



HAL
open science

Jumping Evaluation of Nested Regular Path Queries

Joachim Niehren, Sylvain Salvati, Rustam Azimov

► **To cite this version:**

Joachim Niehren, Sylvain Salvati, Rustam Azimov. Jumping Evaluation of Nested Regular Path Queries. 38th International Conference on Logic Programming (ICLP'2022), Jul 2022, Haifa, Israel. hal-02492780v3

HAL Id: hal-02492780

<https://inria.hal.science/hal-02492780v3>

Submitted on 16 Mar 2022 (v3), last revised 4 Aug 2022 (v6)

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

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

Jumping Evaluation of Nested Regular Path Queries

JOACHIM NIEHREN, SYLVAIN SALVATI

Inria Lille, Université de Lille, France

RUSTAM AZIMOV

JetBrains Research, Saint Petersburg State University, Russia

Abstract

Nested regular path queries are used for querying graph databases and RDF triple stores. We propose a new algorithm for evaluating nested regular path queries. Not only does it evaluate path queries from a given set of start nodes in combined linear time, but also this complexity bound depends only on the size of the query’s *top-down needed* subgraph, a notion that we introduce. For many queries relevant in practice, the top-down needed subgraph is way smaller than the whole datagraph. Our algorithm is based on a novel compilation schema from nested regular path queries to monadic datalog queries. We prove that the top-down evaluation of the datalog program visits only the top-down needed subgraph for the path query. Thereby, the combined linear time complexity depending on the size of the top-down needed subgraph is implied by a general complexity result for top-down datalog evaluation. As an application, we show that our algorithm permits to reformulate in simple terms a variant of a very efficient automata-based algorithm proposed by Maneth and Nguyen that evaluates navigational path queries in datatrees based on indexes and jumping. Moreover, our variant overcomes some limitations of Maneth and Nguyen’s: it is not bound to trees and applies to graphs; it is not limited to forward navigational XPath but can treat any nested regular path query and it can be implemented efficiently without any dedicated techniques, by using any efficient datalog evaluator.

KEYWORDS: Graph databases, path queries, propositional dynamic logic, XPath, Datalog

1 Introduction

Regular path queries (Martens and Trautner 2018) are regular expressions for navigating in edge labeled graphs. They belong to the core of various query languages for graph databases and RDF triple stores. Nested regular path queries (NRPQs) (Libkin et al. 2013) extend on regular expressions by adding filters with logical operators, that in turn may contain regular path queries. They were first invented as the programs of propositional dynamic logic (PDL) (Fischer and Ladner 1979), when restricted to query datatrees, they constitute the navigational core of regular XPATH, and they are also part of *n*SPARQL for querying knowledge stores in the semantic Web (Pérez et al. 2010).

The set of nodes that can be reached by an NRPQ P on a graph G with a set of start nodes S can be computed in combined linear time, i.e. in $\mathcal{O}(|P||G|)$. This is folklore in the context of PDL, XPATH, and *n*SPARQL but was first shown for the richer alternation-free modal μ -calculus (Cleaveland and Steffen 1991). However, this complexity upper bound alone is far too high in practice: if the graph is a database then a complete traversal

$q_0(x) :- q_1(x), q_2(x).$
 $q_1(y) :- \text{start}(x), \text{edge}_a(x, y).$
 $q_2(x) :- \text{edge}_b(x, y), q_3(y).$
 $q_3(y) :- \text{edge}_c(y, z).$

Fig. 1. The Datalog program M_0 for the nested regular path query $P_0 = \text{edge}_a[\text{edge}_b/\text{edge}_c]$.

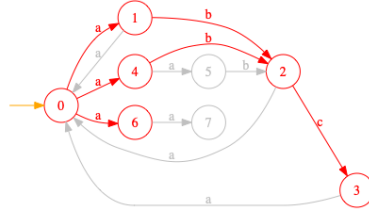


Fig. 2. Graph G_0 , start set $S_0 = \{0\}$, and the top-down needed subgraph for P_0 in red.

should be avoided for answering various queries. Which fraction of the graph is relevant for answering the query may depend on the query answering algorithm. Therefore, we formalize a notion of needed subgraph coined as *top-down needed subgraph*, as the subgraph that is traversed with a top-down evaluation of the query. We propose a query answering algorithm with combined linear complexity *with respect to the top-down needed subgraph*, instead of the whole graph which we consider as too expensive.

For regular path queries, a canonical notion of the top-down needed subgraph seems quite intuitive. It contains all nodes and edges that are traversed when considering the path query as a description for navigation while starting in the given set of start nodes. The presence of the Kleene star makes memoization mandatory for otherwise the algorithm may loop infinitely. The part of the graph that is traversed this way is what we call the top-down needed subgraph. The notion of top-down needed nodes can then be lifted from regular path queries to NRPQs rather naturally. What becomes more tedious is to find an evaluation algorithm for NRPQs that satisfies our complexity requirement. The existing proposals in (Pérez et al. 2010; Arenas and Pérez 2011; Gottlob et al. 2003) achieve combined linear time complexity by pre-evaluating the filters all over the graph in a bottom-up manner and then running an evaluation algorithm for regular path queries. Evaluating the filters top-down seems more difficult, since one would have to jump back to the starting node, requiring to compute a binary relation. However, the bottom-up pre-computation of the filters over all the graph may visit nodes that are *not* needed for top-down evaluation of the NRPQ so these algorithms do not satisfy the envisaged complexity bound. As an example, consider the graph G_0 in Fig. 2 with edge labels $\{a, b, c\}$, the NRPQ $P_0 = \text{edge}_a[\text{edge}_b/\text{edge}_c]$, and set of start nodes $S_0 = \{0\}$. Query P_0 started at S_0 selects all those nodes of G_0 that are connected to the start node 0 by an a -edge, and have a path over a b -edge followed by a c -edge. The top-down algorithm with pre-evaluation of filters will first compute the answer set of the filter $[\text{edge}_b/\text{edge}_c]$, which is $\{1, 4, 5\}$. It will then compute the set of nodes that are reached from the start node 0 over an a -edge, which is $\{1, 4, 6\}$. The answer set is the intersection, which is $\{1, 4\}$. This algorithm, however, will inspect some nodes and edges for the pre-evaluation of the filters that are *not* top-down needed, namely the node 5 and the b -edge from 5 to 2.

We will show that this complexity problem can be avoided by enhancing the naive top-down evaluator with memoization – instead of precomputing the filter queries. The right kind of memoization can be obtained by compiling the path query into a monadic datalog program, and then evaluating this datalog program in a top-down manner. Even though monadic, the datalog program may still use extensional predicates of higher arities. In the case of P_0 we obtain the datalog program in Fig. 1, which for talking about graphs

with edge labels in $\{a, b, c\}$ uses the binary extensional predicates $\text{edge}_a, \text{edge}_b, \text{edge}_c$. Furthermore, there is the monadic extensional predicate *start* for representing the start set. We note that a filter query such as $[\text{edge}_b/\text{edge}_c]$ is compiled quite differently (see the rules of the intensional predicates q_2 and q_3) to how one would compile the path query $\text{edge}_b/\text{edge}_c$. The reason is that a filter query returns the node where the path starts – under the condition that some node is reached at the end – while the path query selects all the nodes reached at the end.

Our first contribution is an algorithm that answers NRPQs in the combined linear time $\mathcal{O}(|\text{tdn}_{G,S}(P)||P|)$ with respect to the size of top-down needed subgraph $\text{tdn}_{G,S}(P)$. For this, we present a linear time compilation scheme for mapping path queries to datalog queries. For the sake of presentation, we treat only negation-free NRPQs, so that stratified negation is not needed. We prove that the compiler is correct in that if it transforms a query P and a start set S into a datalog query M , then top-down needed subgraph $\text{tdn}_{G,S}(P)$ is the part of the graph’s database that is visited by top-down evaluation of the datalog query M on the database. Furthermore, the datalog queries produced are monadic, and restricted in such a way that their top-down evaluation can be done in combined linear time depending on the size of the top-down visited subdatabase (Theorem 3 of Ullman (1989) and Tekle and Liu (2010) for an extensions to Datalog with stratified negation). It follows that the answer set of the NRPQ on the graph with start set S can indeed be computed in time $\mathcal{O}(|\text{tdn}_{G,S}(P)||P|)$.

Our algorithm can be extended to a jumping algorithm for answering NRPQs on graphs with indexes. The indexes are binary relations defined by other NRPQs that allow the algorithm to jump in the graph. For instance, when given an index for the NRPQ $I = \text{edge}^*/a?$ on the input graph, the evaluation algorithm can always jump to all a -labeled nodes accessible from the current node, without visiting the intermediates. We consider that the indexes are given with the input, since they are usually pre-computed elsewhere. Therefore, the indexes can simply be integrated into the graph as new edges that are labeled by the index’s name, which is I in our example. Furthermore, the NRPQ is then rewritten by substituting all occurrences of I as a subquery in the NRPQ by edge_I , so that we can apply the previous machinery. An efficient implementation of our algorithm can be based on any efficient top-down datalog evaluator, since it is sufficient to evaluate the monadic datalog program produced by our compiler.

Our graph jumping algorithm permits to reformulate without specialized techniques a very efficient automata-based algorithm proposed by Maneth and Nguyen (2010) that evaluates NRPQs on datatrees with indexes based on jumping. More precisely, their algorithm covers navigational forwards XPATH queries on XML documents. It is based on alternating tree automata with selection states (which can be seen as binary datalog programs while ours are monadic). Our generic approach overcomes the limitations of their algorithm: it is not bound to trees but applies to graphs; it is not limited to forward navigational XPath but can treat any NRPQs also with backward steps.

Outline. In Section 2, we recall the definition of NRPQs. In Section 3, we formally define their notion of top-down needed subgraphs. In Section 4, we recall preliminaries on datalog queries, while discussing the complexity of top-down evaluation in Section 5. In Section 6, we give our compiler from NRPQs to datalog queries with its complexity theorem. Proofs can be found in the appendix. Section 7 presents the jumping evaluation algorithm for NRPQs on graphs with indexes, and Section 8 an experimental evaluation.

<i>filters</i>	$F \in \mathcal{F}_\Sigma$	$::=$	$[P \mid \text{node} \mid \text{node}_a \mid F \wedge F' \mid F \vee F' \mid \neg F$
<i>paths</i>	$P \in \mathcal{P}_\Sigma$	$::=$	$F? \mid \text{edge}_a \mid \text{edge}_a^{-1} \mid P/P' \mid P \cup P' \mid P^+ \mid \text{goto}(F)$

Fig. 3. The syntax of NRPQs with labels $a \in \Sigma$.

2 Nested Regular Path Queries

Regular path queries on labeled graphs (Libkin et al. 2013) can be extended to NRPQs by adding filters with logical operators (Martens and Trautner 2018). CoreXPath (Gottlob et al. 2003) is a sublanguage of NRPQs with limited recursion where the interpretation is restricted to an unranked tree. NRPQs were known even much earlier as the propositional dynamic logic (PDL) of Fischer and Ladner (1979).

We start from a finite set of labels Σ . A (finite) Σ -labeled digraph is a tuple $G = (V, (V_a)_{a \in \Sigma}, (E_a)_{a \in \Sigma})$ where V is a finite set of nodes, $V_a \subseteq V$ a finite subset of a -labeled nodes, and $E_a \subseteq V \times V$ a finite set of a -labeled edges where $a \in \Sigma$. Note that nodes may have multiple labels or none, while each edge has a unique label. Between two nodes there may be multiple edges with different labels though. An example for a labeled graph G_0 with labels in $\Sigma = \{a, b, c\}$ was given graphically in Fig. 2. The set of nodes of the graph is $V = \{0, \dots, 7\}$. Here, the nodes are not labeled, so $V_a = V_b = V_c = \emptyset$. Each of the edge has a unique label. There are eight a -labeled edges in E_a , three b -labeled edges in E_b and one c -labeled edge in E_c .

The syntax of NRPQs with labels in Σ is presented in Fig. 3. It consists of a set of filters \mathcal{F}_Σ that select a set of graph nodes, and a set of paths \mathcal{P}_Σ that select a set of pairs of graph nodes. The filter *node* selects all the nodes, while the filter *node_a* selects all a -labeled nodes. The set of nodes nodes that are both a -labeled and b -labeled but not c -labeled is queried by filter $\text{node}_a \wedge \text{node}_b \wedge \neg \text{node}_c$. Path *edge_a* selects all a -labeled edges and path $\text{edge} =_{df} \cup_{a \in \Sigma} \text{edge}_a$ the set of all edges. The path *node?* selects the identify on nodes $\{(v, v) \mid v \in V\}$. Path composition P/P' , path union $P \cup P'$ are supported as well as repeated path composition P^+ . The Kleene star on paths can be defined by $P^* =_{df} P^+ \cup \text{node}?$. Backwards edges can be queried by edge_a^{-1} , so that general backwards path P^{-1} can be defined, where $(P_1/P_2)^{-1} = P_2^{-1}/P_1^{-1}$ and $F?^{-1} = F?$. Finally, the path *goto*(F) permits to jump to any node of the graph satisfying filter F . In particular, if there is a label *root* $\in \Sigma$ that distinguishes a set of roots, than path $\text{goto}(\text{node}_{\text{root}})/P$ first jumps to some root node before executing path P .

A little more complex example for an NRPQ with signature $\Sigma = \{a, b, c\}$ is the path query $P_2 = \text{node}_a?/(\text{edge}^+ / [\text{edge}_b / \text{edge}_c]?)^*$. The evaluation of P_2 on a given graph from a start node tests whether the start node is a -labeled, and if so, it navigates from there repeatedly, over a sequence of edges to some node for which there exists an outgoing path over edges with labels b and then c . The set of all nodes reached this way is selected.

The semantics of paths P on labeled digraphs G is the binary relation $\llbracket P \rrbracket_G \subseteq V \times V$ defined in Fig. 4 in mutual recursion with the semantics of filters $\llbracket F \rrbracket_G \subseteq V$. Despite its binary semantics, we will use paths for defining sets of nodes by fixing a start set S for the navigation. So let G be a labeled graph and S a subset of the nodes of G . For any $P \in \mathcal{P}_\Sigma$, the set $\llbracket P \rrbracket_G(S) = \{v \mid \exists v' \in S. (v', v) \in \llbracket P \rrbracket_G\}$ contains all nodes that can be reached when starting at some node of the start set S and navigating over the path P . Similarly, the set $\llbracket F \rrbracket_G(S) = \llbracket F \rrbracket_G \cap S$ contains all nodes from S that satisfy the filter F .

$$\begin{array}{ll}
 \llbracket [P] \rrbracket_G = \{v \mid \exists v'. (v, v') \in \llbracket P \rrbracket_G\} & \llbracket [F?] \rrbracket_G = \{(v, v) \mid v \in \llbracket F \rrbracket_G\} \\
 \llbracket [node] \rrbracket_G = V & \llbracket [edge_a] \rrbracket_G = E_a \\
 \llbracket [node_a] \rrbracket_G = V_a & \llbracket [edge_a^{-1}] \rrbracket_G = E_a^{-1} \\
 \llbracket [\neg F] \rrbracket_G = V \setminus \llbracket [F] \rrbracket_G & \llbracket [P/P'] \rrbracket_G = \llbracket [P] \rrbracket_G \circ \llbracket [P'] \rrbracket_G \\
 \llbracket [F \wedge F'] \rrbracket_G = \llbracket [F] \rrbracket_G \cap \llbracket [F'] \rrbracket_G & \llbracket [P^+] \rrbracket_G = \llbracket [P] \rrbracket_G^+ \\
 \llbracket [F \vee F'] \rrbracket_G = \llbracket [F] \rrbracket_G \cup \llbracket [F'] \rrbracket_G & \llbracket [P \cup P'] \rrbracket_G = \llbracket [P] \rrbracket_G \cup \llbracket [P'] \rrbracket_G \\
 & \llbracket [goto(F)] \rrbracket_G = \{(v, v') \mid v' \in \llbracket [F] \rrbracket_G\}
 \end{array}$$

 Fig. 4. Semantics of NRPQs on a Σ -labeled digraph $G = (V, (V_a)_{a \in \Sigma}, (E_a)_{a \in \Sigma})$.

$$\begin{array}{ll}
 tdn_{G,S}(node) = \{node(v) \mid v \in S\} & tdn_{G,S}([P]) = tdn_{G,S}(P) \\
 tdn_{G,S}(node_a) = \{node(v) \mid v \in S\} & tdn_{G,S}(F \wedge F') = tdn_{G,S}(F) \cup tdn_{G,\llbracket [F] \rrbracket_G(S)}(F') \\
 \cup \{node_a(v) \mid v \in V_a \cap S\} & tdn_{G,S}(F \vee F') = tdn_{G,S}(F) \cup tdn_{G,S}(F') \\
 tdn_{G,S}(F?) = tdn_{G,S}(F) & tdn_{G,S}(P/P') = tdn_{G,S}(P) \cup tdn_{G,\llbracket [P] \rrbracket_G(S)}(P') \\
 tdn_{G,S}(edge_a) = \{node(v) \mid v \in S\} & tdn_{G,S}(P^+) = tdn_{G,\llbracket [P^+] \rrbracket_G(S)}(P) \\
 \cup \{edge_a(v, v'), node(v') \mid v \in S, (v, v') \in E_a\} & tdn_{G,S}(P \cup P') = tdn_{G,S}(P) \cup tdn_{G,S}(P') \\
 tdn_{G,S}(edge_a^{-1}) = \{node(v) \mid v \in S\} & tdn_{G,S}(goto(F)) = tdn_G(F) \quad (\text{see Fig. A 1}) \\
 \cup \{edge_a(v', v), node(v) \mid v' \in S, (v, v') \in E_a\} &
 \end{array}$$

Fig. 5. Facts of top-down needed subgraphs for negation-free paths and filters.

3 Top-Down Needed Subgraphs

We are interested in the top-down evaluation of path queries, starting the navigation at the beginning of the path with a set of start nodes, and then moving along the path to other sets of nodes, reaching the end of the path.

We next define the subgraph that will be visited by such a traversal for a path query, and call it the *top-down needed subgraph*. For doing so we consider labeled graphs as extensional databases, i.e., as the sets of relational facts constructed from a relational signature and a set of constants. More concretely, we map any Σ -labeled graph $G = (V, (V_a)_{a \in \Sigma}, (E_a)_{a \in \Sigma})$ to the following set of database facts:

$$\begin{aligned}
 db(G) = \{ & node(v) \mid v \in V\} \cup \{node_a(v) \mid v \in V_a, a \in \Sigma\} \\
 & \cup \{edge_a(v, v') \mid (v, v') \in E_a, a \in \Sigma\}.
 \end{aligned}$$

The facts are build from the monadic predicates $node$ and $node_a$ and the binary predicates $edge_a$ for all $a \in \Sigma$, and the graph nodes $v \in V$ as constants. Conversely, consider a set of facts D with the following properties: 1. if $node_a(v) \in D$ then $node(v) \in D$ and 2. if $edge_a(v, v') \in D$ then $node(v) \in D$ and $node(v') \in D$. For any such set D there exists a unique graph G such that $db(G) = D$. We can therefore identify any graph G with the sets of facts $D = db(G)$.

For any Σ -labeled digraph G and set of start nodes S we define in Fig. 5 the set of *facts of top-down needed subgraph* $tdn_{G,S}(P)$ and $tdn_{G,S}(F)$ for negation-free paths P and filters F in mutual recursion. In the case of goto expressions, Fig. A 1 defines $tdn_{G,S}(goto(F)) = tdn_G(F)$ for restarting the computation with all nodes satisfying F . The natural algorithm for computing the answer set of filter $node_a$ at start set S will filter for all nodes $v \in S$ such that $v \in V_a$. Therefore all nodes in S need to be visited, as well as the a -label of all nodes in $V_a \cap S$. The extensional database of the top-down needed subgraph $tdn_{G,S}(a)$ therefore contains the facts in $\{node(v) \mid v \in S\}$ and $\{node_a(v) \mid v \in V_a \cap S\}$. The definition of $tdn_{G,S}(F \wedge F')$ is sequential from the

left to the right. When the filter query F is failing for a node v then there is no need to check the filter query F' so as to know that the filter query $F \wedge F'$ is not verified by v . In contrast, the definition of $tdn_{G,S}(F \vee F')$ is done a parallel manner, so that both subfilters need to be evaluated from the start nodes. The sequential alternative would lead to smaller top-down needed subgraphs, which might seem advantageous:

$$tdn_{G,S}^{seq}(F \vee F', S) = tdn_{G,S}(F) \cup tdn_{G, \llbracket \neg F \rrbracket_{\mathcal{C}}(S)}(F').$$

However, obtaining an evaluator with this sequential behavior by compilation to datalog would require us to use stratified negation, that we prefer to avoid for the sake of presentation. For the same reason, we restrict the definition of top-down needed subgraphs to negation-free path queries.

The definition of $tdn_{G,S}(P^+)$ is made of every attempt to construct a path of P starting from the nodes of S or the nodes that can be reached from S with a path of P^+ . In the case of goto expressions, we have defined $tdn_{G,S}(\text{goto}(F)) = tdn_G(F)$ for restarting the computation with all nodes satisfying F . We could set $tdn_G(F)$ to $tdn_{G,V}(F)$, but this would not be optimal since all nodes of V would be top-down needed even for most simple filter $F = \text{node}_a$. A better definition where only the nodes of V_a are top-down needed is given in Fig. A 1.

Example 1. Consider the query $P_0 = \text{edge}_a[\text{edge}_b/\text{edge}_c]$ on the graph G_0 with signature $\Sigma_0 = \{a, b, c\}$ in Fig. 2 with the start set $S_0 = \{0\}$. The set of top-down needed facts $tdn_{G_0, S_0}(P_0)$ is then $\{\text{edge}_a(0, 1), \text{edge}_a(0, 4), \text{edge}_a(0, 6), \text{edge}_b(1, 2), \text{edge}_b(4, 2), \text{edge}_c(2, 3)\}$. The top-down needed subgraph $\text{graph}(tdn_{G_0, \{0\}}(P_0))$ which is annotated in red in Fig. 2 is thus $(\{0, \dots, 6\}, (V_\ell)_{\ell \in \Sigma_0}, (E_\ell)_{\ell \in \Sigma_0})$ where $V_a = V_b = V_c = \emptyset$, $E_a = \{(0, 1), (0, 4), (0, 6)\}$, $E_b = \{(1, 2), (4, 2)\}$, and $E_c = \{(2, 3)\}$.

4 Datalog Queries

We recall preliminaries on the syntax and semantics of datalog programs without negation and how to use them to define datalog queries on extensional databases.

The syntax of datalog is parametrized by a finite set of predicates $p, q, r \in \mathcal{P}$ and a disjoint finite set of constants $a, b, c \in \mathcal{C}$. The set of predicates is partitioned into a subset of *extensional predicates* \mathcal{P}_{ext} and a disjoint subset of *intensional predicates* \mathcal{P}_{int} , so $\mathcal{P} = \mathcal{P}_{ext} \cup \mathcal{P}_{int}$. Constants will serve as database elements and extensional predicates for naming database relations. An (*extensional*) *database* is a subsets of ground literals of the form $p(a_1, \dots, a_n)$ where $p \in \mathcal{P}_{ext}$ has arity $n \geq 0$ and $a_1, \dots, a_n \in \mathcal{C}$.

We fix a set of variables $\mathcal{V} = \{x, y, z, \dots\}$ distinct from the constants and predicates. A *term* $u, s, t \in \mathcal{T}_{\mathcal{C}} = \mathcal{V} \uplus \mathcal{C}$ is either a variable or a constant. The set of (*positive*) *literals* \mathcal{L} is a subset of terms of the form $q(u_1, \dots, u_n)$ where $q \in \mathcal{P}$ has arity n and $u_1, \dots, u_n \in \mathcal{T}_{\mathcal{C}}$. A vector of terms is denoted by $\vec{t} \in \mathcal{T}_{\mathcal{C}}^*$. The set of all literals with extensional predicates is denoted by \mathcal{L}_{ext} and those with intensional predicates by \mathcal{L}_{int} . A *goal* is a vector of literals $\vec{\ell} \in \mathcal{L}^*$ that is to be understood as a conjunction. The set of free variables $fv(\vec{t}), fv(\vec{\ell}) \subseteq \mathcal{V}$ are defined as usual. Similarly for the sets of occurring constants $cst(\vec{t}), cst(\vec{\ell}) \subseteq \mathcal{C}$. A *clause* is a pair of the form $q(\vec{t}) :- \vec{\ell}$. where $q(\vec{t}) \in \mathcal{L}_{int}$ and $\vec{\ell} \in \mathcal{L}^*$. We call $q(\vec{t})$ the *head* and $\vec{\ell}$ the *body* of the clause. The clause $q(\vec{t}) :- \vec{\ell}$ is *safe* if $fv(\vec{t}) \subseteq fv(\vec{\ell})$. We only work with safe clauses throughout this paper.

$$\begin{aligned}
 \llbracket \epsilon \rrbracket_{M,D} &= \{\emptyset\} \\
 \llbracket \ell \rrbracket_{M,D} &= \{\Pi_{fv(\ell)}(\sigma \bowtie \sigma') \mid \sigma = \text{unif}(\ell, \ell'), \ell' :- \vec{\ell} \text{ in } \text{ren}(M), \sigma' \in \llbracket \sigma(\vec{\ell}) \rrbracket_{M,D}\} \text{ if } \ell \in \mathcal{L}_{int} \\
 \llbracket \ell \rrbracket_{M,D} &= \{\Pi_{fv(\ell)}(\sigma) \mid \sigma = \text{unif}(\ell, \ell'), \ell' \in D\} \text{ if } \ell \in \mathcal{L}_{ext} \\
 \llbracket \ell_1 \dots \ell_n \rrbracket_{M,D} &= \{\sigma' \bowtie \sigma \mid \sigma \in \llbracket \ell_1 \rrbracket_{M,D}, \sigma' \in \llbracket \sigma(\ell_2 \dots \ell_n) \rrbracket_{M,D}\}
 \end{aligned}$$

Fig. 6. The semantics of datalog query $?\vec{\ell}.M$ on database D where $n \geq 2$, $\ell, \ell_1, \dots, \ell_n \in \mathcal{L}$.

A (*safe*) *datalog program* is a finite subset M of safe clauses. A (*safe*) *datalog query* has the form $?\vec{\ell}.M$, where $\vec{\ell} \in \mathcal{L}^*$ is a datalog goal and M a safe datalog program M . We now turn our attention to the semantics of datalog queries. Given a datalog query $?\vec{\ell}.M$ and an extensional database D , we need to define the set of substitutions that answer the query. A *substitution* is a finite partial function σ from \mathcal{V} to $\mathcal{T}_{\mathcal{C}}$. We write \square for the empty substitution. Any substitution can be lifted to a total function on all variables by defining $\sigma(x) = x$ for all $x \notin \text{dom}(\sigma)$. We lift substitutions further to total functions $\sigma : \mathcal{T}_{\mathcal{C}}^* \rightarrow \mathcal{T}_{\mathcal{C}}^*$ such that for all $n \geq 0$, $t_1, \dots, t_n \in \mathcal{T}_{\mathcal{C}}$ and $a \in \mathcal{C}$:

$$\sigma(t_1 \dots t_n) = \sigma(t_1) \dots \sigma(t_n) \quad \text{and} \quad \sigma(a) = a$$

Similarly, substitutions are lifted to total functions $\sigma : \mathcal{L}^* \rightarrow \mathcal{L}^*$ such that for all $\vec{t} \in \mathcal{T}_{\mathcal{C}}^*$, $n \geq 0$, and $\ell_1, \dots, \ell_n \in \mathcal{L}$:

$$\sigma(q(\vec{t})) = q(\sigma(\vec{t})), \quad \text{and} \quad \sigma(\ell_1 \dots \ell_n) = \sigma(\ell_1) \dots \sigma(\ell_n)$$

The *renaming closure* of a program is the set of all clauses that can be obtained from the clauses of the program by renaming variables bijectively:

$$\text{ren}(M) = \{\sigma(\ell) :- \sigma(\vec{\ell}) \mid \ell :- \vec{\ell} \text{ in } M, \sigma \text{ is one-to-one substitution, } \text{ran}(\sigma) \subseteq \mathcal{V}\}$$

We define joins and projections on substitutions as for the relational algebra: for any two substitutions σ and σ' and any finite subset of variables $V \subseteq \mathcal{V}$:

$$\sigma \bowtie \sigma' = \begin{cases} \sigma \cup \sigma' & \text{if } \sigma \cup \sigma' \text{ is functional} \\ \text{undefined} & \text{otherwise} \end{cases} \quad \Pi_V(\sigma) = \sigma|_V$$

For any two literals ℓ, ℓ' we define $\text{unif}(\ell, \ell')$ as the most general unifier σ such that $\sigma(\ell) = \sigma(\ell')$ if it exists, and leave it undefined otherwise.

We define the semantics $\llbracket \vec{\ell} \rrbracket_{M,D}$ of a datalog query $?\vec{\ell}.M$ on an extensional database D as the least fixpoint that satisfies the equations in Fig. 6. Notice that whenever we use the operation $\sigma \bowtie \sigma'$ then we have $\text{dom}(\sigma) \cap \text{dom}(\sigma') = \emptyset$, so that $\sigma \bowtie \sigma' = \sigma \cup \sigma'$ is a well-defined substitution. Each query answer $\sigma \in \llbracket \vec{\ell} \rrbracket_{M,D}$ has domain $fv(\vec{\ell})$ and always maps to constants since we work with safe datalog programs, so $\sigma : fv(\vec{\ell}) \rightarrow \mathcal{C}$. The semantics that we have given mimics the top-down datalog evaluation, which starts with the goal in the query and generates subgoals by unfolding the clauses of the datalog program, while instantiating the variables, until it reaches some ground facts from the extensional database. In general, this process may enter into infinite loops if not controlled by memoization. The whole top-down evaluation can always be represented as a join tree as we illustrate by example in Fig. 7. In the case of infinite loops, the join tree is infinite.

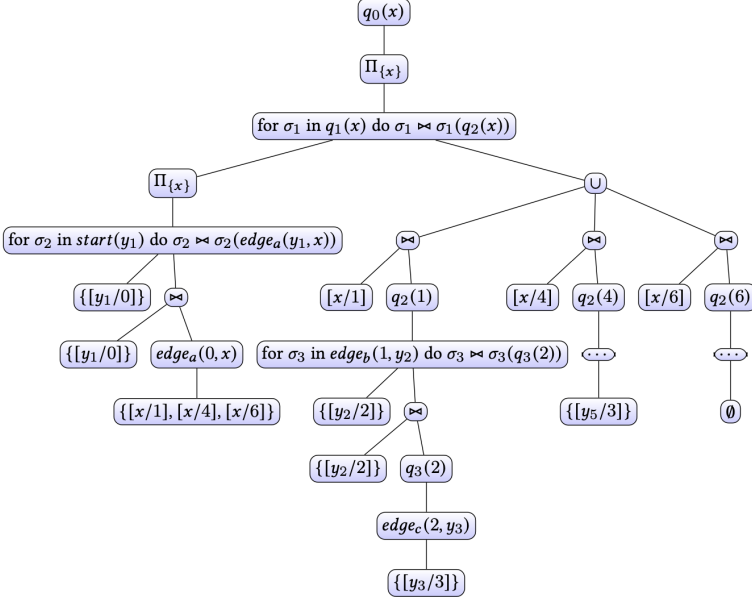


Fig. 7. Top-down evaluation of $\llbracket q_0(x) \rrbracket_{M_0, db(G_0) \cup \{start(0)\}} = \{[x/1], [x/4]\}$ where M_0 is the datalog program from Fig. 1 for $P_0 = edge_a[edge_b/edge_c]$, and G_0 the graph from Fig. 2.

$$tdv_{M,D}(\epsilon) = \emptyset$$

$$tdv_{M,D}(\ell) = \{node(a) \mid a \in cst(\ell)\} \cup \bigcup \{\ell' \mid unif(\ell, \ell') \text{ defined, } \ell' \text{ in } D\} \text{ if } \ell \in \mathcal{L}_{ext}$$

$$tdv_{M,D}(\ell) = \{node(a) \mid a \in cst(\ell)\} \cup \bigcup \{tdv_{M,D}(\sigma(\vec{\ell})) \mid \sigma = unif(\ell, \ell'), \ell' :- \vec{\ell} \text{ in } ren(M)\} \text{ if } \ell \in \mathcal{L}_{int}$$

$$tdv_{M,D}(\ell_1 \dots \ell_n) = tdv_{M,D}(\ell_1) \cup \bigcup \{tdv_{M,D}(\sigma(\ell_2 \dots \ell_n)) \mid \sigma \in \llbracket \ell_1 \rrbracket_{M,D}\}$$

Fig. 8. The top-down visited sub-database $tdv_{M,D}(\vec{\ell})$ where $\ell, \ell_1, \dots, \ell_n \in \mathcal{L}$ and $n \geq 2$.

5 Complexity of Top-Down Evaluation of Datalog Queries

Known results on the complexity of top-down datalog evaluation give us the formal tools to prove for particular datalog queries, that the complexity of the top-down evaluation is in combined linear time but with respect to the top-down visited sub-database, rather than with respect to the full database.

For any datalog query $?-\vec{\ell}$. M and extensional database D we next define the part of D that is visited by the top-down evaluation of the datalog query. For this we assume that the set of extensional predicates of D contains a monadic predicate $node \in \mathcal{P}_{ext}$ such that $node^D = \mathcal{C}$. We define the *top-down visited sub-database* $tdv_{M,D}(\vec{\ell})$ as the extensional database over \mathcal{P}_{ext} – following the semantics of datalog queries – as the least fixed point of equations in Fig. 8.

Definition 1. We call a datalog goal $\vec{\ell}$ simply combined linear (SCL) if $fv(\vec{\ell})$ is a singleton or empty, or if all its variables are guarded by a single extensional literal of $\vec{\ell}$. We call a datalog query $?-\vec{\ell}$. M SCL if the datalog goal $\vec{\ell}$ is SCL and for each of the clauses $\ell :- \vec{\ell}$ in datalog program M , the datalog goal $\vec{\ell}$ is SCL.

For example, let $p, q \in \mathcal{P}_{int}$ be monadic and $r \in \mathcal{P}_{ext}$ be binary. The goal $p(x), r(x, y), q(y)$ is then SCL, since both of its variables x and y are guarded by the extensional literal

$r(x, y)$. The goal $p(x), r(x, x), q(y)$ on the contrary is not SCL, as it contains two variables of which y is not guarded by a single extensional literal. The goal $p(x), q(x)$ is SCL since it contains no more than a single free variable.

Given an extensional database D , any SCL goal $\vec{\ell}$ has a number of ground instances that is linear in the size of D . Even better the number of ground instances inspected by top-down evaluation of the datalog query $?-\vec{\ell}$. M is linear in the size of the top-down visited database $tdv_{M,D}(\vec{\ell})$. In the case where $fv(\vec{\ell})$ contains at most one variable, this variable must be instantiated by some node of the top-down visited sub-database. Otherwise, the set of free variables $fv(\vec{\ell})$ is guarded by a single extensional literal of $\vec{\ell}$, say $p(\vec{t}) \in \mathcal{L}_{ext}$. In this case, any ground instance of $\vec{\ell}$ visited by top-down evaluation of M is determined by $unif(p(\vec{t}), p(\vec{v}))$ for some fact $p(\vec{v}) \in tdv_{M,D}(\vec{\ell})$.

Theorem 2 (Theorem 3 of Ullman (1989)). *The answer set $\llbracket \vec{\ell} \rrbracket_{M,D}$ of a safe SCL query $?-\vec{\ell}$. M on an extensional database D can be computed in time $\mathcal{O}(|M| |tdv_{M,D}(\vec{\ell})|)$.*

For safe SCL datalog queries, the query answering by top-down evaluation with memoization is thus indeed in combined linear time with respect to the size of the top-down visited sub-database. This can be extended to stratified datalog (Tekle and Liu 2010).

6 Compiler to SCL Datalog Queries

We now contribute the compiler from negation-free path queries P and start set S to SCL datalog queries $?-\vec{\ell}$. M , such that for any graph G with nodes subsuming S , the extensional database of the top-down needed subgraph $tdn_{G,S}(P)$ is equal to the top-down visited sub-database $tdv_{M,db(G)}(\vec{\ell})$. The top-down evaluation of the datalog query $?-\vec{\ell}$. M on the graph's database $db(G)$ thus yields the expected upper complexity bound for the evaluation of path queries by Theorem 2.

For any set of start nodes S and monadic predicate $i \in \mathcal{P}_{int}$, we define a datalog program $Start^i(S) = \{i(v) :- . \mid v \in S\}$. The compilation scheme for path queries follows the structure of paths and filters by mutual recursion. It is given by the datalog programs $Acc^{i,f}(P)$ in Fig. 9, $Filt^c(F)$ in Fig. 10 and $Ex^{c,r}(P)$ in Fig. 11. Path queries outside filters need to compute all accessible nodes by $Acc^{i,f}(P)$, while path queries within filters need to check the existence of accessible nodes by $Ex^{c,r}(P)$. The compiler introduces fresh monadic predicates for all subexpressions: *initial predicates* $i, i', i'' \in \mathcal{P}_{int}$, *final predicates* $f, f', f'' \in \mathcal{P}_{int}$ *final*, *checks* $c, c', c'' \in \mathcal{P}_{int}$, and *continuations* $r, r', r'' \in \mathcal{P}_{int}$.

Given a graph G and with a start set $S \subseteq V$ of graph nodes, the answer set of the datalog query $?-f(x)$. $Acc^{i,f}(P) \cup Start^i(S)$ on the extensional database $db(G)$ is $\{[x/v] \mid v \in \llbracket P \rrbracket_G(S)\}$, assigning the free variable x to some node v reachable from S over P in G . The initial predicate i captures the set of start nodes, and the *final* predicate f the answer set of the path query P started from there. The fresh monadic predicates make the datalog programs for the subexpressions able to communicate. For instance, we have $Acc^{i,f}(P'/P'') = Acc^{i,f'}(P') \cup Acc^{f',f}(P'')$. Here the final predicate $f' \in \mathcal{P}_{int}$ represents the answer set of path P' started at node set i , but also the start set for the path P'' . This is since the start nodes of P'' in the query P'/P'' are the nodes that are reached with the query P' . For the recursive path queries P^+ we have $Acc^{i,f}(P^+) = Acc^{i,f}(P) \cup \{i(x) :- f(x)\}$. Here the rule $i(x) :- f(x)$. represents the fact

$$\begin{array}{ll}
Acc^{i,f}(\text{edge}_a) = \{f(x) :- i(y), \text{edge}_a(y, x).\} & Acc^{i,f}(P' \cup P'') = Acc^{i,f}(P') \cup Acc^{i,f}(P'') \\
Acc^{i,f}(\text{edge}_a^{-1}) = \{f(x) :- i(y), \text{edge}_a(x, y).\} & Acc^{i,f}(\text{goto}(F')) = Filt^{f'}(F') \cup \\
Acc^{i,f}(P'/P'') = Acc^{i,f'}(P') \cup Acc^{f',f}(P'') & \{f(x) :- j(), f'(x). \quad j() :- i(x).\} \\
Acc^{i,f}(P^+) = Acc^{i,f}(P) \cup \{i(x) :- f(x).\} & Acc^{i,f}(F'?) = Filt^{f'}(F') \cup \{f(x) :- i(x), f'(x).\}
\end{array}$$

Fig. 9. The datalog program $Acc^{i,f}(P)$ for path P and monadic predicates $i, f \in \mathcal{P}_{int}$.

$$\begin{array}{ll}
Filt^c(a) = \{c(x) :- \text{node}_a(x).\} & Filt^c(F' \wedge F'') = Filt^{c'}(F') \cup Filt^{c''}(F'') \cup \\
Filt^c(\text{node}) = \{c(x) :- \text{node}(x).\} & \{c(x) :- c'(x), c''(x).\} \\
Filt^c(F' \vee F'') = Filt^{c'}(F') \cup Filt^{c''}(F'') \cup & Filt^c([P]) = Ex^{c,r}(P) \cup \{r(x) :- \text{node}(x).\} \\
\{c(x) :- c'(x). \quad c(x) :- c''(x).\} &
\end{array}$$

Fig. 10. The datalog program $Filt^c(F)$ for filter F and monadic predicate $c \in \mathcal{P}_{int}$.

$$\begin{array}{ll}
Ex^{c,r}(\text{edge}_a) = \{c(x) :- \text{edge}_a(x, y), r(y).\} & Ex^{c,r}(P' \cup P'') = Ex^{c,r}(P') \cup Ex^{c,r}(P'') \\
Ex^{c,r}(\text{edge}_a^{-1}) = \{c(x) :- \text{edge}_a(y, x), r(y).\} & Ex^{c,r}(\text{goto}(F')) = Filt^{c'}(F') \cup \\
Ex^{c,r}(P'/P'') = Ex^{c,f}(P') \cup Ex^{f,r}(P'') & \{c(x) :- j(). \quad j() :- c'(y), r(y).\} \\
Ex^{c,r}(P^+) = Ex^{c,r}(P) \cup \{r(x) :- c(x).\} & Ex^{c,r}(F'?) = Filt^{c'}(F') \cup \{c(x) :- c'(x), r(x).\}
\end{array}$$

Fig. 11. The datalog program $Ex^{c,r}(P)$ for path P with monadic predicates $c, r \in \mathcal{P}_{int}$.

that once a node is reached by the query P^+ it becomes a possible start node for the same query.

We next consider the datalog programs $Filt^c(F)$ defined in Fig. 10. For any graph G the answer set of the datalog query $?-c(x)$. $Filt^c(F)$ on the extensional database $db(G)$ is $\{[x/v] \mid v \in \llbracket F \rrbracket_G\}$, so that the free variables x may be bound to any node selected by the filter. Hence, for any start set S , the answer set of $?-i(x), c(x)$. $Filt^c(F) \cup Start^i(S)$ is $\{[x/v] \mid v \in \llbracket F \rrbracket_G(S)\}$. The filter for all nodes is compiled to $Filt^c(\text{node}) = \{c(x) :- \text{node}(x)\}$. Thereby, the check c is called for all nodes of the graph. Note that node is an extensional predicate, so this clause is safe. A conjunction of filters $Filt^c(F' \wedge F'')$ is compiled by adding the clause $c(x) :- c'(x), c''(x)$ to the datalog programs $Filt^{c'}(F')$ and $Filt^{c''}(F'')$. The added clause checks sequentially, whether a node x is filtered by F' and if so whether it is also filtered by F'' . A disjunction of filters $Filt^c(F' \vee F'')$ is compiled by adding the two clause $c(x) :- c'(x)$. and $c(x) :- c''(x)$. to the datalog programs $Filt^{c'}(F')$ and $Filt^{c''}(F'')$. The two added clauses check in parallel whether a node x is filtered by F' or whether x is filtered by F'' .

In Fig. 11 we define the datalog programs $Ex^{c,r}(P)$ for evaluating paths P existentially as needed when paths are used in filters, that is $Filt^c([P]) = Ex^{c,r}(P) \cup \{r(x) :- \text{node}(x)\}$. The check predicate c denotes the set of source nodes, from which some target node can be reached over P , while r is the continuation to which the target node must belong. Given a graph G and a start set S , the answer set of the datalog query $?-c(x)$. $Ex^{c,r}(P) \cup \{r(x) :- \text{node}(x)\}$ on the extensional database $db(G)$ is $\{[c/v] \mid (v, v') \in \llbracket P \rrbracket_G\}$. The continuation predicate r is required to allow us to compile path concatenations in filters, i.e., in $Ex^{c,r}(P'/P'') = Ex^{c,f}(P') \cup Ex^{f,r}(P'')$. Note that the interplay of the predicate c and r is similar to the one between i and f in $Acc^{i,f}(P)$.

Lemma 3. *For any path P , filter F , graph G , start set S , and monadic predicates $i, f, c, r \in \mathcal{P}_{int}$, the programs $Start^i(S)$, $Acc^{i,f}(P)$, $Filt^c(F)$, $Ex^{c,r}(P)$ are safe and SLC.*

$$\begin{aligned}
 reach_{M,r}(\epsilon) &= \emptyset \\
 reach_{M,r}(r(v), \vec{\ell}_1) &= \{v\} \cup reach_{M,r}(\sigma(\vec{\ell}_2, \vec{\ell}_1)) \mid \sigma = unif(r(v), \ell'), \ell' :- \vec{\ell}_2. \text{ in } ren(M)\} \\
 reach_{M,r}(\ell, \vec{\ell}_1) &= reach_{M,r}(\sigma(\vec{\ell}_2, \vec{\ell}_1)) \mid \sigma = unif(\ell, \ell'), \ell' :- \vec{\ell}_2. \text{ in } ren(M)\} \quad \text{if } \ell \neq r(v)
 \end{aligned}$$

Fig. 12. The $reach_{M,r}$ function for the datalog query $?-\vec{\ell}. M$.

The function $reach_{M,r}(\vec{\ell})$ defined in Fig. 12 returns the set of all nodes v , such that $r(v)$ is queried in the process of the top-down evaluation of the datalog query $?-\vec{\ell}. M$. Now, we provide the Propositions 4 and 5 for dividing the correctness proof into two parts. First — about subpaths and subfilters of some filter. Concomittantly with Theorem 2 they will imply the main efficiency Theorem 6.

Proposition 4. *For any filter query $F \in \mathcal{F}_\Sigma$, path query $P \in \mathcal{P}_\Sigma$, label $a \in \Sigma$, labeled graph G , subset $S \subseteq V$ of nodes of G , distinct monadic predicates $i, c, r \in \mathcal{P}_{int}$ and $x \in \mathcal{V}$.*

1. *if $M = Fill^c(F) \cup Start^i(S)$ and $\vec{\ell} = i(x), c(x)$ then:*
2. *if $M = Ex^{c:r}(P) \cup Start^i(S) \cup \{r(x) :- node(x).\}$ and $\vec{\ell} = i(x), c(x)$ then:*

$$\begin{aligned}
 \bullet \llbracket \vec{\ell} \rrbracket_{M,db(G)} &= \{[x/v] \mid v \in \llbracket F \rrbracket_G(S)\} & \bullet \llbracket \vec{\ell} \rrbracket_{M,db(G)} &= \{[x/v] \mid v \in S, \llbracket P \rrbracket_G(\{v\}) \neq \emptyset\} \\
 \bullet tdv_{M,db(G)}(\vec{\ell}) &= tdn_{G,S}(F) & \bullet tdv_{M,db(G)}(\vec{\ell}) &= tdn_{G,S}(P) \\
 \bullet tdv_{M,db(G)}(\vec{\ell}) &= tdn_{G,S}(F) & \bullet reached_M(\vec{\ell}) &= \llbracket P \rrbracket_G(S)
 \end{aligned}$$

Proposition 5. *For any path query $P \in \mathcal{P}_\Sigma$, labeled graph G , subset S of nodes of G , distinct intensional predicates $i, f \in \mathcal{P}_{int}$ and $x \in \mathcal{V}$, if $M = Acc^{i:f}(P) \cup Start^i(S)$ then:*

- $\llbracket f(x) \rrbracket_{M,db(G)} = \{[x/v] \mid v \in \llbracket P \rrbracket_G(S)\}$
- $tdv_{M,db(G)}(f(x)) = tdn_{G,S}(P)$

Theorem 6. *For any graph G with subset of nodes S and any path query $P \in \mathcal{P}_\Sigma$ the answer set $\llbracket P \rrbracket_G(S)$ can be computed in time $\mathcal{O}(|P| \cdot tdn_{G,S}(P))$.*

7 Jumping in Graphs

Preprocessing is mandatory for sharing efforts when evaluating multiple queries on the same large graph. Most typically, one can pre-compute indexes that give efficient access to some particular relations of the graph. Here we consider indexes, which are binary relations defined by NRPQs themselves.

For instance, we might want to from any node of the graph to the next a -labeled node in some fixed total order. In this case, one would like to have a jumping algorithm that visits only the top-down needed subgraph, but taken with respect to the graph, that is enriched with extra edges labeled by the names of the indexes.

Let us next consider a little more complex example. For this we suppose that we have an index for the NRPQ $acc_a = edge^*/a?$. We can then extend the signature Σ with a new label acc_a , the graph G with acc_a -labeled edges for all pairs in $\llbracket acc_a \rrbracket_G$, and rewrite the target path query by substituting all its subqueries acc_a by $edge_{acc_a}$. This has the advantage that fewer nodes are top-down needed after the rewriting on the enriched graph. For instance, a top-down evaluator for the path query acc_a without jumping needed to inspect *all* nodes of the graph accessible from S , since all of them needed to be tested for whether they satisfied the filter query a . After the rewriting to $edge_{acc_a}$, a top-down algorithm can jump directly from the start nodes in S to the accessible a -labeled nodes by using the index, so only accessible a -labeled nodes will be visited.

The general jumping algorithm starts with a set of indexes for NRPQs say for P_1, \dots, P_n . For answering a query P on a graph G with these indexes the jumping algorithm enriches the signature Σ by new labels P_1, \dots, P_n , the original graph G with new labeled edges $E_{P_j} = \llbracket P_j \rrbracket_G$ where $1 \leq j \leq n$, and then substitutes in the target query P all occurrences of the subqueries P_j by edge_{P_j} . The order of the substitution can be chosen arbitrarily, depending on the intended jumping strategy. In this way, the top-down needed subgraph of the enriched graph for the rewritten query is intuitively exactly the subgraph of the original graph that a top-down evaluation algorithm with jumping needs to visit.

This jumping algorithm can be used to reformulate in simple terms a variant of the efficient automata-based algorithm from Maneth and Nguyen (2010) that evaluates navigational path queries. More precisely, their algorithm covers navigational forward XPATH queries on XML documents, and is based on alternating tree automata with selection states (which can be seen as binary datalog programs, while ours are monadic). XML documents can be seen as labeled graphs, with two edge labels *firstchild* and *nextsibling*. Their algorithm can be based on the indexes for jumping to the a -labeled children, that is $\text{edge}/a?$, and for jumping to the top-most a -labeled descendants, i.e., $\text{top}_a = (\text{edge}/-a?)^*/\text{edge}/a?$. An XPATH query such as `descendant::a` can be rewritten as the NRPQ $(\text{top}_a)^+$. The evaluation of the query $(\text{top}_a)^+$ can then take advantage of the index $\text{edge}_{\text{top}_a}$. The main difference between both approaches is that ours doesn't try to produce the answer set in document order, while theirs does so. Therefore, binary indexes are sufficient for our purpose, while they need to use a ternary index (for relating following a -labeled nodes x of y below z). Moreover, our algorithm traverses the same part of the XML document as theirs and will thus be as efficient while being much simpler in terms of presentation. Our general graph approach overcomes the main limitations of Maneth and Nguyen's: it is not bound to trees and is not limited to forward navigational XPath but can treat any NRPQs also with backward steps.

8 Experiments

We implemented in OCaml our compiler from PDL to Datalog and also a compiler from navigational XPath queries to PDL queries on the graphs of XML documents. The edges of these graphs are labeled by *element*, *document* and string, and the edges by *first*, *next*, *name*, and attribute names.

We selected in Fig. 13 two typical benchmark XPATH queries from Maneth and Nguyen (2010) that can be applied to the scalable XML-documents from the XPath-Mark benchmark: query $Q01$ composes two child axis, and query $Q05$ two descendant axis `//listitem//keyword`. The translations for these XPath queries to the PDL queries $\text{pdl}.Q01$ and $\text{pdl}.Q05$ that can also be found there.

Query $\text{pdl}.Q01$ is easier for top-down evaluation, since it does not contain vertically recursive axis, so that its top-down needed subgraph remains small on the benchmark documents. $Q05$ is more difficult since using descendant axis, so that the top-down needed subgraph $\text{pdl}.Q05$ is the whole graph if not using indexes. So we also computed the indexes $\text{top}_{\text{keyword}}$ and $\text{top}_{\text{listitem}}$ for the descendant axis of $Q05$ and added them as extra edges to the graphs. Furthermore, the optimized query $\text{pdl}.Q05.\text{index}$ obtained from $\text{pdl}.Q05$ by using the index edges is given in Fig. 13 too.

Q01	/site/regions	Q05	//listitem//keyword
pdl.	$node_{document}?(edge_{first}/(edge_{next})^*/node_{element}?[edge_{name}/node_{site}]/edge_{first}/(edge_{next})^*/node_{element}?[edge_{name}/node_{regions}]$	pdl.	$node_{document}?(edge_{first}/(edge_{next})^*)^+ / node_{element}?[edge_{name}/node_{listitem}]/(edge_{first}/(edge_{next})^*)^+ / node_{element}?[edge_{name}/node_{keywords}]$
Q01		Q05	
		pdl.	$node_{document}?(edge_{toplistitem})^+ / (edge_{topkeyword})^+$
		Q05	
		index	

Fig. 13. Two benchmark XPATH queries from Maneth and Nguyen (2010), their translation to PDL, and the indexed PDL queries.

27KB	Q01	Q05	Q05.index	100MB	Q01	Q05	Q05.index
Saxon	0.000206	indexing	0.000315	Saxon		indexing	0.0016
XSB	0.001	0.006	0.001	XSB		35.857	5.029
SWI	0.004	0.189	0.004	SWI	-	-	-
LogicBlox	0.0045	0.0054	0.0045	LogicBlox	0.0124	-	0.0974

Fig. 14. Time in seconds for querying the 27KB XML-document with indexes.

Fig. 15. Time in seconds for querying the 100MB XML-document with indexes.

As gold standard, we compare with is the Saxon evaluator for XPath queries via XSLT. In order to measure the time we run the same XPath query 100 times in the same XSLT program with Saxon 10.5, subtract the time needed to load and index the XML document and divide by 100. It turns out, that Saxon has the best performance in all our tests, confirming our conjecture that it performs jumping evaluation with indexing for descendant axis.

We then implemented and tested our jumping algorithm based on existing top-down Datalog evaluators. We started with OCaml’s Datalog 0.6, but had to notice that the top-down evaluator did not always produce the correct results. We then experimented with the Prolog engines XSB 4.0 and SWI 8.4.1.1. On a small XML-document of 27KB, both engines perform decently, even though not as quick as Saxon. On Q01, they are one order of magnitude slower. The same holds for Q05 but only when using indexing.

We then considered a much bigger XML-document of 100MB. With this size we had to give up with SWI. XSB in contrast could read the graph of the XML document, but needed more than 30 minutes. Once done it could answer query *pdl.Q05* without indexing in 35 seconds. With indexing the time for answering *pdl.Q05.index* went down to 5 seconds. Saxon, in contrast, can load the graph in 15 seconds and answer query *Q05* in 1.6 milliseconds. So for answering the query *Q05*, Saxon showed 4 orders of magnitude more efficient than XSB.

We finally investigated the LogicBlox system (Aref et al. 2015), a more recent deductive database system which implements the language LogiQL extending on Datalog. With version 4.38 of LogicBlox we could read the graph of 100 MB in 19 seconds (rather than in more than 30 minutes as with XSB). LogiQL is a typed language implying some minor syntactic differences to standard datalog. Finally, LogicBlox has a transaction

level, that permits to interact with graphs dynamically, so that it can be queried many times without being reloaded. The earlier versions of LogicBlox supported bottom-up evaluation only. But since recently, top-down evaluation can be chosen by adding On-Demand annotations for all extensional predicates. When doing so, we could answer the query *pdl.Q05.index* in 97.4 milliseconds on the 100MB document. This is 2 orders of magnitudes better than with XSB! Nevertheless it is still by a factor of 75 slower than with Saxon. Figure 4 of Maneth and Nguyen (2010) reports 65 milliseconds for Q05 with optimal jumping, but on a slightly larger 116MB document. So the question is how the efficiency of our implementation could be increased further: with better indexes, early completion during Datalog evaluation, or by using special features of XPATH queries?

Conclusion and Future Work. The definition of the top-down needed subgraph allows us to prove that our algorithm for answering negation-free NRPQs visits only the interesting part of the graph. We believe that the restriction to negation-freeness can be relieved by compiling to stratified datalog. The new notion of top-down needed subgraphs may also allow the design of algorithms that transform NRPQs into equivalent ones that have a smaller top-down needed subgraph, for instance by inverting the path, or starting with some filter. It thus sets the stage for query optimization. In particular, the *goto* instructions permit the algorithm to jump directly to nodes with *rare* properties in the graph first and then compute the queries more efficiently. Another line of improvement would be to stop the evaluation of filters when it has been proven correct.

References

- AREF, M., TEN CATE, B., GREEN, T. J., KIMELFELD, B., OLTEANU, D., PASALIC, E., VELD-HUIZEN, T. L., AND WASHBURN, G. Design and implementation of the logicblox system. In *SIGMOD International Conference on Management of Data* 2015, pp. 1371–1382. ACM.
- ARENAS, M. AND PÉREZ, J. Querying semantic web data with sparql. In *30th ACM Symposium on Principles of Database Systems* 2011, pp. 305–316.
- CLEAVELAND, R. AND STEFFEN, B. A linear-time model-checking algorithm for the alternation-free modal mu-calculus. In *Proceedings of the 3rd International Workshop on Computer Aided Verification* 1991, CAV '91, 48–58, Springer.
- FISCHER, M. J. AND LADNER, R. E. 1979. Propositional dynamic logic of regular programs. *J. Comput. Syst. Sci.*, 18, 2, 194–211.
- GOTTLOB, G., KOCH, C., AND PICHLER, R. The complexity of XPath query evaluation. In *22nd ACM Symposium on Principles of Database Systems* 2003, pp. 179–190.
- LIBKIN, L., MARTENS, W., AND VRGOVC, D. Querying graph databases with XPath. In *16th International Conference on Database Theory* 2013, ICDT '13, 129–140. ACM.
- MANETH, S. AND NGUYEN, K. 2010. Xpath whole query optimization. *PVLDB*, 3, 1, 882–893.
- MARTENS, W. AND TRAUTNER, T. Evaluation and Enumeration Problems for Regular Path Queries. In *21st International Conference on Database Theory (ICDT)* 2018, volume 98 of *LIPICs*, pp. 19:1–19:21.
- PÉREZ, J., ARENAS, M., AND GUTIÉRREZ, C. 2010. nSPARQL: A navigational language for RDF. *J. Web Semant.*, 8, 4, 255–270.
- TEKLE, K. AND LIU, Y. Precise complexity analysis for efficient datalog queries. In *PPDP'10 - Symposium on Principles and Practice of Declarative Programming* 2010, pp. 35–44.
- ULLMAN, J. D. Bottom-up beats top-down for datalog. In *8th ACM Symposium on Principles of Database Systems (PODS)*, 1989, 140–149. ACM.

$$\begin{array}{ll}
 tdn_G(\text{node}) = \{\text{node}(v) \mid v \in V\} & tdn_G([P]) = tdn_G(P) \\
 tdn_G(\text{node}_a) = \{\text{node}_a(v) \mid v \in V_a\} & tdn_G(F \wedge F') = tdn_G(F) \cup tdn_{G, \llbracket F \rrbracket_G}(F') \\
 tdn_G(F?) = tdn_G(F) & tdn_G(F \vee F') = tdn_G(F) \cup tdn_G(F') \\
 tdn_G(\text{edge}_a) = \{\text{edge}_a(v, v') \mid (v, v') \in E_a\} & tdn_G(P/P') = tdn_G(P) \cup tdn_{G, \llbracket P \rrbracket_G(V)}(P') \\
 tdn_G(\text{edge}_a^{-1}) = \{\text{edge}_a(v', v) \mid (v, v') \in E_a\} & tdn_G(P^+) = tdn_{G, \llbracket P^+ \rrbracket_G(V)}(P) \\
 & tdn_G(P \cup P') = tdn_G(P) \cup tdn_G(P') \\
 & tdn_G(\text{goto}(F)) = tdn_G(F)
 \end{array}$$

Fig. A 1. Top-down needed subgraphs without start sets as needed for goto expressions.

Appendix A Proofs for Section 3 (Top-Down Needed Subgraphs)

Appendix B Proofs for Section 6 (Compiler to SCL Datalog Queries)

Our compiler could be extended to negated filters, but this would require to compile to stratified datalog. Stratified negation would also allow to compile disjunctions sequentially, so that only a smaller subgraph would get visited.

The accessibility of nodes over goto-paths could be expressed by

$$Acc^{i,f}(\text{goto}(F')) = Filt^{f'}(F') \cup \{f(x) :- i(y), f'(x).\}$$

where $Filt^{f'}(F')$ is the datalog program that represents the computation of the filter query F' and where f' is the predicate which captures the answer set. The clause $f(x) :- i(y), f'(x)$. means that a node x satisfying the filter F' is in the answer set if there is some node y in the start set S . So if $S \neq \emptyset$ then the datalog query $?-f'(x)$. $Filt^c(F') \cup Start^i(S)$ can be reduced to the datalog query $?-c(x)$. $Filt^c(F')$, which on the extensional database $db(G)$ has the answer set $\{[c/v] \mid v \in \llbracket F' \rrbracket_G\}$. We note that the top-down evaluation of the latter datalog query avoids visiting the nodes of the graph that are not top-down needed for the path $\text{goto}(F')$. If for instance $F' = \text{node}_a$ then only the a -labeled nodes of the graph are visited. The problem with the definition discussed so far, however, is that the clause $f(x) :- i(y), f'(x)$ is not SCL (see Definition 1). Even worse, it would lead to a quadratic evaluation time. Therefore, we replace it by the equivalent datalog program with two clauses $\{f(x) :- j(), f'(x).\} \cup \{j() :- i(y).\}$ which is SCL, leading to the definition of $Acc^{i,f}(\text{goto}(F'))$ in Fig. 9. Here $j \in \mathcal{P}_{int}$ is a fresh nullary predicate that is true for $Start^i(S)$ if and only if $S \neq \emptyset$.

Lemma 3. *For any path P , filter F , graph G , start set S , and monadic predicates $i, f, c, r \in \mathcal{P}_{int}$, the programs $Start^i(S)$, $Acc^{i,f}(P)$, $Filt^c(F)$, $Ex^{c,r}(P)$ are safe and SLC.*

Proof

Elementary by inspection of all cases of the definitions of these datalog programs. \square

Proposition 4. *For any filter query $F \in \mathcal{F}_\Sigma$, path query $P \in \mathcal{P}_\Sigma$, label $a \in \Sigma$, labeled graph G , subset $S \subseteq V$ of nodes of G , distinct monadic predicates $i, c, r \in \mathcal{P}_{int}$ and $x \in V$.*

1. *if $M = Filt^c(F) \cup Start^i(S)$ and $\vec{\ell} = i(x), c(x)$ then:*
 - $\llbracket \vec{\ell} \rrbracket_{M, db(G)} = \{[x/v] \mid v \in \llbracket F \rrbracket_G(S)\}$
 - $tdv_{M, db(G)}(\vec{\ell}) = tdn_{G, S}(F)$
2. *if $M = Ex^{c,r}(P) \cup Start^i(S) \cup \{r(x) :- \text{node}(x).\}$ and $\vec{\ell} = i(x), c(x)$ then:*
 - $\llbracket \vec{\ell} \rrbracket_{M, db(G)} = \{[x/v] \mid v \in S, \llbracket P \rrbracket_G(\{v\}) \neq \emptyset\}$
 - $tdv_{M, db(G)}(\vec{\ell}) = tdn_{G, S}(P)$
 - $reached_M(\vec{\ell}) = \llbracket P \rrbracket_G(S)$

Proof

By simultaneous induction on the structures of $F \in \mathcal{F}_\Sigma$ and $P \in \mathcal{P}_\Sigma$. We discuss a few cases for the possible forms of filters and paths only. For warming up, we consider the proof of property (1) for filters $F = F' \wedge F''$. The definition of the compiler yields:

$$\text{Filt}^c(F' \wedge F'') = \text{Filt}^{c'}(F') \cup \text{Filt}^{c''}(F'') \cup \{c(x) :- c'(x), c''(x).\}$$

Let $M = \text{Filt}^c(F' \wedge F'') \cup \text{Start}^i(S)$, $M' = \text{Filt}^{c'}(F') \cup \text{Start}^i(S)$, and $M'' = \text{Filt}^{c''}(F'') \cup \text{Start}^i(\llbracket F' \rrbracket_G(S))$. The top-down evaluator for $?-i(x), c(x)$. M will visit $\text{tdn}_{G,S}(F')$ in order to get all facts corresponding to the filter query F' . After that, the top-down evaluator will start from all nodes in S that satisfy filter F' to find those that also satisfy the filter F'' . First, we show that the top-down evaluation $\llbracket i(x), c(x) \rrbracket_M$ yields $\{[x/v] \mid v \in \llbracket F' \wedge F'' \rrbracket_G(S)\}$. By induction hypothesis,

$$\llbracket i(x), c'(x) \rrbracket_{M'} = \{[x/v] \mid v \in \llbracket F' \rrbracket_G(S)\}$$

and

$$\llbracket i(x), c''(x) \rrbracket_{M''} = \{[x/v] \mid v \in \llbracket F'' \rrbracket_G(\llbracket F' \rrbracket_G(S))\}.$$

Therefore,

$$\begin{aligned} \llbracket i(x), c(x) \rrbracket_M &= \{[x/v] \mid v \in \llbracket F'' \rrbracket_G \cap (\llbracket F' \rrbracket_G \cap S)\} \\ &= \{[x/v] \mid v \in \llbracket F' \wedge F'' \rrbracket_G \cap S\} \\ &= \{[x/v] \mid v \in \llbracket F' \wedge F'' \rrbracket_G(S)\}. \end{aligned}$$

Second, we show that the top-down visited sub-database is equal to $\text{tdn}_{G,S}(F' \wedge F'')$. By induction hypothesis,

$$\text{tdv}_{M', \text{db}(G)}(i(x), c'(x)) = \text{tdn}_{G,S}(F')$$

and

$$\text{tdv}_{M'', \text{db}(G)}(i(x), c''(x)) = \text{tdn}_{G, \llbracket F' \rrbracket_G(S)}(F'').$$

Therefore,

$$\begin{aligned} \text{tdv}_{M, \text{db}(G)}(i(x), c(x)) &= \text{tdn}_{G,S}(F') \cup \text{tdn}_{G, \llbracket F' \rrbracket_G(S)}(F'') \\ &= \text{tdn}_{G,S}(F' \wedge F''). \end{aligned}$$

We now turn to the most complicated case, which is the proof of property 2 for paths $P = P'/P''$. The definition of the compiler yields:

$$\text{Ex}^{c,r}(P'/P'') = \text{Ex}^{c,f}(P') \cup \text{Ex}^{f,r}(P'')$$

Let $M = \text{Ex}^{c,r}(P'/P'') \cup \text{Start}^i(S) \cup \{r(x) :- \text{node}(x).\}$, $M' = \text{Ex}^{c,f}(P') \cup \text{Start}^i(S) \cup \{f(x) :- \text{node}(x).\}$, and $M'' = \text{Ex}^{f,r}(P'') \cup \text{Start}^i(R) \cup \{r(x) :- \text{node}(x).\}$. We can divide the top-down evaluation of the datalog query $?-i(x), c(x)$. M into two steps. First — the top-down evaluation of the datalog query $?-i(x), c(x)$. M' reaches the set R of facts of the form $f(v)$ where $R = \text{reached}_{M',f}(i(x), c(x))$. The second step is equivalent to the evaluation of the datalog query $?-i(x), f(x)$. M'' where the set R is used as the set of starting nodes. By induction hypothesis, $R = \llbracket P' \rrbracket_G(S)$. Therefore, the result of top-down evaluation $\llbracket i(x), c(x) \rrbracket_{M, \text{db}(G)}$ is equal to the join of the result of evaluation $\llbracket i(x), c(x) \rrbracket_{M', \text{db}(G)}$ and result of evaluation $\llbracket i(x), f(x) \rrbracket_{M'', \text{db}(G)}$. By induction hypothesis,

$$\llbracket i(x), c(x) \rrbracket_{M', \text{db}(G)} = \{[x/v] \mid v \in S, \llbracket P' \rrbracket_G(\{v\}) \neq \emptyset\}$$

and

$$\llbracket i(x), f(x) \rrbracket_{M'', db(G)} = \{[x/v] \mid v \in R, \llbracket P'' \rrbracket_G(\{v\}) \neq \emptyset\}.$$

Let $\vec{\ell} = i(x), c(x)$. Therefore,

$$\begin{aligned} \llbracket \vec{\ell} \rrbracket_{M, db(G)} &= \{ \sigma \bowtie \Pi_\emptyset(\sigma') \mid \\ &\quad \sigma \in \{[x/v] \mid v \in S, \llbracket P' \rrbracket_G(\{v\}) \neq \emptyset\}, \\ &\quad \sigma' \in \{[x/v, y/v'] \mid v \in S, \\ &\quad (v, v') \in \llbracket P' \rrbracket_G, \llbracket P'' \rrbracket_G(\{v'\}) \neq \emptyset\} \} \\ &= \{[x/v] \mid v \in S, \exists v', v''. (v, v') \in \llbracket P' \rrbracket_G, \\ &\quad (v', v'') \in \llbracket P'' \rrbracket_G\} \\ &= \{[x/v] \mid v \in S, \llbracket P'/P'' \rrbracket_G(\{v\}) \neq \emptyset\}. \end{aligned}$$

Next, we show that the top-down visited sub-database is equal to $tdn_{G,S}(P'/P'')$. By induction hypothesis, we obtain

$$tdv_{M', db(G)}(i(x), c(x)) = tdn_{G,S}(P')$$

and

$$tdv_{M'', db(G)}(i(x), f(x)) = tdn_{G, \llbracket P' \rrbracket_G(S)}(P'').$$

Therefore,

$$\begin{aligned} tdv_{M, db(G)}(\vec{\ell}) &= tdn_{G,S}(P') \cup tdn_{G, \llbracket P' \rrbracket_G(S)}(P'') \\ &= tdn_{G,S}(P'/P''). \end{aligned}$$

Finally, we show that the set of all nodes v , such that the $r(v)$ is queried during the top-down evaluation of the datalog query $?-i(x), c(x)$. M , is equal to $\llbracket P'/P'' \rrbracket_G(S)$. The $r(v)$ can be queried only in the second step of the top-down evaluation which is equivalent to the evaluation of the datalog query $?-i(x), f(x)$. M'' . By induction hypothesis,

$$reached_{M'', r}(i(x), f(x)) = \llbracket P'' \rrbracket_G(\llbracket P' \rrbracket_G(S)).$$

Thus,

$$\begin{aligned} reached_{M, r}(i(x), c(x)) &= reached_{M'', r}(i(x), f(x)) \\ &= \llbracket P'' \rrbracket_G(\llbracket P' \rrbracket_G(S)) \\ &= \llbracket P'/P'' \rrbracket_G(S). \end{aligned}$$

□

Proposition 5. For any path query $P \in \mathcal{P}_\Sigma$, labeled graph G , subset S of nodes of G , distinct intensional predicates $i, f \in \mathcal{P}_{int}$ and $x \in \mathcal{V}$, if $M = Acc^{i,f}(P) \cup Start^i(S)$ then:

- $\llbracket f(x) \rrbracket_{M, db(G)} = \{[x/v] \mid v \in \llbracket P \rrbracket_G(S)\}$
- $tdv_{M, db(G)}(f(x)) = tdn_{G,S}(P)$

Proof

Let G be a graph with a node set V . The proof is induction on the structure of paths $P \in \mathcal{P}_\Sigma$. We only consider the case $P = P'/P''$. The definition of the compiler yields:

$$Acc^{i,f}(P'/P'') = Acc^{i,f'}(P') \cup Acc^{f',f}(P'')$$

Let $M = Acc^{i,f}(P'/P'') \cup Start^i(S)$, $M' = Acc^{i,f'}(P') \cup Start^i(S)$ and $M'' = Acc^{f',f}(P'') \cup Start^{f'}(R)$. We can divide the top-down evaluation of the datalog query $?-f(x)$. M into two steps. First — the top-down evaluation of the datalog query $?-f'$. M' provides

the set of nodes R equal to the set of all nodes v , such that $f'(v)$ is inferred after the first step of evaluation. The second step is equivalent to the evaluation of the datalog query $?-f(x)$. M'' where the set R is used as the set of starting nodes. By induction hypothesis, $R = \llbracket P' \rrbracket_G(S)$. Therefore,

$$\begin{aligned} \llbracket f(x) \rrbracket_{M, db(G)} &= \llbracket f(x) \rrbracket_{M'', db(G)} = \{[x/v] \mid \\ &\quad v \in \llbracket P'' \rrbracket_G(\llbracket P' \rrbracket_G(S))\} \\ &= \{[x/v] \mid v \in \llbracket P'/P'' \rrbracket_G(S)\}. \end{aligned}$$

Next, we show that the top-down visited sub-database is equal to $tdn_{G,S}(P'/P'')$. By induction hypothesis, we obtain

$$tdv_{M', db(G)}(f'(x)) = tdn_{G,S}(P')$$

and

$$tdv_{M'', db(G)}(f(x)) = tdn_{G, \llbracket P' \rrbracket_G(S)}(P'').$$

Therefore,

$$\begin{aligned} tdv_{M, db(G)}(f'(x)) &= tdn_{G,S}(P') \cup tdn_{G, \llbracket P' \rrbracket_G(S)}(P'') \\ &= tdn_{G,S}(P'/P''). \end{aligned}$$

□

Theorem 6. *For any graph G with subset of nodes S and any path query $P \in \mathcal{P}_\Sigma$ the answer set $\llbracket P \rrbracket_G(S)$ can be computed in time $\mathcal{O}(|P| |tdn_{G,S}(P)|)$.*

Proof

This follows from Proposition 5, Proposition 4 and Theorem 2. □

Appendix C Proofs for Section 7 (Jumping in Graphs)

It should be noticed that avoiding indexes of quadratic size may be relevant in practice, but more difficult to reach without restrictions. The index top_a , for instance, may be of the quadratic size, but only for XML documents that do not occur in practice. The choice of appropriate indexes raise many interesting research questions that are out of the scope of the present paper. It should also be mentioned that one may want to represent binary indexes in a more concise manner, rather than by enumeration of node pairs. For instance, for being able to jump to a -labeled nodes it is sufficient to store all a -labeled nodes, rather than pair of nodes (x, y) such that y is a -labeled.

Appendix D Proofs for Section 8 (Experiments)

Another line of improvement would be to stop the evaluation of filters when it has been proven correct. This effect may only be obtained if we use a datalog top-down evaluator that follows the *early completion* strategy, i.e. stops whenever a ground predicate (such as filter queries in our case) is proven true. But this strategy does not survive magic-set rewriting of Datalog programs, in order to mimic top-down evaluation in a bottom-up manner. Moreover, *early completion* does not allow us to define clearly a notion of needed nodes in a graph for a given query. A way out of this problem is to implement directly in

datalog what it means to traverse a set of nodes sequentially. For this, we need to assume that outgoing edges in graphs are ordered. This local order extends to a total order on paths starting at a given node by using a lexicographic order. Then we may implement a depth-first left-to-right traversal of the graph following this lexicographic order.