$. By the definition of the latter, there exist a rule $r \in \Pi_r^{s,i}$ and substitution σ such that $\mathbf{b}^+(r)\sigma \subseteq I^p \cup J$, $\mathbf{b}^-(r)\sigma \cap I^n = \emptyset$, $\mathbf{b}^+(r)\sigma \cap (\Delta \cup J) \neq \emptyset$, and $F = \mathbf{h}(r)\sigma$ all hold. Let k be the smallest index such that $\mathbf{b}^+(r)\sigma \subseteq I^p \cup J_k$ and $\mathbf{b}^+(r)\sigma \cap (\Delta \cup J_k) \neq \emptyset$ hold. Then, $\mathbf{b}^+(r)\sigma \cap \Delta_k \neq \emptyset$ must hold, or k is not the smallest such index. This implies $F \in \Pi_r^{s,i}[I^p \cup J_k, I^n : \Delta_k] \subseteq I^p \cup J$, as required.

Now we show that, for each dataset J' where $\Pi_r^{s,i}[I^p \cup J', I^n : \Delta \cup J'] \subseteq I^p \cup J'$ holds, we have $J \subseteq J'$ —that is, J is the smallest such dataset. To this end, we show by induction on k that $\Delta_k \subseteq J'$ holds for $k > 0$. For the induction base, $\Delta_0 = \Delta$ and $J_0 = \emptyset$ imply $\Delta_1 = \Pi_r^{s,i}[I^p, I^n : \Delta] \setminus I^p \subseteq \Pi_r^{s,i}[I^p \cup J', I^n : \Delta \cup J'] \setminus I^p \subseteq J'$. For the inductive step, consider arbitrary $k > 1$ where $\Delta_{k'} \subseteq J'$ holds for each k' with $1 \leq k' < k$. Then, we clearly have $J_{k-1} = \bigcup_{1 \leq k' < k} \Delta_{k'} \subseteq J'$. But then, $\Delta_k = \Pi_r^{s,i}[I^p \cup J_{k-1}, I^n : \Delta_{k-1}] \setminus (I^p \cup J_{k-1}) \subseteq \Pi_r^{s,i}[I^p \cup J', I^n : J'] \subseteq I^p \cup J'$ holds, as required.

We now proceed with the proof of our main claim. Let $I^0 = \emptyset$. Moreover, for each $1 \leq s \leq S$ where S is the largest stratum index, let I_0^s, I_1^s, \dots be the sequence of sets where $I_0^s = I^{s-1} \cup (E \cap \mathcal{O}^s)$, and $I_i^s = I_{i-1}^s \cup \Pi^s[I_{i-1}^s]$ for each $i > 0$. Index k clearly exists for which the fixpoint is reached (i.e., $I_k^s = I_{k+1}^s$ holds), so we let $I^s = I_k^s$. Finally, let $I = I^S$. It is straightforward to see that $I = \text{mat}(\Pi, E)$ —that is, I is the materialisation of Π w.r.t. E .

Consider a run of Algorithm 2 on Π, λ , and E . Let $I^0|_{\text{mod}} = \emptyset$, and for each $1 \leq s \leq S$, let $I^s|_{\text{mod}}$ be the value of I after the loop of lines 8–15 finishes for stratum s . We show by induction on s that property (7) holds for $0 \leq s \leq S$. Then, property (8) for $s = S$ and $I^S = I = \text{mat}(\Pi, E)$ jointly imply the correctness of the algorithm.

$$I^s|_{\text{mod}} = I^s \tag{7}$$

The base case where $s = 0$ is trivial since both sets are empty. For the inductive step, consider an arbitrary s with $1 \leq s \leq S$ such that (7) holds for $s - 1$, and we show that (7) holds for s as well. To this end, consider the execution of lines 9–15 for stratum s . For each $j > 0$, let $\Delta_i|_j$ and $\Delta|_j$ be the values of Δ_i (for $1 \leq i \leq n(s)$) and Δ when the j th iteration of lines 12–15 starts. We show that property (8) holds. Then, the way I is updated in line 12 ensures that property (7) holds.

$$I^{s-1}|_{\text{mod}} \cup \bigcup_{j>0} \Delta|_j = I^s \tag{8}$$

For the \subseteq direction of (8), $I^{s-1}|_{\text{mod}} = I^{s-1} \subseteq I^s$ holds by the induction assumption for (7). Next we prove $\bigcup_{j>0} \Delta|_j \subseteq I^s$ by induction on j .

- For the base case where $j = 1$, line 10 ensures that $\Delta_1 = (E \cap \mathcal{O}^s) \cup \Pi_{\text{nr}}^s[I^{s-1}|_{\text{mod}}]$. But then, the induction assumption $I^{s-1}|_{\text{mod}} = I^{s-1}$ and the definition of I^s jointly imply $\Delta_1 \subseteq I^s$.
- For the inductive step, consider arbitrary $j > 1$ such that $\Delta_k \subseteq I^s$ holds for each $1 \leq k < j$. Then, line 14 and the induction assumption for (7) ensure that $\Delta_i|_j = \text{Add}^{s,i}[I^{s-1} \cup \bigcup_{1 \leq k < j} \Delta|_k, I^{s-1} \cup \bigcup_{1 \leq k < j} \Delta|_k, \Delta|_{j-1} \setminus \Delta_i|_{j-1}]$. By property (6) we have $\Delta_i|_j = \text{Semi}[\Pi_r^{s,i}, I^{s-1} \cup \bigcup_{1 \leq k < j} \Delta|_k, I^{s-1} \cup \bigcup_{1 \leq k < j} \Delta|_k, \Delta|_{j-1} \setminus \Delta_i|_{j-1}]$. Then, the induction assumption and the definition of I^s imply $I^{s-1} \cup \bigcup_{1 \leq k < j} \Delta|_k \subseteq I^s$. Now let the sequences of Δ_m and J_m with $m \geq 0$ be defined in the same way as in the definition for the Semi function. We prove by induction on m that $\Delta_m \subseteq I^s$ and $J_m \subseteq I^s$ hold, then the definition of Semi implies $\Delta_i|_j \subseteq I^s$.
 - We have $\Delta_0 = \Delta|_{j-1} \setminus \Delta_i|_{j-1} \subseteq \Delta|_{j-1} \subseteq I^s$ and $J_0 = \emptyset \subseteq I^s$, so the induction base where $m = 0$ clearly holds.
 - For the inductive step, consider arbitrary $m > 0$ such that $\Delta_{m-1} \subseteq I^s$ and $J_{m-1} \subseteq I^s$ hold. But then, by definition we have $\Delta_m \subseteq \Pi_r^{s,i}[I^{s-1} \cup \bigcup_{1 \leq k < j} \Delta|_k \cup J_{m-1}, I^{s-1} \cup \bigcup_{1 \leq k < j} \Delta|_k] \subseteq \Pi_r^{s,i}[I^s, I^{s-1} \cup \bigcup_{1 \leq k < j} \Delta|_k]$. Facts in $I^s \setminus I^{s-1}$ all belong to stratum s , so they will not affect the evaluation of negative body atoms from rules in stratum s . Therefore we have $\Delta_m \subseteq \Pi_r^{s,i}[I^s, I^s] = \Pi_r^{s,i}[I^s] \subseteq I^s$, as required. Furthermore, by definition $J_m = J_{m-1} \cup \Delta_m$, together with the induction assumption $J_{m-1} \subseteq I^s$ this ensures that $J_m \subseteq I^s$ holds as well.

Now line 15 and the fact that $D_i|_j \subseteq I^s$ holds for each $1 \leq i \leq n(s)$ jointly imply $\Delta|_j = \bigcup_{1 \leq i \leq n(s)} \Delta_i|_j \subseteq I^s$, as required. This completes our proof for $\bigcup_{j>0} \Delta|_j \subseteq I^s$.

For the \supseteq direction of property (8), we prove by induction on i that $I_i^s \subseteq I^{s-1}|_{\text{mod}} \cup \bigcup_{j>0} \Delta|_j$ holds for $i \geq 0$.

- For the base case, we have $I_0^s = I^{s-1} \cup (E \cap \mathcal{O}^s)$. But then, line 10 ensures $E \cap \mathcal{O}^s \subseteq \Delta|_1$, which together with the induction assumption for (7) implies $I_0^s \subseteq I^{s-1}|_{\text{mod}} \cup \bigcup_{j>0} \Delta|_j$.
- For the induction step, consider arbitrary $i > 0$ such that $I_{i-1}^s \subseteq I^{s-1}|_{\text{mod}} \cup \bigcup_{j>0} \Delta|_j$ holds, and we would like to show that $I_i^s \subseteq I^{s-1}|_{\text{mod}} \cup \bigcup_{j>0} \Delta|_j$ holds as well. By the induction assumption for $i-1$ and the fact that $I_i^s = I_{i-1}^s \cup \Pi^s[I_{i-1}^s]$, it is enough to prove $\Pi^s[I_{i-1}^s] \subseteq I^{s-1}|_{\text{mod}} \cup \bigcup_{j>0} \Delta|_j$. To this end, consider arbitrary $F \in \Pi^s[I_{i-1}^s]$. There are two cases here. If $F \in \Pi_{\text{nr}}^s[I_{i-1}^s]$ —that is, F can be derived by a nonrecursive rule, then we have $F \in \Pi_{\text{nr}}^s[I^{s-1}]$. But then, the induction assumption for (7) and line 10 of the algorithm ensure $F \in \Delta|_1$. If $F \in \Pi_r^s[I_{i-1}^s]$, then there exists a module with index k such that $F \in \Pi_r^{s,k}[I_{i-1}^s]$. By the definition of rule application, there exist rule r and its instance r' such that $r \in \Pi_r^{s,k}$, $\text{b}^+(r') \subseteq I_{i-1}^s$, and $\text{b}^-(r') \cap I_{i-1}^s = \emptyset$ all hold. Since $I_{i-1}^s \subseteq I^{s-1}|_{\text{mod}} \cup \bigcup_{j>0} \Delta|_j$ holds by the induction assumption, let j' be the largest index j such that $\text{b}^+(r') \cap \Delta|_j \neq \emptyset$. Then, if $\text{b}^+(r') \cap (\Delta|_{j'} \setminus \Delta_k|_{j'}) \neq \emptyset$, then Definition 3 ensures that F is added to $\Delta_k|_{j'+1}$ in line 14 during the execution of the j' th iteration of lines 12–15; if $\text{b}^+(r') \cap (\Delta|_{j'} \setminus \Delta_k|_{j'}) = \emptyset$, then we have $\text{b}^+(r') \subseteq I^{s-1}|_{\text{mod}} \cup \bigcup_{0 < j < j'} \Delta|_j \cup \Delta_k|_{j'}$, so Definition 3 ensures that F is added to $\Delta_k|_{j'}$ in line 14 during the execution of the $(j' - 1)$ th iteration of lines 12–15. Either way, we have $F \in I^{s-1}|_{\text{mod}} \cup \bigcup_{j>0} \Delta|_j$. Since the choice of F is arbitrary, we have $\Pi_r^s[I_{i-1}^s] \subseteq I^{s-1}|_{\text{mod}} \cup \bigcup_{j>0} \Delta|_j$, as required.

This completes our proof for the correctness of the algorithm.

Next we show that if all $\text{Add}^{s,i}$ use the seminaïve strategy—that is, each $\text{Add}^{s,i}[I^p, I^n, \Delta]$ is computed in the same way that $\text{Semi}[\Pi_r^{s,i}, I^p, I^n, \Delta]$ is constructed, then each applicable rule instance is considered at most once. To this end, consider a run of Algorithm 2. First, please note that the program is processed in a stratum-by-stratum manner, so no applicable rule instance will be considered in two distinct iterations of lines 8–15. Now consider an arbitrary stratum index s and the iteration of lines 8–15 for s , the only places that consider rule instances are lines 10 and 14. Line 10 handles nonrecursive rules whereas line 14 handles recursive rules, so no rule instance will be considered in both places. Line 10 is only executed once for s while line 14 can be executed multiple times. Thus it is sufficient to show that for stratum s , line 14 never repeatedly consider an applicable rule instance. The way $\text{Semi}[\Pi_r^{s,i}, I^p, I^n, \Delta]$ is constructed ensures that one application of $\text{Add}^{s,i}$ does not repeat rule instances itself. Now consider the m th and the n th iterations of lines 12–15 where we have $m \neq n$. We would like to show that for each i , the application of $\text{Add}^{s,i}[I^{s-1}|_{\text{mod}} \cup \bigcup_{0 < j \leq m} \Delta|_j, I^{s-1}|_{\text{mod}} \cup \bigcup_{0 < j \leq m} \Delta|_j, \Delta|_m \setminus \Delta_i|_m]$ and the application of $\text{Add}^{s,i}[I^{s-1}|_{\text{mod}} \cup \bigcup_{0 < j \leq n} \Delta|_j, I^{s-1}|_{\text{mod}} \cup \bigcup_{0 < j \leq n} \Delta|_j, \Delta|_n \setminus \Delta_i|_n]$ do not repeat rule instances. Without loss of generality assume that $m < n$ holds. By the construction of $\text{Semi}[\Pi_r^{s,i}, I^p, I^n, \Delta]$ we know that for the former, each applicable rule instance must have at least one positive body atom in $(\Delta|_m \setminus \Delta_i|_m) \cup \Delta_i|_{m+1}$, and for the latter, each applicable rule instance must have one positive body atom in $(\Delta|_n \setminus \Delta_i|_n) \cup \Delta_i|_{n+1}$. It is straightforward to see that these two sets are disjoint, so no applicable rule instance will be considered by these applications of Add , and this completes our proof for the second half of Theorem 4. \square

A.2 Proof of Theorem 8

Theorem 8. *Algorithm 3 updates I from $\text{mat}(\Pi, E)$ to $\text{mat}(\Pi, (E \setminus E^-) \cup E^+)$ if functions $\text{Add}^{s,i}$, $\text{Del}^{s,i}$, $\text{Diff}^{s,i}$, and $\text{Red}^{s,i}$ capture $\Pi_r^{s,i}$ in all of their calls.*

Proof. For a program Π and datasets I^p , I^n , and Δ with $\Delta \subseteq I^p$, let $\Delta_0 = \Delta$ and $J_0 = \emptyset$; moreover, for $i > 0$, let $\Delta_{i+1} = \Pi[I^p \setminus J_i, I^n : \Delta_i] \setminus (\Delta \cup J_i)$ and $J_{i+1} = J_i \cup \Delta_i$; let $\text{InvSemi}[\Pi, I^p, I^n, \Delta] = J_i \setminus \Delta$ for i such that $J_i = J_{i+1}$ holds. We show that $J = \text{InvSemi}[\Pi, I^p, I^n, \Delta]$ is the smallest dataset satisfying $\Pi[I^p, I^n : \Delta \cup J] \subseteq \Delta \cup J$.

We show by induction on $i \geq 0$ that $\Pi[I^p, I^n : \Delta \cup J_i] \subseteq \Delta \cup J$ holds. The base case where $i = 0$ trivially holds. For the inductive step, consider arbitrary $i > 0$ such that $\Pi[I^p, I^n : \Delta \cup J_{i-1}] \subseteq \Delta \cup J$ holds. For each $F \in \Pi[I^p, I^n : \Delta \cup J_i]$, there exist rule $r \in \Pi$ and substitution σ such that $\text{b}^+(r)\sigma \in I^p$, $\text{b}^-(r)\sigma \cap I^n = \emptyset$, $\text{b}^+(r)\sigma \cap (\Delta \cup J_i) \neq \emptyset$, and $F = \text{h}(r)\sigma$ all hold. If $\text{b}^+(r)\sigma \cap (\Delta \cup J_{i-1}) \neq \emptyset$, then the induction assumption ensures $F \in \Delta \cup J$. Otherwise $\text{b}^+(r)\sigma \cap \Delta_{i-1} \neq \emptyset$ holds. Now there are two possibilities: if $\text{b}^+(r)\sigma \subseteq I^p \setminus J_{i-1}$, then F is derived in the construction of Δ_i ; if $\text{b}^+(r)\sigma \cap J_{i-1} \neq \emptyset$, then the induction assumption ensures $F \in \Delta \cup J$. Therefore, $\Pi[I^p, I^n : \Delta \cup J] \subseteq \Delta \cup J$ holds, as required.

To see that J is the smallest such set, let J' be an arbitrary set satisfying $\Pi[I^p, I^n : \Delta \cup J'] \subseteq \Delta \cup J'$ and we prove by induction on i that $\bigcup_{0 \leq j \leq i} \Delta_j \setminus \Delta \in J'$ holds. The base case where $i = 0$ clearly holds since $\Delta_0 \setminus \Delta = \emptyset$. For the inductive step, consider arbitrary index $i > 0$ such that $\bigcup_{0 \leq j \leq i-1} \Delta_j \setminus \Delta \subseteq J'$ holds. This implies $\Delta_{i-1} \subseteq \Delta \cup J'$, which together with

the definition of Δ_i ensures $\Delta_i \subseteq \Delta \cup J'$, so the inductive step holds. Therefore, $J = \text{InvSemi}[\Pi, I^p, I^n, \Delta]$ is the smallest dataset satisfying $\Pi[I^p, I^n : \Delta \cup J] \subseteq \Delta \cup J$.

For a program Π , datasets I^p, I^n , and Δ with $\Delta \subseteq I^p$, and a mapping C_{nr} of facts to nonnegative integers, let $\Delta_0 = \Delta$ and $J_0 = \emptyset$; moreover, for $i > 0$, let

$$\Delta_{i+1} = \{F \in \Pi[I^p \setminus J_i, I^n : \Delta_i] \setminus (\Delta \cup J_i) \mid C_{nr}(F) = 0\} \quad \text{and} \quad J_{i+1} = J_i \cup \Delta_i;$$

and let $\text{InvSemi}^c[\Pi, I^p, I^n, \Delta, C_{nr}] = J_i \setminus \Delta$ for i such that $J_i = J_{i+1}$. We show that $J = \text{InvSemi}^c[\Pi, I^p, I^n, \Delta, C_{nr}]$ is the smallest dataset that satisfies the following: for each $F \in \Pi[I^p, I^n : \Delta \cup J]$, either $F \in \Delta \cup J$ or $C_{nr}(F) > 0$ holds.

We prove induction on $i \geq 0$ that for each $F \in \Pi[I^p, I^n : \Delta \cup J_i]$, either $F \in \Delta \cup J$ or $C_{nr}(F) > 0$ holds. The base case where $i = 0$ clearly holds. For the inductive step, consider arbitrary i such that for each $F \in \Pi[I^p, I^n : \Delta \cup J_{i-1}]$, either $F \in \Delta \cup J$ or $C_{nr}(F) > 0$ holds. Then, for each fact $G \in \Pi[I^p, I^n : \Delta \cup J_i]$, there exist a rule $r \in \Pi$ and a substitution σ such that $b^+(r)\sigma \subseteq I^p$, $b^-(r)\sigma \cap I^n = \emptyset$, $b^+(r)\sigma \cap (\Delta \cup J_i) \neq \emptyset$, and $G = h(r)\sigma$ all hold. Now if $b^+(r)\sigma \cap (\Delta \cup J_{i-1})$, then the induction assumption ensures that either $G \in \Delta \cup J$ or $C_{nr}(G) > 0$ holds. Otherwise we have $b^+(r)\sigma \cap (J_i \setminus J_{i-1}) = b^+(r)\sigma \cap \Delta_{i-1} \neq \emptyset$. There are two possibilities: if $b^+(r)\sigma \subseteq I^p \setminus J_{i-1}$, then G is derived in the construction of Δ_i and it is either guaranteed to be in $\Delta \cup J$ or we have $C_{nr}(G) > 0$; if $b^+(r)\sigma \cap J_{i-1} \neq \emptyset$, then the induction assumption ensures that either $G \in \Delta \cup J$ or $C_{nr}(G) > 0$ holds.

To see that J is the smallest such set, let J' be an arbitrary set satisfying the following: for each $F \in \Pi[I^p, I^n : \Delta \cup J']$, either $F \in \Delta \cup J'$ or $C_{nr}(F) > 0$ holds. We prove by induction on i that $\Delta_i \setminus \Delta \subseteq J'$ holds. The base case where $i = 0$ clearly holds since $\Delta_0 \setminus \Delta = \emptyset$. For the inductive step, consider arbitrary index $i > 0$ such that $\Delta_{i-1} \setminus \Delta \subseteq J'$ holds. This implies $\Delta_{i-1} \subseteq \Delta \cup J'$, which together with the definition of Δ_i and the induction assumption ensures $\Delta_i \subseteq \Delta \cup J'$, so the inductive step holds. Therefore, $J = \text{InvSemi}^c[\Pi, I^p, I^n, \Delta, C_{nr}]$ is indeed the smallest dataset satisfying the above property.

To see that our main claim holds, note that $\text{InvSemi}[\Pi, I^p, I^n, \Delta]$ captures overdeletion in DRed and corresponds to the upper bound J_u in Definition 6, whereas $\text{InvSemi}^c[\Pi, I^p, I^n, \Delta, C_{nr}]$ captures overdeletion in DRed^c and corresponds to the lower bound J_l in Definition 6. Then, the correctness of Algorithm 3 follows from the correctness of DRed and DRed^c . \square

A.3 Proof of Theorem 9

Theorem 9. *In each call in Algorithms 2 and 3, functions $\text{Add}^{\text{tc}(R)}$, $\text{Del}^{\text{tc}(R)}$, $\text{Diff}^{\text{tc}(R)}$, and $\text{Red}^{\text{tc}(R)}$ capture a datalog program that axiomatises relation R as transitive.*

We consider each of these functions in combination with each of the relevant algorithms in a separate claim. We first consider Algorithm 2. For I a dataset and R a predicate, let $R[I]$ denote the set of all R facts in I .

Claim 11. *If $\text{mat}(\Pi^{\text{tc}(R)}, X_R) = R[I^p \setminus \Delta]$ holds before each call to function $\text{Add}^{\text{tc}(R)}[I^p, I^n, \Delta]$ in a run of Algorithm 2, then $\text{mat}(\Pi^{\text{tc}(R)}, X_R \cup \Delta) = R[I^p] \cup J$ holds, where $J = \text{Add}^{\text{tc}(R)}[I^p, I^n, \Delta]$. Moreover, $\text{Semi}[\Pi^{\text{tc}(R)}, I^p, I^n, \Delta] = J$ holds.*

Proof. First we show that $\text{mat}(\Pi^{\text{tc}(R)}, X_R \cup \Delta) = R[I^p] \cup J$ holds. For the \subseteq direction of the equation, consider arbitrary $R(u, v) \in \text{mat}(\Pi^{\text{tc}(R)}, X_R \cup \Delta)$. If $R(u, v) \in \text{mat}(\Pi^{\text{tc}(R)}, X_R)$, then $R(u, v) \in R[I^p \setminus \Delta] \subseteq R[I^p] \cup J$ clearly holds. Otherwise, there exists a shortest chain of R facts in $X_R \cup \Delta$ that connects u and v . We show by induction on the length of this chain that $R(u, v) \in R[I^p] \cup J$ holds and $R(u, v)$ is added to Q at some point during the execution of the algorithm. For the base case where the length is one, since $R(u, v) \notin \text{mat}(\Pi^{\text{tc}(R)}, X_R)$, we have $R(u, v) \notin X_R$. Thus, $R(u, v) \in \Delta \in I^p$ holds, and $R(u, v)$ is added to Q in line 41. For the inductive step, consider a length $i + 1$ chain $R(a_0, a_1), \dots, R(a_i, a_{i+1})$. We consider two cases. In the first case, we have $R(a_0, a_1) \in \Delta$ and $R(a_1, a_{i+1}) \in \text{mat}(\Pi^{\text{tc}(R)}, X_R) = R[I^p \setminus \Delta]$. But then, $R(a_0, a_{i+1})$ is added to Q and J in line 43. In the second case, we have $R(a_0, a_1) \in X_R \cup \Delta$, $R(a_1, a_{i+1}) \in \text{mat}(\Pi^{\text{tc}(R)}, X_R \cup \Delta)$, and $R(a_1, a_{i+1}) \notin \text{mat}(\Pi^{\text{tc}(R)}, X_R)$. But then, there exists a chain of length i that derives $R(a_1, a_{i+1})$, which by the induction assumption ensures that $R(a_1, a_{i+1})$ is added to Q at some point during the execution of the algorithm. Lines 46 and 47 then ensure that $R(a_0, a_{i+1})$ is added to Q and that $R(a_0, a_{i+1}) \in R[I^p] \cup J$ holds. We next prove $J = \text{Semi}[\Pi^{\text{tc}(R)}, I^p, I^n, \Delta]$. Then $\text{mat}(\Pi^{\text{tc}(R)}, X_R) = R[I^p \setminus \Delta]$ implies $\text{mat}(\Pi^{\text{tc}(R)}, I^p \setminus \Delta) = R[I^p \setminus \Delta]$; similarly, $\text{mat}(\Pi^{\text{tc}(R)}, X_R \cup \Delta) = R[I^p] \cup J$ implies $\text{mat}(\Pi^{\text{tc}(R)}, I^p) = \text{mat}(\Pi^{\text{tc}(R)}, I^p \cup J) = R[I^p] \cup J$. Moreover, it can be easily verified by induction on the construction of $\text{Semi}[\Pi^{\text{tc}(R)}, I^p, I^n, \Delta]$ that $\text{mat}(\Pi^{\text{tc}(R)}, I^p \setminus \Delta) \cup R[\Delta] \cup \text{Semi}[\Pi^{\text{tc}(R)}, I^p, I^n, \Delta] = \text{mat}(\Pi^{\text{tc}(R)}, I^p)$ —that is, the seminaïve computation correctly closes I^p with respect to $\Pi^{\text{tc}(R)}$. Thus, $J = \text{Semi}[\Pi^{\text{tc}(R)}, I^p, I^n, \Delta]$ holds, as required. \square

That $\text{Add}^{\text{tc}(R)}$ captures $\Pi^{\text{tc}(R)}$ during the execution of Algorithm 2 follows from Claim 11 and property (6). We next show that, during the execution of Algorithm 3, functions $\text{Add}^{\text{tc}(R)}$, $\text{Del}^{\text{tc}(R)}$, $\text{Diff}^{\text{tc}(R)}$, and $\text{Red}^{\text{tc}(R)}$ capture $\Pi^{\text{tc}(R)}$. Note that the order of the function calls is important for correctness, so we examine $\text{Diff}^{\text{tc}(R)}$ first. For each call to function $\text{Diff}^{\text{tc}(R)}[I^p, \Delta^p, \Delta^n]$ in line 25, we have $\Delta^p = D \setminus A$, which contains only facts from previous strata. Thus, $\Pi^{\text{tc}(R)}[I^p : \Delta^p, \Delta^n] = \emptyset$ clearly holds since the only rule in $\Pi^{\text{tc}(R)}$ is recursive. Our implementation for $\text{Diff}^{\text{tc}(R)}[I^p, \Delta^p, \Delta^n]$ always return empty set as well, so by

Definition 5 $\text{Diff}^{\text{tc}(R)}$ captures $\Pi^{\text{tc}(R)}$ for the calls in line 25. For the same reason, $\text{Diff}^{\text{tc}(R)}$ captures $\Pi^{\text{tc}(R)}$ for the calls in line 35 as well. We next focus on $\text{Del}^{\text{tc}(R)}$, $\text{Red}^{\text{tc}(R)}$, and $\text{Add}^{\text{tc}(R)}$.

Claim 12. *If $\text{mat}(\Pi^{\text{tc}(R)}, X_R) = \text{mat}(\Pi^{\text{tc}(R)}, I^P)$ holds before each call to $\text{Del}^{\text{tc}(R)}[I^P, I^n, \Delta, C_{nr}]$ in a run of Algorithm 3, then $\text{mat}(\Pi^{\text{tc}(R)}, X_R \setminus \Delta) = \text{mat}(\Pi^{\text{tc}(R)}, (I^P \setminus \Delta) \setminus J)$ holds, where $J = \text{Del}^{\text{tc}(R)}[I^P, I^n, \Delta, C_{nr}]$. Moreover, $J_l \subseteq J \subseteq J_u$ holds, where J_l and J_u are the lower and upper bounds from Definition 6, respectively.*

Proof. The proof is analogous to the proof of Claim 11 and is based on the intuition that the function implements seminaïve evaluation for the rule $X(x, y) \wedge R(y, z) \rightarrow R(x, z)$. \square

That $\text{Del}^{\text{tc}(R)}$ captures $\Pi^{\text{tc}(R)}$ during the execution of Algorithm 3 immediately follows from Claim 12 and Definition 6.

Claim 13. *If $\text{mat}(\Pi^{\text{tc}(R)}, X_R) = \text{mat}(\Pi^{\text{tc}(R)}, I^P \setminus \Delta)$ holds before each call to $\text{Red}^{\text{tc}(R)}[I^P, I^n, \Delta]$ in a run of Algorithm 3, then $\text{mat}(\Pi^{\text{tc}(R)}, X_R) = R[(I^P \setminus \Delta)] \cup J$ holds, where $J = \text{Red}^{\text{tc}(R)}[I^P, I^n, \Delta]$. Moreover, J is the smallest dataset satisfying $\Pi^{\text{tc}(R)}[(I^P \setminus \Delta) \cup J, I^n] \cap \Delta \subseteq J$.*

Proof. First we show that $\text{mat}(\Pi^{\text{tc}(R)}, X_R) = R[(I^P \setminus \Delta)] \cup J$ holds. The \supseteq direction of the property is trivial: each fact $R(u, v)$ added to J in line 62 corresponds to a chain of facts in X_R , so $R(u, v) \in \text{mat}(\Pi^{\text{tc}(R)}, X_R)$ holds. For the \subseteq direction, consider arbitrary fact $R(u, v)$ such that $R(u, v) \in \text{mat}(\Pi^{\text{tc}(R)}, X_R)$ and $R(u, v) \notin R[(I^P \setminus \Delta)]$ hold. Since the original materialisation I is passed as I^P when the function gets called, we have $R(u, v) \in I^P \setminus (I^P \setminus \Delta) = \Delta$. But then, lines 60–62 ensure that a chain of facts in X_R deriving $R(u, v)$ will be found and that $R(u, v)$ is added to J . We now show that J is the smallest dataset satisfying $\Pi^{\text{tc}(R)}[(I^P \setminus \Delta) \cup J, I^n] \cap \Delta \subseteq J$. Consider arbitrary fact $R(u, v) \in \Pi^{\text{tc}(R)}[(I^P \setminus \Delta) \cup J, I^n] \cap \Delta$. By $\text{mat}(\Pi^{\text{tc}(R)}, X_R) = R[(I^P \setminus \Delta)] \cup J$, we have $R(u, v) \in \Delta \cap \text{mat}(\Pi^{\text{tc}(R)}, X_R)$. But then, lines 60–62 ensure $R(u, v) \in J$. Moreover, it is straightforward to show by induction that any J' satisfying $\Pi^{\text{tc}(R)}[(I^P \setminus \Delta) \cup J', I^n] \cap \Delta \subseteq J'$ must at least contain all facts in $\text{mat}(\Pi^{\text{tc}(R)}, I^P \setminus \Delta) \cap \Delta = \text{mat}(\Pi^{\text{tc}(R)}, X_R) \setminus (I^P \setminus \Delta) = J$ —in other words, J is indeed the smallest dataset satisfying $\Pi^{\text{tc}(R)}[(I^P \setminus \Delta) \cup J, I^n] \cap \Delta \subseteq J$. \square

That $\text{Red}^{\text{tc}(R)}$ captures $\Pi^{\text{tc}(R)}$ during the execution of Algorithm 3 immediately follows from Claim 13 and Definition 7. Finally, the proof of $\text{Add}^{\text{tc}(R)}$ capturing $\Pi^{\text{tc}(R)}$ during the execution of Algorithm 3 is analogous to the proof of Claim 11 so we omit the details for the sake of brevity.

A.4 Proof of Theorem 10

Theorem 10. *In each call in Algorithms 2 and 3, functions $\text{Add}^{\text{stc}(R)}$, $\text{Del}^{\text{stc}(R)}$, $\text{Diff}^{\text{stc}(R)}$, and $\text{Red}^{\text{stc}(R)}$ capture a datalog program that axiomatises R as symmetric–transitive.*

We consider each of these functions in combination with each of the relevant algorithms in a separate claim. We first consider Algorithm 2. For C_R a set of sets (representing a set of connected components of edges consisting of relation R), let $\text{Close}(C_R) = \bigcup_{U \in C_R} \bigcup_{u, v \in U} \{R(u, v)\}$.

Claim 14. *If $\text{Close}(C_R) = R[I^P \setminus \Delta]$ holds before each call to $\text{Add}^{\text{stc}(R)}[I^P, I^n, \Delta]$ in a run of Algorithm 2, then C_R is updated so that $\text{Close}(C_R) = R[I^P] \cup J$ holds, where $J = \text{Add}^{\text{stc}(R)}[I^P, I^n, \Delta]$. Moreover, $J = \text{Semi}[\Pi^{\text{stc}(R)}, I^P, I^n, \Delta]$ holds.*

Proof. First we show that $\text{Close}(C_R) = R[I^P] \cup J$ holds after the function call. For the \subseteq direction, please note that C_R can be updated in only three places—lines 69, 71, and 73. In line 69 the first command adds a new component $U = \{u\}$ to C_R . This will add $R(u, u)$ to $\text{Close}(C_R)$. But then, the second command in line 69 ensures that $R(u, u) \in I^P \cup J$ holds. The same reasoning applies to line 71. In line 73 two components U and V are merged; but then, lines 74–75 ensure that the affected facts are added to $I^P \cup J$. Thus, $\text{Close}(C_R) \subseteq R[I^P] \cup J$ holds after the update. For the \supseteq direction, please note that the algorithm ensures that for each fact $R(u, v) \in R[\Delta] \cup J$, u and v are in the same component in C_R after the update; moreover, $R[I \setminus \Delta] \subseteq \text{Close}(C_R)$ already holds before the update; so the \supseteq direction of the property $\text{Close}(C_R) = R[I^P] \cup J$ also holds after the update. Next we show that $J = \text{Semi}[\Pi^{\text{stc}(R)}, I^P, I^n, \Delta]$ holds. Before the function is executed, we have $\text{Close}(C_R) = R[I^P \setminus \Delta]$, which implies $\text{mat}(\Pi^{\text{stc}(R)}, I^P \setminus \Delta) = R[I^P \setminus \Delta]$; similarly, after the function call we have $\text{Close}(C_R) = R[I^P] \cup J$, which implies $\text{mat}(\Pi^{\text{stc}(R)}, I^P) = \text{mat}(\Pi^{\text{stc}(R)}, I^P \cup J) = R[I^P] \cup J$. Moreover, it can be easily verified by induction on the construction of $\text{Semi}[\Pi^{\text{stc}(R)}, I^P, I^n, \Delta]$ that $\text{mat}(\Pi^{\text{stc}(R)}, I^P \setminus \Delta) \cup R[\Delta] \cup \text{Semi}[\Pi^{\text{stc}(R)}, I^P, I^n, \Delta] = \text{mat}(\Pi^{\text{stc}(R)}, I^P)$ —that is, the seminaïve computation correctly closes I^P with respect to $\Pi^{\text{stc}(R)}$. Therefore, we have $J = \text{Semi}[\Pi^{\text{stc}(R)}, I^P, I^n, \Delta]$, as required. \square

The fact that $\text{Add}^{\text{stc}(R)}$ captures $\Pi^{\text{stc}(R)}$ during the execution of Algorithm 2 directly follows from Claim 14 and property (6). Next we show that during the execution of Algorithm 3, functions $\text{Add}^{\text{stc}(R)}$, $\text{Del}^{\text{stc}(R)}$, $\text{Diff}^{\text{stc}(R)}$, and $\text{Red}^{\text{stc}(R)}$ capture $\Pi^{\text{stc}(R)}$. Note that the order of the function calls is important for correctness, so we examine $\text{Diff}^{\text{stc}(R)}$ first. For each call to function $\text{Diff}^{\text{stc}(R)}[I^P, \Delta^p, \Delta^n]$ in line 25, we have $\Delta^p = D \setminus A$, which contains only facts from previous strata. Thus

$\Pi^{\text{stc}(R)}[I^p : \Delta^p, \Delta^n] = \emptyset$ clearly holds since both rules in $\Pi^{\text{stc}(R)}$ are recursive. Our implementation for $\text{Diff}^{\text{stc}(R)}[I^p, \Delta^p, \Delta^n]$ always return empty set as well, so by Definition 5 $\text{Diff}^{\text{stc}(R)}$ captures $\Pi^{\text{stc}(R)}$ for the calls in line 25. For the same reason, $\text{Diff}^{\text{stc}(R)}$ captures $\Pi^{\text{stc}(R)}$ for the calls in line 35 as well. We next focus on $\text{Del}^{\text{stc}(R)}$, $\text{Red}^{\text{stc}(R)}$, and $\text{Add}^{\text{stc}(R)}$.

Claim 15. *If $\text{Close}(C_R) \cup Y_R = R[I^p]$ holds before each call to $\text{Del}^{\text{stc}(R)}[I^p, I^n, \Delta, C_{nr}]$ in a run of Algorithm 3, then C_R and Y_R are updated so that $\text{Close}(C_R) \cup Y_R = R[I^p \setminus \Delta] \setminus J$ holds, where $J = \text{Del}^{\text{stc}(R)}[I^p, I^n, \Delta, C_{nr}]$. Moreover, $J_l \subseteq J \subseteq J_u$ holds, where J_l and J_u are the lower and upper bounds from Definition 6, respectively.*

Proof. First we show that $\text{Close}(C_R) \cup Y_R = R[I^p \setminus \Delta] \setminus J$ holds after the update. For the \subseteq direction, consider arbitrary $R(u, v) \in \text{Close}(C_R) \cup Y_R$ after the update. If $R(u, v) \in \text{Close}(C_R)$, since the algorithm only removes components from C_R , we have $R(u, v) \in I^p$. Moreover, lines 78 and 82, and the fact that the components in C_R are disjoint ensure that the component containing u and v is not removed from C_R during the update. But then, by line 78 and line 80 we know that $R(u, v) \notin \Delta \cup J$. Therefore, $R(u, v) \in R[I^p \setminus \Delta] \setminus J$ holds. If $R(u, v) \in Y_R$ after the update, then we have $C_{nr}(R(u, v)) > 0$, so line 80 ensures that $R(u, v) \notin J$ holds; moreover, $R(u, v) \notin \Delta$ holds since otherwise there exists another $\text{Del}^{\text{stc}(R)}$ that violates the lower bound constraint. Therefore, $\text{Close}(C_R) \cup Y_R \subseteq R[I^p \setminus \Delta] \setminus J$ holds after the update.

For the \supseteq direction, consider arbitrary $R(u, v) \in R[I^p \setminus \Delta] \setminus J$. By $R[I^p \setminus \Delta] \setminus J \subseteq R[I^p]$ and $\text{Close}(C_R) \cup Y_R = R[I^p]$ we have $R(u, v) \in \text{Close}(C_R) \cup Y_R$ before the update. No fact is removed from Y_R during the update, so if $R(u, v) \in Y_R$ before the update, the same still holds after the update. Now if $R(u, v) \in \text{Close}(C_R)$ before the update, then there are two cases. If the component in C_R containing u and v is not removed, then clearly we still have $R(u, v) \in \text{Close}(C_R)$ after the update. If the component is indeed removed, then we have $C_{nr}(R(u, v)) > 0$ since otherwise we would have $R(u, v) \in \Delta \cup J$, which leads to a contradiction. But then, $R(u, v)$ is added to Y_R , so $R(u, v) \in \text{Close}(C_R) \cup Y_R$ holds as well after the function call.

Now we show that $\text{InvSemi}^c[\Pi^{\text{stc}(R)}, I^p, I^n, \Delta, C_{nr}] = J_l \subseteq J \subseteq J_u = \text{InvSemi}[\Pi^{\text{stc}(R)}, I^p, I^n, \Delta]$ holds. For the right-hand inclusion $J \subseteq J_u$, consider arbitrary $R(u, v) \in J$. The fact is added to J in line 80 since there exists a fact $R(u', v')$ such that $R(u', v') \in \Delta$ holds, and u', v' are in the same component as u, v before the function call. But then, by $\text{Close}(C_R) \cup Y_R = R[I^p]$ we clearly have $R(u, u') \in I^p$ and $R(v', v) \in I^p$. It is straightforward to verify by induction on the construction of $\text{InvSemi}[\Pi^{\text{stc}(R)}, I^p, I^n, \Delta]$ that applicable rule instances $R(u, u') \wedge R(u', v') \rightarrow R(u, v')$ and $R(u, v') \wedge R(v', v) \rightarrow R(u, v)$ will be considered, so $R(u, v) \in J_u$ holds. For the left-hand inclusion $J_l \subseteq J$, consider arbitrary $R(u, v) \in \text{InvSemi}^c[\Pi^{\text{stc}(R)}, I^p, I^n, \Delta, C_{nr}]$, by the definition of InvSemi and the fact that $\text{Close}(C_R) \cup Y_R = R[I^p]$ holds before the update, we know that there exists a component $U \in C_R$ that contains u, v, u', v' such that $R(u', v') \in \Delta$. But then, line 78 ensures that U is deleted and $R(u, v)$ is checked in line 80. $R(u, v) \in J_l$ ensures $C_{nr}(R(u, v)) = 0$, since otherwise J_l would not satisfy the lower bound condition in Definition 6. Therefore, line 80 ensures that $R(u, v)$ is in J (and it is not in Δ due to $R(u, v) \in J_l$ and the definition of InvSemi). \square

That $\text{Del}^{\text{stc}(R)}$ captures $\Pi^{\text{stc}(R)}$ during the execution of Algorithm 3 immediately follows from Claim 15 and Definition 6.

Claim 16. *If $\text{Close}(C_R) \cup Y_R = R[I^p \setminus \Delta]$ holds before each call to $\text{Red}^{\text{stc}(R)}[I^p, I^n, \Delta]$ in a run of Algorithm 3, then C_R is updated so that $\text{Close}(C_R) = R[I^p \setminus \Delta] \cup J$ holds, where $J = \text{Red}^{\text{stc}(R)}[I^p, I^n, \Delta]$. Moreover, J is the smallest dataset satisfying $\Pi^{\text{stc}(R)}[(I^p \setminus \Delta) \cup J, I^n] \cap \Delta \subseteq J$.*

Proof. $\text{Close}(C_R) \cup Y_R = R[I^p \setminus \Delta]$ ensures $\text{Close}(C_R) = R[(I^p \setminus \Delta) \setminus Y_R]$. But then, $\text{Close}(C_R) = R[I^p \setminus \Delta] \cup J$ holds after the update in the same way as in the proof for Claim 14. Moreover, J is the smallest dataset that satisfies

$$\Pi^{\text{stc}(R)}[(I^p \setminus \Delta) \cup J, I^n : Y_R \cup J] \subseteq (I^p \setminus \Delta) \cup J.$$

Now to see that J is also the smallest dataset satisfying $\Pi^{\text{stc}(R)}[(I^p \setminus \Delta) \cup J, I^n] \cap \Delta \subseteq J$, assume for the sake of a contradiction that a smaller dataset J' satisfies the above—that is, $\Pi^{\text{stc}(R)}[(I^p \setminus \Delta) \cup J', I^n] \cap \Delta \subseteq J'$ holds. Now consider arbitrary $F \in \Pi^{\text{stc}(R)}[(I^p \setminus \Delta) \cup J', I^n : Y_R \cup J'] \subseteq \Pi^{\text{stc}(R)}[(I^p \setminus \Delta) \cup J', I^n]$. Since the original materialisation I is passed as I^p when the function gets called, it is clear that $F \in I^p = I$ holds. Now if $F \notin \Delta$, then by our assumption we have $F \in J'$. But then, J' also satisfies $\Pi^{\text{stc}(R)}[(I^p \setminus \Delta) \cup J', I^n : Y_R \cup J'] \subseteq (I^p \setminus \Delta) \cup J'$ and J' is smaller than J , which is a contradiction. \square

That $\text{Red}^{\text{stc}(R)}$ captures $\Pi^{\text{stc}(R)}$ during the execution of Algorithm 3 follows from Claim 16 and Definition 7. Finally, the proof of $\text{Add}^{\text{stc}(R)}$ capturing $\Pi^{\text{stc}(R)}$ during the execution of Algorithm 3 is analogous to the proof of Claim 14 so we omit the details for the sake of brevity.