

Hybrid Reasoning on OWL RL

Jacopo Urbani ^{a,*}, Robert Piro ^b Frank van Harmelen ^a Henri Bal ^a

^a *Department of Computer Science, Vrije Universiteit Amsterdam, The Netherlands*

Email: {jacopo,frankh,bal}@cs.vu.nl

^b *Department of Computer Science, University of Oxford, United Kingdom*

Email: robert.piro@cs.ox.ac.uk

Abstract. Both materialization and backward-chaining as different modes of performing inference have complementary advantages and disadvantages. Materialization enables very efficient responses at query time, but at the cost of an expensive up front closure computation, which needs to be redone every time the knowledge base changes. Backward-chaining does not need such an expensive and change-sensitive precomputation, and is therefore suitable for more frequently changing knowledge bases, but has to perform more computation at query time.

Materialization has been studied extensively in the recent semantic web literature, and is now available in industrial-strength systems. In this work, we focus instead on backward-chaining, and we present a general hybrid algorithm to perform efficient backward-chaining reasoning on very large RDF datasets.

To this end, we analyze the correctness of our algorithm by proving its completeness using the theory developed in deductive databases and we introduce a number of techniques that exploit the characteristics of our method to execute efficiently (most of) the OWL RL rules. These techniques reduce the computation and hence improve the response time by reducing the size of the generated proof tree and the number of duplicates produced in the derivation.

We have implemented these techniques in an experimental prototype called QueryPIE and present an evaluation on both realistic and artificial datasets of a size that is between five and ten billion of triples. The evaluation was performed using one machine with commodity hardware and it shows that (i) with our approach the initial precomputation takes only a few minutes against the hours (or even days) necessary for a full materialization and that (ii) the remaining overhead introduced by reasoning still allows atomic queries to be processed with an interactive response time. To the best of our knowledge our method is the first that demonstrates complex rule-based reasoning at query time over an input of several billion triples and it takes a step forward towards truly large-scale reasoning by showing that complex and large-scale OWL inference can be performed without an expensive distributed hardware architecture.

1. Introduction

The amount of RDF data available on the Web calls for RDF applications that can process this data in an efficient and scalable way.

One of the advantages of publishing RDF data is that applications are able to infer implicit information by applying a reasoning algorithm on the input data. To this end, a predefined set of inference rules, which is complete w.r.t. some underpinning logic, can be applied in order to derive additional data.

Several approaches that perform rule-based inference were presented in the literature [19,11,27] and demonstrated reasoning upon several billion of triples. These methods apply the rules in a forward-chaining fashion, so that all the possible derivations are produced and stored together with the original input. While these methods exhibit good scalability because they can efficiently exploit computational parallelism, they have several disadvantages which compromise their use in real-world scenarios. First, they cannot efficiently deal with small incremental updates since they have to compute the complete materialization anew. Second, they become inefficient if the user is only in-

*Corresponding author

interested in a small portion of the entire input because forward-chaining needs to calculate all derivations.

Unlike forward-chaining, backward-chaining applies only inference rules depending on a given query. In this case, the computations required to determine the rules that need to be executed often become too expensive for interactive applications. Thus, backward-chaining has until now been limited to either small datasets (usually in the context of expressive DL reasoners) or weak logics (RDFS inference).

In this paper, we propose a method which materializes a fixed set of selected queries, *before* query time, whilst applying backward chaining *during* query time. This *hybrid* approach is a trade-off between a reduction in rule applications at query time and a small, query independent computation of data before query time. Our method relies on a backward-chaining algorithm to calculate the inference. The backward-chaining algorithm exploits the parallel computing power of modern architectures and uses the fact that a partial materialization of some selected queries is available to reduce the computation.

We thus have to tackle several problems. The first is to show that our reasoning algorithm is correct, i.e. it terminates, is sound and complete. We shall argue that the correctness is not dependent on a particular rule set, but holds for any Datalog program.

For the implementation and evaluation, however, we apply our method considering the semantics of the OWL RL fragment, which is the most recently standardized OWL profile designed to work on a large scale.

To this end, we have implemented these techniques in an experimental prototype called QueryPIE and tested the performance using artificial and realistic datasets of a size between five and ten billion triples. The evaluation shows that we are able to perform OWL reasoning using one machine equipped with commodity hardware which keeps the response time often below one second.

This paper is a revised and improved version of our initial work that was presented in [21]. More specifically, it extends the initial version that targets the *pD** fragment to one that supports most of the OWL RL rules, which are officially standardized by W3C. Also, this paper provides a theoretical analysis of the approach proving its correctness w.r.t. the considered rule set and presents an improved explanation and evaluation over larger datasets.

The remainder of this paper is organized as follows: in Section 2 we introduce the reader to our problem and provide a high level overview of our approach.

Next, in Section 3, we describe the backward-chaining algorithm that is used in our method to calculate the inference. Section 4 formalizes the precomputation algorithm of hybrid reasoning and proves its correctness. Section 5 focuses on the execution of most of the OWL 2 RL/RDF rule set¹ (henceforth called OWL RL rule set) presenting a series of optimizations to improve the performance on a large input. In Section 6 we present an evaluation of our approach using atomic queries on both realistic and artificial data. In Section 7 we report on related work. Finally, Section 8 concludes and gives directions for future work.

2. Hybrid reasoning: Overview

In principle, there are two different approaches to infer answers in a database with a given rule set: One is to compute the complete extension of a database under some given rule set *before* query time and the other is to infer only the necessary entries needed to yield a complete answer from the rule set on-demand, i.e. *at* query time.

The former's advantage is that querying reduces, after the full materialisation, to a mere lookup in the database and is therefore very fast compared to the latter approach, where for each answer a proof tree has to be built.

On the other hand, if the underlying database changes frequently, then a complete materialisation before query time has a severe disadvantage as the whole extension must be recomputed with each update. In this case, an on-demand approach has a clear advantage.

The approach presented in this paper positions itself in between: the answers for a carefully chosen set of queries are materialized before query time and added to the database. Answers to queries later posed by the user are inferred at query time.

Traditionally, each approach has been associated with an algorithmic method to retrieve the results: Backward-chaining was specifically aimed at on-demand retrieval of answers, only materialising as little information as necessary to yield a complete set of

¹Note that we do not consider all the rules, and this means we implement *incomplete* OWL RL reasoning.

Abbreviation	Full text
TYPE	<i>rdf:type</i>
SCO	<i>rdfs:subClassOf</i>
SPO	<i>rdfs:subPropertyOf</i>
EQC	<i>owl:equivalentClass</i>
EQP	<i>owl:equivalentProperty</i>
INV	<i>owl:inverseOf</i>
SYM	<i>owl:SymmetricProperty</i>
TRANS	<i>owl:TransitiveProperty</i>

Table 1

List of abbreviations for common URIs used in this paper.

answers, whilst forward-chaining applies the rules of the given rule set until the closure is reached.

Since we want to avoid complete materialization of the database, and therefore are only interested in specific answers, we use backward-chaining in both instances: we use backward-chaining to materialize only the necessary information for the carefully chosen queries which we then add to the database, and we use backward-chaining to answer the user queries.

To this end, we introduce a backward-chaining algorithm which exploits parallel computing power and the fact that some triple patterns are pre-materialized to improve the performance. For example, if one of these pre-materialized queries is requested at query-time, then the backward-chaining algorithm does not need to build the proof tree, but a lookup suffices. In case the pre-materialized patterns frequently appear at user query-time, such optimization is particularly effective.

To give an idea how this works, consider the following example.

Example 1. Consider the two following rules from the OWL RL rule set:

$$\begin{aligned} T(a, p1, b) &\leftarrow T(p, SPO, p1) \wedge T(a, p, b) \\ T(x, SPO, y) &\leftarrow T(x, SPO, w) \wedge T(w, SPO, y) \end{aligned}$$

where *SPO* is an abbreviation of “*rdfs:subPropertyOf*” and *a, b, p, p1, x, y, w* represent generic variables.

Assume we want to suppress the unfolding of all atoms of the form $T(x, SPO, y)$, modulo variable renaming. If we use Datalog to implement these rules in a program, then we can replace each atom by some new atom, say *S* that never appears in the head of the rules. After the substitution, Example 1 would be-

come:

$$\begin{aligned} T(a, p1, b) &\leftarrow S(p, SPO, p1) \wedge T(a, p, b) \\ T(x, SPO, y) &\leftarrow S(x, SPO, w) \wedge S(w, SPO, y) \end{aligned}$$

In general, the two programs do not yield the same answers for *T* anymore. To restore this equality for a given database \mathcal{I} we need to calculate all “ $T(x, SPO, y)$ ”-triples derivable from \mathcal{I} and add them to the auxiliary relation named *S* in \mathcal{I} . In our example this would mean that *S* contains the transitive closure of all “ $T(x, SPO, y)$ ”-triples which are inferable under the rule set in \mathcal{I} .

Note that “ $T(x, SPO, y)$ ”-triples can also be derived with the first rule if $p1 = SPO$. Furthermore, if *S* indeed contains the transitive closure of all “ $T(x, SPO, y)$ ”-triples the second rule can be rewritten as $T(x, SPO, y) \leftarrow S(x, SPO, y)$. \square

In the following, we discuss the backward-chaining algorithm that we use and explain how it exploits the pre-materialization for an efficient execution. After this, we will show that our method to replace the original rules with others is, after a small precomputation, harmless in the sense that everything which could be inferred under the original program can be inferred under the altered program and vice versa.

3. Hybrid Reasoning: Backward-chaining

We organized the content of this section as follows. In the following Section 3.1 we briefly recall some well-known notions of Datalog that will be frequently used. After this, we provide a theoretical description and analysis of our backward-chaining algorithm in sections 3.2 and 3.3 and discuss how its implementation can exploit a pre-materialization in section 3.4.

To ease the understanding of our explanation, we enrich the explanation of our approach with brief examples in order to facilitate the comprehension in case the reader is not completely familiar with the concepts that are being used.

In the current and following sections, we use the notation and notions that come from the Datalog theory, and more in particular from [1, Chapter 12], to formalize and to prove the correctness of our method.

We chose the Datalog theory to define our algorithms since the OWL RL inference rules are formulated in Datalog style; therefore, they can trivially be rendered into a Datalog program as already witnessed in Example 1.

3.1. Preliminaries

In this paper, we use abbreviations to indicate well-known URIs for reasons of space.² In our notations, we use as a convention fixed-width characters to denote constant terms (e.g. `SPO`) and italics for Datalog variables and predicate names (e.g. *a* or *T*). Note that in SPARQL [16] and the official OWL RL documentation variables are indicated with a preceding “?” (e.g. ?a).

For two functions $f : A \rightarrow B$ and $g : B' \rightarrow C$ with $B \subseteq B'$ we denote with $g \circ f$, pronounced *g after f*, the function $A \rightarrow C : x \mapsto g(f(x))$. For a set of functions $G := \{g \mid g : B_g \rightarrow C_g\}$ with $B \subseteq B_g$ for all $g \in G$ we define $G \circ f := \{g \circ f \mid g \in G\}$.

In Datalog, a signature SIG is a finite set of symbols which is the disjoint union of the set CONS of constants and the set PRED of predicate symbols. Each predicate symbol is associated with its *arity*, a positive natural number. A *database* \mathcal{I} for SIG is a pair consisting of a finite set $dom \mathcal{I}$, the *domain*, and an interpretation function $\cdot^{\mathcal{I}}$ whose domain is SIG.

If $R \in PRED$ is an n -ary predicate symbol, then $R^{\mathcal{I}}$, where $R^{\mathcal{I}} \subseteq (dom \mathcal{I})^n$, denotes the n -ary *relation* named R in the database \mathcal{I} . If $c \in CONS$ then $c^{\mathcal{I}} \in dom \mathcal{I}$. For our purposes, we assume $CONS = dom \mathcal{I}$ and $c^{\mathcal{I}} = c$ for all $c \in CONS$.

With VAR we denote a countably infinite set of variables where $VAR \cap SIG = \emptyset$. Thus, the set of all terms is $TERM = VAR \cup CONS$. For every term tuple \bar{t} we denote with $Var(\bar{t})$ the set of all variables in \bar{t} .

If $R \in PRED$ is an n -ary predicate symbol and \bar{t} is an n -ary term tuple, then $R(\bar{t})$ is called *atom*. An atom $R(\bar{t})$ is called *ground atom* if $\bar{t} \subseteq CONS$. If $R \in SIG$ and \mathcal{I} is a database over SIG we write $R(\bar{a}) \in \mathcal{I}$ if $R(\bar{a})$ is ground atom and $\bar{a} \in R^{\mathcal{I}}$. We continue *Var* on atoms by setting $Var(R(\bar{t})) := Var(\bar{t})$ for every atom $R(\bar{t})$.

A *substitution* is a mapping $\theta : V \rightarrow TERM$ where $V \subseteq VAR$. We call the substitution θ *assignment*, if $\theta(V) \subseteq CONS$ and θ is called *variable renaming* if $\theta(V) \subseteq VAR$ and θ is injective. $\theta_\varepsilon : \emptyset \rightarrow \emptyset$ is the *empty substitution*. Every substitution θ has a continuation $\tilde{\theta}$ on terms, where $\tilde{\theta}(t) = \theta(t)$ if $t \in V$ and $\tilde{\theta}(t) = t$ if $t \in TERM \setminus V$. Henceforth, we will denote θ but always implicitly refer to its continuation $\tilde{\theta}$. Similarly we apply substitutions or rather their continuations to term tuples $\theta(t_1, \dots, t_n) := (\theta(t_1), \dots, \theta(t_n))$ and atoms $\theta(R(\bar{t})) := R(\theta(\bar{t}))$.

We allow ourselves to represent a substitution θ as a set $\{t_0/t_1 \mid t_0 \in dom \theta \text{ and } \theta(t_0) = t_1\}$. The set-representation of θ_ε is simply \emptyset . We now define \bowtie that we will use later in Algorithm 1, line 21: A substitution θ_0 is *compatible* with a substitution θ_1 if $\theta_0(t) = \theta_1(t)$ for all $t \in dom \theta_0 \cap dom \theta_1$. For each pair of compatible substitutions θ_0, θ_1 we set $\theta_0 \cup \theta_1$ to be the substitution which corresponds to the union of their set-representation. θ_ε is compatible with every substitution θ and is neutral in the sense that $\theta_\varepsilon \cup \theta = \theta \cup \theta_\varepsilon = \theta$.

For sets Θ_0 and Θ_1 of substitutions we define $\Theta_0 \bowtie \Theta_1 := \{\theta_0 \cup \theta_1 \mid \theta_0 \in \Theta_0 \text{ compatible with } \theta_1 \in \Theta_1\}$. Using substitutions as results for database look-ups is not unusual and joins over sets of substitutions are particularly used in SPARQL [16].

Let $R_0(\bar{t}_0)$ and $R_1(\bar{t}_1)$ be two atoms such that $Var(\bar{t}_0) \cap Var(\bar{t}_1) = \emptyset$. A *unifier* for $R_0(\bar{t}_0)$ and $R_1(\bar{t}_1)$ is a substitution $\theta : (Var(\bar{t}_0) \cup Var(\bar{t}_1)) \rightarrow (Var(\bar{t}_0) \cup Var(\bar{t}_1))$ such that $\theta(R_0(\bar{t}_0)) = \theta(R_1(\bar{t}_1))$. A *most general unifier* (MGU) of two atoms $R_0(\bar{t}_0)$ and $R_1(\bar{t}_1)$ is a unifier θ for $R_0(\bar{t}_0)$ and $R_1(\bar{t}_1)$ such that for each unifier θ' of $R_0(\bar{t}_0)$ and $R_1(\bar{t}_1)$ there is a substitution σ with $\theta' = \sigma \circ \theta$.³ It is decidable whether or not a unifier for two given atoms exists. If it exists, an MGU exists and can be computed.

A Datalog *query* is an expression $q(\bar{t}_0) \leftarrow R_1(\bar{t}_1) \wedge \dots \wedge R_n(\bar{t}_n)$ where \bar{t}_0 is a term tuple, $R_i(\bar{t}_i)$ is an atom for each $i \in \{1, \dots, n\}$ and $Var(\bar{t}_0) \subseteq \bigcup_{1 \leq i \leq n} Var(\bar{t}_i)$. We omit \bar{t}_0 from $q(\bar{t}_0)$ if it is clear or not of interest. A tuple \bar{a} is an *answer* to $q(\bar{t}_0)$ w.r.t. \mathcal{I} if there is an assignment β with $\beta(\bar{t}_0) = \bar{a}$ and for all $i \in \{1, \dots, n\}$ we have $\beta(\bar{t}_i) \in R_i^{\mathcal{I}}$. The set of all answers to $q(\bar{t}_0)$ w.r.t. \mathcal{I} is denoted as $q(\bar{t}_0)^{\mathcal{I}}$.

We call a query *atomic* if $n = 1$ and will refer to it by its sole atom $R_1(\bar{t}_1)$. Answers in \mathcal{I} to such an atom $R_1(\bar{t}_1)$ are ground atoms $R_1(\bar{a}) \in \mathcal{I}$ such that \bar{t}_1 and \bar{a} unify.

Let \bar{t} and \bar{t}' be term tuples of the same length. Then \bar{t} is an *instance* of \bar{t}' , $\bar{t} \sqsubseteq \bar{t}'$, if there is a substitution σ such that $\sigma(\bar{t}') = \bar{t}$. Additionally, if $R(\bar{t})$ and $R(\bar{t}')$ are atoms we define $R(\bar{t}) \sqsubseteq R(\bar{t}')$ and say $R(\bar{t}')$ is *at least as general as* $R(\bar{t})$ iff $\bar{t} \sqsubseteq \bar{t}'$. $R(\bar{t})$ equals $R'(\bar{t}')$ up to variable renaming iff $R(\bar{t}) \sqsubseteq R'(\bar{t}')$ and $R'(\bar{t}') \sqsubseteq R(\bar{t})$.

Example 2. $(x, \text{TYPE}, \text{SYM}) \sqsubseteq (x, \text{TYPE}, y)$ where x, y are variables and TYPE and SYM are the abbreviation of Table 1. Also $(x, \text{TYPE}, y) \sqsubseteq (y, \text{TYPE}, x)$.

²Table 1 reports a list of all the abbreviations used in this paper.

³In literature, substitutions are used in post-fix notation, hence the condition for being an MGU is denoted as $\theta' = \theta\sigma$.

□

A Datalog rule has the form $R_0(\bar{t}_0) \leftarrow R_1(\bar{t}_1) \wedge \dots \wedge R_n(\bar{t}_n)$ such that for each $i \in \{0, \dots, n\}$, $R_i(\bar{t}_i)$ is an atom and $\text{Var}(\bar{t}_0) \subseteq \bigcup_{1 \leq i \leq n} \text{Var}(\bar{t}_i)$. We call $R_0(\bar{t}_0)$ *head atom* and all others *body atom*. A *Datalog program* is a finite set of Datalog rules. With $\text{Var}(r) \subseteq \text{VAR}$ we denote the set of all variables occurring in a Datalog rule r .

For any concrete given Datalog program P or Datalog query q which is applied to \mathcal{I} , we always assume that predicate and constant symbols occurring in P , and q respectively, are elements in SIG. We will rarely mention the signature since SIG is for a given database and program implicitly determined.

In our formalization, we denote the *immediate consequence operator* of a Datalog program P as T_P . T_P maps a database \mathcal{I} to the database $T_P(\mathcal{I})$, where $T_P(\mathcal{I})$ is \mathcal{I} extended by all facts that could be inferred from facts in \mathcal{I} under P . More formally, for each rule r with head atom $R_0(\bar{t}_0)$ we define $q_r(\bar{t}_0) \leftarrow R_1(\bar{t}_1) \wedge \dots \wedge R_n(\bar{t}_n)$ where the $R_i(\bar{t}_i)$ are the body atoms of r . For each $R \in \text{PRED}$ let $P \upharpoonright R$ be the set of rules whose predicate symbol in the head atom is R . We set $R^{T_P(\mathcal{I})} := R^{\mathcal{I}} \cup \bigcup_{r \in P \upharpoonright R} q_r^{\mathcal{I}}$. We define $T_P^0(\mathcal{I}) := \mathcal{I}$ and $T_P(\mathcal{I})$ to be the database \mathfrak{H} where $\text{dom } \mathfrak{H} = \text{dom } \mathcal{I}$ and $R^{\mathfrak{H}} := R^{T_P(\mathcal{I})}$ for all $R \in \text{PRED}$. We set further $T_P^{n+1}(\mathcal{I}) := T_P \circ T_P^n(\mathcal{I})$.

With ω we indicate the first infinite limit ordinal and set $T_P^\omega(\mathcal{I})$ to be the database \mathfrak{H} where $\text{dom } \mathfrak{H} = \text{dom } \mathcal{I}$ and $R^{\mathfrak{H}} := \bigcup_{n < \omega} R^{T_P^n(\mathcal{I})}$ for all $R \in \text{PRED}$. With $P(\mathcal{I})$ we denote the fully materialised database of \mathcal{I} under the program P . According to [1, Chapter 12], we have $P(\mathcal{I}) = T_P^\omega(\mathcal{I})$ and in particular there is $n < \omega$ such that $R(\bar{a}) \in T_P^n(\mathcal{I})$ for every ground atom $R(\bar{a}) \in P(\mathcal{I})$.

Let V be a set of *vertices* and $E \subseteq V \times V$ an *edge relation*, then $G := (V, E)$ is called *graph*. For vertices $v_0, v_1 \in V$ we call v_0 *predecessor* of v_1 , and v_1 *successor* of v_0 respectively, iff $(v_0, v_1) \in E$. A vertex that does not have a successor is called *leaf*. We define $E(v) := \{v' \in V \mid (v, v') \in E\}$ for each $v \in V$. For all $v \in V$ we define $\langle v \rangle_G$, the *v-subgraph* of G , to be the Graph $G' := (V', E')$ with $E' = E \cap (V' \times V')$ where V' is the smallest set such that $v \in V'$ and whenever $v' \in V'$ then $E(v') \subseteq V'$.

A *tree* is a graph such that each vertex has exactly one predecessor, except for one, the *root*, which has no predecessor. Trees can be considered to be recursively defined, where a tree G is either a single root, i.e. $V = \{v\}$, or a tree consists of a root v and each

of its successors $v' \in E(v)$ is a root of the tree $\langle v' \rangle_G$. The height of a finite tree is recursively derived as follows: if the root v is a leaf, then its height is 0, otherwise it is the maximal height of all trees $\langle v' \rangle$, where $v' \in E(v)$ plus 1. For arbitrary trees, the height might not be defined.

Let P be a Datalog program and \mathcal{I} a database and $R(\bar{a})$ a ground atom in $P(\mathcal{I})$. We call the pair (G, ℓ) a *Datalog proof-tree* for $R(\bar{a})$ in $P(\mathcal{I})$ if all of the following is satisfied: $G = (V, E)$ is a tree with a finite set V and $\ell : V \rightarrow \text{Atom}$ is a function, the *labelling function*, where Atom is the set of all ground atoms over SIG. Furthermore the root $v_0 \in V$ is labelled with $R(\bar{a})$, i.e. $\ell(v_0) = R(\bar{a})$, and for every leaf $v \in V$ we have $\ell(v)$ is a ground atom in \mathcal{I} . For every vertex $v \in V$ which is not a leaf there is a rule $r := R_0(\bar{t}_0) \leftarrow R_1(\bar{t}_1) \wedge \dots \wedge R_n(\bar{t}_n)$ and an assignment $\beta : \text{Var}(r) \rightarrow \text{dom } \mathcal{I}$, where $\text{Var}(r)$ denotes the set of all variables in r , such that $\ell(v) = \beta(R_0(\bar{t}_0))$ and there is a bijection $\iota : E(v) \rightarrow B$ between the successors of v and the set of body atoms $B := \{R_i(\bar{t}_i) \mid 1 \leq i \leq n\}$ of r such that $\ell(v') = \beta \circ \iota(v')$ for all $v' \in E(v)$.

A Datalog proof-tree thus represents a certain choice of rules whose application leads form atoms in \mathcal{I} to (possibly) derived atoms in $P(\mathcal{I})$. A proof by induction upon $m < \omega$ shows that every atom in $T_P^m(\mathcal{I})$ has a Datalog proof-tree of height at most m .

In the pseudocode of Algorithm 1 we use the function *lookup*: For an atomic query $Q = R(\bar{t})$ and a given database \mathcal{I} , or rather a given set of atoms \mathcal{I} , we define $\text{lookup}(R(\bar{t}), \mathcal{I})$ in line 13 to be the set $\{\theta : \text{Var}(\bar{t}) \rightarrow \text{dom } \mathcal{I} \mid R(\theta(\bar{t})) \in \mathcal{I}\}$ of substitutions.

Example 3. For any ground atom $R(\bar{a})$ we have

1. $\text{lookup}(R(\bar{a}), \mathcal{I}) = \emptyset$ iff $R(\bar{a}) \notin \mathcal{I}$
2. $\text{lookup}(R(\bar{a}), \mathcal{I}) = \{\theta_\varepsilon\}$ iff $R(\bar{a}) \in \mathcal{I}$.

□

Hence $\text{lookup}(R_0(\bar{t}_0), \mathcal{I}) \bowtie \text{lookup}(R_1(\bar{t}_1), \mathcal{I})$ is the set containing all assignments $\theta : (\text{Var}(\bar{t}_0) \cup \text{Var}(\bar{t}_1)) \rightarrow \text{dom } \mathcal{I}$ which are assignments for both atoms in the sense that $R_0(\theta(\bar{t}_0)) \in \mathcal{I}$ and $R_1(\theta(\bar{t}_1)) \in \mathcal{I}$.

Finally, we call an atomic query Q *blocked* in the procedure call $\text{infer}(Q, \text{PrevQueries})$ if there is $Q' \in \text{PrevQueries}$ such that $Q \sqsubseteq Q'$ and $Q' \sqsubseteq Q$.

3.2. Backward-Chaining

The purpose of a backward-chaining algorithm is to derive for a given database \mathcal{J} and a Datalog program P all possible ground atoms $R(\bar{a}) \in P(\mathcal{J})$ that are answers to a given query atom Q .

Traditionally, users interact with RDF datasets using the SPARQL language [16] where all the triple patterns that constitute the body of the query are joined together according to some specific criteria. In this paper, we do not consider the problem of efficiently joining the RDF data and focus instead on the process of retrieving the triples that are needed for the query. Therefore, we target our reasoning procedure at *atomic* queries, e.g., $T(?a, SCO, ?b)$.

The algorithm that we present is a variation of the well-known algorithm QSQ (Query-subquery) [24, 1, 6]. The variations that we introduce are meant to exploit the computational parallelism that is possible to obtain by using modern architectures.

The general idea behind the QSQ algorithm is to recursively rewrite the given query into many subqueries until no more rewritings can be performed and the subqueries can only be evaluated against the knowledge base.

Example 4. To give an idea on how QSQ works, suppose that our initial query is

$$T(x, \text{TYPE}, \text{Person})$$

and that we have a generic database \mathcal{J} and the OWL RL rule set as P . Initially, the algorithm will determine which rules can produce a derivation that is part of the input query. For example, it could apply the subclass and subproperties inheritance rules (`cax-sco` and `prp-spo1` in the OWL RL rule set). After it has determined them, it will move to the body of the rules and proceed evaluating them. In case these subqueries will produce some results, the algorithm will execute the rules and return the answers to the upper level.

With this process, the algorithm is creating a tree that has the original query as root and the rules and subqueries that might contribute to derive some answers as the internal nodes.

This tree represents all the derivation steps that are taken to derive answers of our initial query (the root) starting from some existing facts (the leaves). In Figure 1 we report an example of such a tree for our example query. \square

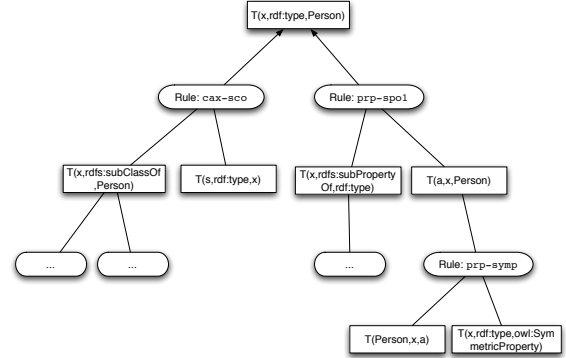


Fig. 1. Example of the execution of backward-chaining for the input query $T(x, \text{TYPE}, \text{Person})$ and the OWL RL rules.

The original QSQ algorithm was introduced in 1986 [23]. Unfortunately, the first version of this algorithm was found incomplete, but a fixed-up version was presented already the year after [24, 6].

In principle, the QSQ algorithm extends the standard SLD resolution technique [22] by applying it to a set of tuples instead of single ones [6]. An important problem of backward-chaining algorithms such as QSQ concerns the execution of recursive rules. Recursive rules and more in general cycles in the proof tree are an important threat since they could create infinite loops in the computation.

To solve this problem, QSQ extends the standard SLD resolution adding two techniques: an *Admissibility test* and a *Lemma resolution*. The resulting method, called SLD-AL, rests on the method which QSQ and all its variants is derived from. By analyzing the properties of the SLD-AL resolution, the author of QSQ has proved in [24] that this algorithm always terminates and is complete (i.e. is able to calculate all the answers).

In this paper, we do not re-propose a complete description of SLD-AL. We kindly refer the reader to the original explanation [24] or alternatively in [25], if the reader is interested in this topic. Here, we will simply sketch some characteristics of the technique which is important in our context.

In very general terms, SLD-AL is an abstract technique to construct a computational *tree* to answer a given query. This tree, which is called *SLD tree*, resembles the tree in Figure 1, and is explored by an algorithm such as QSQ in order to find all the answers (called atomic lemmas) that satisfy the input query. It is crucial that the SLD tree is *finite* and *complete*, otherwise, it would be impossible for an algorithm to ter-

minate and compute all the derivations. These properties do hold for SLD-AL trees, as shown in [25].

The main idea behind the *AL* technique is the following: when the algorithm computes a tree and needs to evaluate a new query, it applies an “admissability test”, to verify whether this query can be resolved using the rules. If the new query is “similar” to a previous one, then the query is evaluated using only the answers produced so far (“lemma resolution”). In this way, the program does not get trapped in a loop generating an infinite tree.

The SLD tree can be explored in several ways. In the general case, the search strategy can determine the completeness of the approach, but in Datalog any strategy is equivalent. For example, the tree can be constructed using an iterative depth-first strategy as proposed with the original QSQ algorithm in [23], or using a more sophisticated constructive search as in QoSaq (QSQ+gObAlloptimization) [25].

In the early research on these methods, much emphasis was put on optimizing the computation to avoid redundancy. For example, the original QSQ algorithm traverses the tree using a dept-first strategy, and repeats the query execution at every node until it retrieves all answers for that subgoal [6]. In this way, it is able cache the results for that subquery and reuse them during the evaluation of the rest of the tree. QoSaq [25] goes further than the original QSQ algorithm proposing one more solution to this problem: here, a copy of the tree is maintained in main memory and a “waking” mechanism is used to feed new answers derived in one node to other equivalent queries that appear in other branches of the tree.

These techniques are very efficient in optimizing the resolution process, but have the drawback that are very difficult to implemented in a parallel (and possibly distributed) environment. In fact, the original depth-first strategy used by QSQ require a sequential search of the tree, otherwise expensive synchronization mechanisms must be used. The “waking” mechanism proposed by QoSaq cannot be applied in a distributed environment, since it requires to maintain a copy of the tree in main memory and this mechanism would be slowed down by a network access.

Therefore, we adapted the original QSQ algorithm so that it can be more easily parallelized paying the price of duplicate derivation and possibly redundant computation. We will present the algorithm in the following section and show that the properties of termination soundness, and completeness are still valid. After this, in Section 3.4, we will describe how we can

exploit the precomputation of some queries to increase the performance of backward-chaining.

3.3. Our approach

We introduced two key differences from the original QSQ algorithm, which aim is to improve the parallelization of the computation:

- Unlike QSQ, our algorithm does not construct the proof-tree sequentially but in parallel by applying the rules on separate threads and in an asynchronous manner. For example, if we look back at Figure 1, the execution of rules `cax-sco` and `prp-spol` is performed concurrently by different threads. This execution strategy makes the implementation and the maintenance of the global data structure, used for caching results of previous queries, difficult and inefficient. We hence chose to replace this mechanism with one that only remembers which queries were already executed along single paths of the tree. Whilst such a choice might lead to some duplicate answers because the same queries can be repeated multiple times in different parts of the proof tree, it allows the computation to be performed in parallel limiting the usage of expensive synchronization mechanisms;
- Because the proof tree is built in parallel, repeating every query until no new results are found is an inefficient operation: the same query can appear multiple times in different parts of the tree. Therefore, we replace it with a global loop that is performed only at the root level of the tree and that stores during every iteration all the intermediate derivations.

We report the algorithm using pseudocode in Algorithm 1. Before explaining the algorithm in detail, we will show for the construction of the MGU in line 16 that we may w.l.o.g. assume that for every Datalog program P and atomic query Q a finite variable set V suffices such that for every query Q' occurring in the computation of $\text{infer}(Q, \emptyset)$ and $r \in P$ a variable renaming ρ exists such that $\text{Var}(Q') \cap \rho \circ \text{Var}(r) = \emptyset$ and $\rho \circ \text{Var}(r) \subseteq V$.

To this end let \mathcal{A} be the set of all atoms occurring in P unified with $\{Q\}$. We set V_0 to be the set of all variables in \mathcal{A} and $V := V_0 \cup V_1$ where V_1 is a disjoint copy of V_0 . We define \mathcal{Q} to be the closure of \mathcal{A} under all substitutions $\theta : V \rightarrow (V \cup \text{CONST})$ where

Algorithm 1 Backward-chaining algorithm:

\mathcal{I} and P are global constants, where \mathcal{I} is a finite set of facts and P is a Datalog program. Tmp and Mat are global variables, where Mat stores results of the previous materialization round and Tmp stores the results of the current round. The parameter Q represents an input pattern, i.e. an atom. Both functions *main* and *infer* return a set of atoms, whilst the function *lookup* returns a set of substitutions. We say $Q \in PrevQueries$ (cf. line 15) iff there is $Q' \in PrevQueries$ s.t. $Q \sqsubseteq Q'$ and $Q' \sqsubseteq Q$.

```

1  function main(Q)
2    New, Tmp, Mat :=  $\emptyset$ 
3    repeat
4      Mat := Mat  $\cup$  New  $\cup$  Tmp
5      New := infer(Q,  $\emptyset$ )
6    until New  $\cup$  Tmp  $\subseteq$  Mat  $\cup$   $\mathcal{I}$ 
7    return New
8  end function
9
10 function infer(Q, PrevQueries)
11
12 //This cycle is executed in parallel
13 all_subst := \lookup(Q,  $\mathcal{I}$   $\cup$  \Mat)
14 for ( $\forall r \in P$  s.t.  $Q$  is unifiable
15   with  $r.HEAD$  and  $Q \notin PrevQueries$ )
16    $\theta_h$  := MGU(Q,  $r.HEAD$ )
17   subst :=  $\{\theta_\varepsilon\}$ 
18   for  $\forall p \in r.BODY$ 
19     tuples := infer( $\theta_h(p)$ , PrevQueries  $\cup$  Q)
20     Tmp := Tmp  $\cup$  tuples
21     subst := subst  $\bowtie$  \lookup( $\theta_h(p)$ , tuples)
22   end for
23   all_subst := all_subst  $\cup$  (subst  $\circ$   $\theta_h$ )
24 end for
25
26 return  $\bigcup_{\theta \in all\_subst} \{\theta(Q)\}$ 
27
28 end function

```

CONST are the constants occurring in \mathcal{A} . Obviously, Q is finite.

Lemma 1.

1. For every $Q' \in \mathcal{Q}$ and every $r \in P$ there is a variable renaming, i.e. an injective function $\rho : Var(r) \rightarrow V$ such that $Var(Q') \cap \rho \circ Var(r) = \emptyset$.
2. We have $Q' \in \mathcal{Q}$ for every subsequent procedure call *infer*(Q' , PrevQueries) of *infer*(Q , \emptyset).

Proof.

1. For a set M we denote with $|M|$ its cardinality. Let $Q' \in \mathcal{Q}$ and $r \in P$ be arbitrary. For every atom $R(\bar{t}) \in \mathcal{A}$ we know $|Var(R(\bar{t}))| \leq \frac{1}{2}|V|$. Since Q' is the result of a variable substitution in an atom in \mathcal{A} , we have $|Var(Q')| \leq \frac{1}{2}|V|$ and hence that $|V'| \geq \frac{1}{2}|V|$, where $V' :=$

$V \setminus Var(Q')$. Similarly we know that $|Var(r)| \leq \frac{1}{2}|V|$. Hence there is an injective substitution $\rho : Var(r) \rightarrow V' \subseteq VAR$.

2. The proof is carried out via the nesting depth $k < \omega$ of precedence call. If $k = 0$ the procedure call is *infer*(Q , \emptyset) and since $Q \in \mathcal{Q}$ the claim is true. Assume we are in a subsequent procedure call *infer*(Q' , *infer*) where $Q' \in \mathcal{Q}$. For every body atom p of r used in a loop pass (line 14–24), item 1. yields a variable renaming ρ such that $Var(Q') \cap \rho \circ Var(r) = \emptyset$ and $\rho \circ Var(r) \subseteq V'$. Let the variables of r be renamed by ρ , then $Var(MGU(Q', $r.HEAD$)(p)) \subseteq V'$ for all body atoms of r and hence $MGU(Q', $r.HEAD$)(p) \in \mathcal{Q}$. □

The procedure *main* is the main function used to invoke the backward-chaining procedure for a given atomic query Q . *main* returns the derived answers for the input query. The procedure consists of a loop in which the recursive function *infer* is invoked with the input query. This function returns all the derived answers for Q that were calculated by applying the rules using backward-chaining (line 5) and all the intermediate answers that were inferred in the process, saved in the global variable *Tmp*. In each loop pass the latest results in *Tmp* and *New* are checked against the accumulated answers of the previous runs in *Mat* and \mathcal{I} . If no new tuple was derived then the loop terminates.

After this loop has terminated, the algorithm returns *New* (line 7) which contains after the last loop pass all answers to the input query (cf. line 13).

The function *infer* is the core of the backward-chaining algorithm. Using the function *lookup*, it first retrieves for the formal parameter Q all answers which are facts in the database or were previously derived (line 13). After this, it determines the rules that can be applied to derive new answers for Q (lines 14–15) and calculates the substitution θ_h to unify the head of the applicable rule with the query Q (line 16).

It proceeds with evaluating the body of the rule (lines 16–23) storing in *tuples* and *Tmp* the retrieved answers (lines 19–20), and performing the joins necessary according to the rule body (line 21).

In line 23, each assignment in *subst* is composed with the unifier under which it has been derived in the for-loop (line 18–22), which renders it into an assignment for Q . All these assignments are eventually copied into *all_subst* from which a set of answers for Q is derived (line 26) The answers are then returned to

the function caller. After the whole recursion tree has been explored exhaustively, *infer* returns control to the function *main*, where the derived answers are copied into the variable *New*. The process is repeated until the closure is reached.

To facilitate the understanding of this algorithm and more in particular of the function *infer*, consider the following example:

Example 5. Suppose that we have a program *P* that consists of a single rule,

$$r_0 := T(x, \text{TYPE}, y) \leftarrow T(z, \text{SCO}, y) \wedge T(x, \text{TYPE}, z)$$

the input query is $Q := (a, \text{TYPE}, u)$ and \mathcal{I} contains solely $T(a, \text{TYPE}, c)$ and $T(c, \text{SCO}, d)$.

The call *main*(*Q*) executes *infer*(*Q*, \emptyset) which performs a lookup in line 13, setting $\text{all_subst} = \{\{u/c\}\}$. It is then checked for all rules *r* whether the head of *r* is unifiable with *Q* (lines 14–15). In our example, only rule r_0 satisfies this condition, and the algorithm computes some MGU θ_h (line 16), e.g. $\theta_h := \{x/a, y/u\}$.

θ_h is applied to each body atom of r_0 and *infer* is recursively invoked. In our example, the first recursive call would be *infer*($T(z, \text{SCO}, u), \{Q\}$). We obtain at first the substitution $\{z/c, u/d\}$ in line 13 which is stored in the local variable *all_subst*. Since no rule head unifies with $T(z, \text{SCO}, u)$ the program jumps to line 26 and returns $\{\{z/c, u/d\}(T(z, \text{SCO}, u))\} = \{T(c, \text{SCO}, d)\}$. Hence the join in line 21 evaluates to $\{\{z/c, u/d\}\}$ and *subst* = $\{\{z/c, u/d\}\}$ after line 21.

The inner for-loop (lines 18–22) moves to the next predicate and launches the second recursive call *infer*($T(a, \text{TYPE}, z), \{Q\}$). Line 13 sets *all_subst* to $\{z/c\}$ and forgoes lines 14–24 since $T(a, \text{TYPE}, z)$ is a blocked query, i.e. $T(a, \text{TYPE}, z)$ equals *Q* up to variable renaming and *PrevQueries* = $\{Q\}$.

Returning $\{T(a, \text{TYPE}, c)\}$ to the computation one recursion level above, the result of the join $\{\{z/c, u/d\}\} \bowtie \{\{z/c\}\} = \{\{z/c, u/d\}\}$ in line 21 will be stored in *subst*. Line 23 sets *all_subst* to $\{\{u/c\}, \{z/c, u/d, x/a, y/d\}\}$ which is the result of $\{\{u/c\}\} \cup \{\{z/c, u/d\}\} \circ \{x/a, y/u\}$.

Each substitution in *all_subst* is applied to *Q* in line 26 and $\text{Result} := \{T(a, \text{TYPE}, c), T(a, \text{TYPE}, d)\}$ is returned to the function *main* where they are stored in *Mat*. The function *main* will repeat the computation once again to ensure that no more triples can be derived. In this last loop-pass *infer*(*Q*, \emptyset) returns the atoms in *Result* which it obtains this time from a

lookup in line 13. *Result* is stored in *New* which is then returned to the user (line 7). \square

We will now discuss the correctness of our algorithm, which are termination, soundness and completeness. To this end, we will furnish further lemmas and definitions:

Consider database \mathcal{I} , a Datalog program *P* and an atom *Q*. Let \mathcal{Q} be the set of atoms defined above Lemma 1. With \supseteq (being a proper super-set) we obtain a well-founded order on the powerset of \mathcal{Q} which we shall use for further inductive arguments.

Lemma 2. *Each call infer(Q, \emptyset) entails at most finitely many recursive calls to infer.*

Proof. An inductive argument over the nesting depth of procedure calls shows that in the *n*-th nested procedure call, *PrevQueries* contains *n* elements. Since $\text{PrevQueries} \subseteq \mathcal{Q}$ (cf. Lemma 1 item 2.) and both for-loops iterate over finite sets, *infer* is called at most finitely many times. \square

Lemma 3. *Let $P(\mathcal{I})$ be the minimal model of \mathcal{I} under *P*. If $\text{Mat}, \text{Tmp} \subseteq P(\mathcal{I})$ then for all $Q \in \mathcal{Q}$ and all subsets $\text{PrevQueries} \subseteq \mathcal{Q}$ we have that*

$$\text{infer}(Q, \text{PrevQueries}) \subseteq P(\mathcal{I})$$

and $\text{Tmp}' \subseteq P(\mathcal{I})$ where Tmp' is *Tmp* after the execution of *infer*(*Q*, *PrevQueries*).

Proof. The proof is carried out by induction. For the base case let $Q \in \mathcal{Q}$ be arbitrary and execute *infer*(*Q*, \mathcal{Q}). Then *lookup*(*Q*, $\mathcal{I} \cup \text{Mat}$) in line 13 returns a set of assignments $\beta : V \rightarrow \text{dom } \mathcal{I}$, where *V* contains all variables in *Q*, such that $\beta(Q) \in (\mathcal{I} \cup \text{Mat})$. Since $Q \in \mathcal{Q}$ lines 14–24 are skipped and hence

$$\text{infer}(Q, \mathcal{Q}) \subseteq P(\mathcal{I})$$

and $\text{Tmp}' = \text{Tmp} \subseteq P(\mathcal{I})$.

Assume as induction hypothesis that for all \mathcal{R} with $\text{PrevQueries} \subsetneq \mathcal{R} \subseteq \mathcal{Q}$ and for all $Q \in \mathcal{Q}$ we have that if $\text{Mat}, \text{Tmp} \subseteq P(\mathcal{I})$ then $\text{infer}(Q, \mathcal{R}) \subseteq P(\mathcal{I})$ and $\text{Tmp}' \subseteq P(\mathcal{I})$ where Tmp' is the updated set *Tmp* after the procedure call *infer*(*Q*, \mathcal{R}).

For the induction step let $\text{PrevQueries} \subseteq \mathcal{Q}$ and $Q \in \mathcal{Q}$ be arbitrary and assume $\text{Mat}, \text{Tmp} \subseteq P(\mathcal{I})$. Execute *infer*(*Q*, *PrevQueries*):

As in the induction base, *all_subst* contains after line 13 only assignments β such that $\beta(Q) \in (\mathcal{I} \cup$

Mat). If $Q \in \text{PrevQueries}$ up to variable renaming, we skip lines 14–24 ending up with the same outcome as in the induction base.

Hence assume $Q \notin \text{PrevQueries}$ modulo variable renamings. In each loop-pass of the outer-loop and inner-loop, the induction hypothesis yields that

$$\text{infer}(\theta_h(p), \text{PrevQueries} \cup \{Q\}) \subseteq P(\mathcal{J}).$$

Hence $\text{Tmp}' \subseteq P(\mathcal{J})$.

Similar to line 13, $\text{lookup}(\theta_h(p), \text{tuples})$ contains exactly those assignments $\beta : V \rightarrow \text{dom } \mathcal{J}$, where V contains all variables in $\theta_h(p)$, such that $\beta \circ \theta_h(p) \in \text{tuples}$. As remarked below the definition of \bowtie , every $\beta \in \text{subst}$ after line 22 satisfies $\beta \circ \theta_h(p) \in \text{tuples}$ for all $p \in r.\text{BODY}$. Datalog requires all variables of the head to be covered by some atom in the body, so we know that each $\beta \in \text{subst}$ is an assignment for $\theta_h(r.\text{HEAD})$ where $\beta \circ \theta_h(r.\text{HEAD}) \in P(\mathcal{J})$. Additionally, since θ_h was the unifier for Q and $r.\text{HEAD}$, we know that every $\beta \circ \theta_h$ in all_subst is an assignment for Q . Hence, we have $\theta(Q) \in P(\mathcal{J})$ for all $\theta \in \text{all_subst}$ which shows the claim. \square

3.3.1. Termination.

We first concentrate on the termination of the procedure call $\text{infer}(Q, \emptyset)$ in the function *main*. Let \mathcal{J} be a database over a finite signature SIG: by definition $\text{dom } \mathcal{J}$ is finite. Lemma 3 yields that in every repeat-loop pass (lines 3–6)

$$\text{New}, \text{Tmp}, \text{Mat} \subseteq P(\mathcal{J}) \quad (*)$$

where $P(\mathcal{J})$ is the minimal model of \mathcal{J} under P . An inductive argument over $\text{PrevQueries} \subseteq \mathcal{Q}$ shows that under the precondition $(*)$ for all $Q \in \mathcal{Q}$ every procedure call $\text{infer}(Q, \text{PrevQueries})$ terminates: Lemma 2 shows that there are only finitely many calls to *infer*, but we can now show that also every call to *lookup* in line 13 and line 21 yields a finite result via finite computation and thus \bowtie in line 21 terminates.

In every repeat-loop pass *Tmp* or *New* grow or the loop is terminated. Since *Tmp* and *New* are bounded by $P(\mathcal{J})$, which is finite, the repeat-loop terminates after finitely many passes. This shows that for every database \mathcal{J} , every Datalog program P and every atomic query Q the function *main*(Q) terminates.

3.3.2. Soundness.

The assingment of *New* in line 5 of the last repeat-loop pass is the return value of *main*(Q). In order to show soundness, we have to show that the return value of $\text{infer}(Q, \emptyset)$ only contains answers to Q from $P(\mathcal{J})$. To this end, assume \mathcal{J} is a database, P is a Datalog program and Q is an atomic query. In every repeat-loop pass, $\text{infer}(Q, \emptyset)$ contains only ground atoms from $P(\mathcal{J})$ that unify with Q : Using Lemma 3, we know that $\text{Tmp}, \text{Mat} \subseteq P(\mathcal{J})$ for every repeat-loop pass (lines 3–6). Thus, line 13 only yields assignments β such that $\beta(Q) \in (\mathcal{J} \cup \text{Mat})$ and hence such that $\beta(Q) \in P(\mathcal{J})$. Using Lemma 3 on line 19, we know that $\beta \circ \theta_h(p) \in P(\mathcal{J})$ for every assignment $\beta \in \text{lookup}(\theta_h(p), \text{tuples})$ in line 21 and every body atom p of a rule whose head unifies with Q . Hence $\beta \in \text{subst}$ after line 22 is an assingment such that $\beta \circ \theta_h(Q) \in P(\mathcal{J})$. Hence the returned set in line 26 only contains ground atoms which are answers to Q and are in $P(\mathcal{J})$.

3.3.3. Completeness.

Let \mathcal{J} be a given database, P a given Datalog program and $Q := R'(\bar{t})$ an atomic query. To show the completeness of Algorithm 1, we need to prove that every atom in $R'^{P(\mathcal{J})}$ which unifies with Q can be derived by *main*(Q). This claim is shown via Proposition 1 below, as the latter holds in particular for all answers to the input query Q derived under P from \mathcal{J} . Proposition 1 however rests on Lemma 4. The lemma states that if a ground atom $R(\bar{a})$ appears as label in a Datalog proof-tree for some answer $R'(\bar{b})$ to Q , then $R(\bar{a})$ is an answer to some query Q_n which will occur during a computation of $\text{infer}(Q, \emptyset)$. Note that it does not show that $\text{infer}(Q, \emptyset)$ derives $R(\bar{a})$.

Lemma 4. *Let $Q := R'(\bar{t})$ be an atomic query and $R(\bar{a})$ a label, which appears in the Datalog proof-tree (G, ℓ) for some answer to Q in $P(\mathcal{J})$. Then there is a subsequent procedure call $\text{infer}(Q_n, \text{PrevQueries})$ of $\text{infer}(Q, \emptyset)$ such that Q_n is a non-blocked atomic query and $R(\bar{a})$ is an answer to Q_n derived under P in \mathcal{J} .*

Proof. Let $R'(\bar{t})$ be an answer to Q which has a Datalog proof-tree (G, ℓ) in $P(\mathcal{J})$ containing a label $R(\bar{a})$. We have to show, that there is a sequence of atomic queries Q_0, \dots, Q_n such that

1. $Q = Q_0$ and $R(\bar{a})$ unifies with Q_n
2. for each $i \in \{0, \dots, n\}$ there is a rule $r \in P$ and $\theta := \text{MGU}(Q_i, r.\text{HEAD})$ such that $Q_{i+1} = \theta(p)$, where p is some body-atom of r

3. no query is blocked, i.e. there is no subsequence $Q_i \dots Q_k$ with $0 \leq i < k \leq n$ such that Q_i is up to variable renaming equal to Q_k ($Q_i \sqsubseteq Q_k$ and $Q_k \sqsubseteq Q_i$).

3. guarantees in particular that the condition $Q \notin \text{PrevQueries}$ in line 15 is true when the procedure call $\text{infer}(Q_i, \{Q_0, \dots, Q_{i-1}\})$ is executed. Hence, if 1–3 holds, $\text{infer}(Q, \emptyset)$ will, according to lines 14–24 of Algorithm 1, call in ascending sequence $\text{infer}(Q_i, \{Q_0, \dots, Q_{i-1}\})$ where $0 \leq i \leq n$.

In a first step we specify a sequence of rules r_0, \dots, r_n which leads from $R'(\bar{b})$ to the atom $R(\bar{a})$ and then we determine a sequence of Queries Q_0, \dots, Q_n .

Let $R'(\bar{b})$ be the atom which unifies with the input query Q and in whose Datalog proof-tree (G, ℓ) $R(\bar{a})$ appears. Then either (G, ℓ) has height 0 and thus $R'(\bar{b}) = R(\bar{a})$ and $Q_n := Q$ is found, or there is a sequence of rule applications r_0, \dots, r_n such that $R'(\bar{b})$ unifies with the head of r_0 via some MGU θ_0 and for all $i \in \{0, \dots, n-1\}$ some unified body-atom $B_{i,k_i} = \theta_i(R_{i,k_i}(\bar{t}_{i,k_i}))$ of r_i unifies via some MGU θ_{i+1} with the head of r_{i+1} and finally $R(\bar{a})$ unifies with some unified body-atom $B_{n,k_n} = \theta_n(R_n(\bar{t}_{n,k_n}))$ of r_n .

Fix this sequence of rules r_0, \dots, r_n . Since $R'(\bar{b})$ was an answer to Q (i.e. a ground atom) and unifies with the head atom H_0 of r_0 , Q unifies with H_0 yielding $\vartheta_0 := \text{MGU}(Q, H_0)$. For all $i \in \{0, \dots, n-1\}$ the unified body-atom $Q_i := \vartheta_i(R_{i,k_i}(\bar{t}_{i,k_i}))$ of r_i unifies with the head H_{i+1} of r_{i+1} yielding $\vartheta_{i+1} := \text{MGU}(Q_i, H_{i+1})$, so that we finally reach the body atom $R_{n,k_n}(\bar{t}_{n,k_n})$ of r_n where $Q_n := \vartheta_n(R_{n,k_n}(\bar{t}_{n,k_n}))$ is the query which unifies with $R(\bar{a})$.

We hence obtain a sequence Q_0, \dots, Q_n satisfying items 1 and 2. We shall show that for every sequence satisfying items 1 and 2 there is a sequence Q'_0, \dots, Q'_m satisfying items 1–2 and 3:

The claim is clear, if the sequence is of length 1: Q_0 is never blocked. Let $Q_0 \dots Q_n$ be a sequence of length $n+1$ with Q_i equals Q_k up to variable renaming where $0 \leq i < k \leq n$. Then the head of r_{k+1} unifies with the query Q_i . The sequence $Q_0, \dots, Q_i, Q_{k+1} \dots Q_n$ is properly shorter than $Q_0 \dots Q_n$ and satisfies items 1–2. The induction hypothesis yields a sequence Q'_0, \dots, Q'_m which satisfies items 1–3. \square

Proposition 1. *Let P, \mathfrak{J} and Q be fixed and \mathcal{Q} as defined above Lemma 2. Let $R'(\bar{t})$ be an answer to*

Q which has a Datalog proof-tree (G, ℓ) in $P(\mathfrak{J})$ containing a label $R(\bar{a})$. Then there is a repeat-loop pass in $\text{main}(Q)$ from which onward for every atomic query $Q' \in \mathcal{Q}$ which unifies with $R(\bar{a})$ and for all $\text{PrevQueries} \subseteq \mathcal{Q}$ every call $\text{infer}(Q', \text{PrevQueries})$ returns $R(\bar{a})$.

Proof. Let $R'(\bar{t})$ be an answer to Q which has a Datalog proof-tree (G, ℓ) in $P(\mathfrak{J})$ containing a label $R(\bar{a})$. We prove by induction upon $k < \omega$, that the atom $R(\bar{a})$ is yielded after the k -th repeat-loop pass by every call $\text{infer}(Q', \text{PrevQueries})$, if Q' unifies with $R(\bar{a})$ and the minimal height of some Datalog proof tree for $R(\bar{a})$ is equal to k .

If there is a Datalog proof tree for $R(\bar{a})$ of height 0, and Q' unifies with $R(\bar{a})$ then $\text{infer}(Q', \text{PrevQueries})$ will produce $R(\bar{a})$ for all $\text{PrevQueries} \subseteq \mathcal{Q}$ in the look-up of line 13 which will be returned (cf. line 26) for all further repeat-loop passes.

Assume there is a Datalog proof tree (G', ℓ') for $R(\bar{a})$ of height $k+1$ and we have started the $k+1$ repeat-loop pass. Lemma 4 shows that there is subsequent procedure call $\text{infer}(Q_n, \text{PrevQueries})$ of $\text{infer}(Q, \emptyset)$ such that Q_n is an unblocked atomic query, i.e. there is no $Q' \in \text{PrevQueries}$ with $Q_n \sqsubseteq Q'$ and $Q' \sqsubseteq Q_n$, and $R(\bar{a})$ unifies with Q_n .

Since (G', ℓ') is of height $k+1$ there is a rule $r : R(\bar{t}) \leftarrow R_1(\bar{t}_1) \wedge \dots \wedge R_m(\bar{t}_m)$ and a variable assignment β such that $R(\beta(\bar{t})) = R(\bar{a})$ and for each $i \in \{1, \dots, m\}$ the fact $R_i(\beta(\bar{t}_i))$ has a proof tree of height at most k under P in \mathfrak{J} .

The head H of r and Q_n unify with the ground atom $R(\bar{a})$, so there is $\theta := \text{MGU}(Q_n, H)$ and each $R_i(\beta(\bar{t}_i))$ unifies with $Q'_i := R_i(\theta(\bar{t}_i))$ with $i \in \{1, \dots, m\}$. Since Q_n is not blocked, the subsequent procedure call $\text{infer}(Q'_i, \text{PrevQueries} \cup \{Q_n\})$ for all $i \in \{1, \dots, m\}$ is issued.

By the induction hypothesis, for all $i \in \{1, \dots, m\}$, $\text{infer}(Q'_i, \text{PrevQueries} \cup \{Q_n\})$ yields $R_i(\beta(\bar{t}_i))$. Hence $R(\bar{a})$ is returned by $\text{infer}(Q_n, \text{PrevQueries})$ at the very latest in the $n+1$ repeat-loop pass and eventually added to Mat (cf. line 4) so that after the $n+1$ repeat-loop pass for every $\text{PrevQueries} \subseteq \mathcal{Q}$ every subsequent call $\text{infer}(Q', \text{PrevQueries})$ that unifies with $R(\bar{a})$ will return this fact as look-up in line 13. \square

By proving Proposition 1, we have shown that Algorithm 1 is complete: Every answer to a given query Q has a Datalog proof-tree and is the label of the root of this proof-tree. Proposition 1 shows that this answer is eventually derived by $\text{infer}(Q, \emptyset)$ (line 5) and thus

returned by $main(Q)$ (line 7). This completes the three steps to prove correctness w.r.t. to retrieval of exactly those answers to Q under $P(\mathcal{I})$.

However, in our approach this algorithm is invoked with a rule set that is different from the one that is given. Therefore, we still need to prove that our entire approach is complete in a sense that an execution of this algorithm with the modified rule set retrieves the same results as an execution that uses the original rule set (provided that a pre-materialization is performed beforehand). This issue will be discussed in Section 4.

In the following section, we will conclude the discussion of our backward-chaining algorithm by describing how a generic set of the precalculated predicates can be used in the implementation to speed up the performance of reasoning at query time.

3.4. On the implementation of backward-chaining with pre-materialized queries

So far, we made no difference in the description of our backward-chaining algorithm between subqueries that are precomputed and those which are not. However, the pre-materialization of a selection of queries allows us to substantially improve the implementation and performance of backward-chaining by exploiting the fact that these queries can be retrieved with a single lookup. In our implementation, these queries are stored in main memory so that the joins required by the rules can be efficiently executed.

Also, the availability of the pre-materialized queries in memory allows us to implement another efficient information passing strategy to reduce the size of the proof tree by identifying beforehand whether a rule can contribute to derive answers for a given query.

In fact, the pre-materialization can be used to determine early failures: Emptiness for queries which are subsumed by the pre-materialized queries can be cheaply derived since a lookup suffices. Therefore, when scheduling the derivation of rule body atoms, we give priority to those body atoms that potentially match these pre-materialized queries so that if these “cheap” body atoms do not yield any answers, the rule will not apply, and we can avoid the computation of the more expensive body atoms of the rule for which further reasoning would have been required.

To better illustrate this concept, we proceed with an example. Suppose we have the proof tree described in Figure 1. In this case, the reasoner can potentially apply rule `prp-symp` (concerning symmetric properties

in OWL) to derive some triples that are part of the second body atom of rule `prp-spol`.

However, in this case, rule `prp-symp` will fire only if some of the subjects (i.e. the first component) of the triples part of $T(x, SPO, TYPE)$ will also be the subject of $T(x, TYPE, SYM)$. If both patterns are pre-computed, then we know beforehand all the possible ‘ x ’, and therefore we can immediately perform an intersection between the two sets. If the intersection is non-empty, the reasoner proceeds executing rule `prp-symp`, otherwise it can skip its execution since the rule will never fire.

It is very unlikely that the same property appears in all the terminological patterns, therefore an information passing strategy that is based on the pre-computed triple patterns is very effective in significantly reducing the tree size and hence improving performance.

Furthermore, our implementation of this algorithm applies a *sideways information passing* strategy [3,2] to improve the execution of the joins required by the rules. This technique consists of “passing” admissible values to the variables of the following queries in order to limit the retrieval to only facts that can contribute to the join.

To illustrate this concept, consider Example 5: Here, when the algorithm needs to invoke the function *infer* with the query $(a, TYPE, w)$, even though the variable w could in principle assume any value, in practice the implementation knows already that only the values of w in *subst* are admissible (in our case c), because only those can lead to a successful join. Thus, the implementation can link these values to the variable w so that in subsequent call of, for example, *lookup*, the computation is limited to retrieve only these values and not all possible w . This technique of passing values between the queries is well-known in rule-based systems, and applied in nearly all implementations.

In the following section, we will focus on the pre-materialization phase, whose purpose is to calculate the results for these subqueries in advance.

4. Hybrid Reasoning: Pre-Materialization

In our approach, we use a rule set that is different from standard OWL. Our approach thus relies on a pre-materialization phase to be performed before the user can query the knowledge base. The purpose of this pre-materialization is to guarantee that whenever the user queries the knowledge base our backward-chaining al-

gorithm is able to infer the same answers as with the original rules.

In order to prove the correctness of our method for hybrid reasoning, we proceed as follows: First, in Section 4.1, we formalize and discuss the completeness of the pre-materialization algorithm, which consists in the ability of calculating all the results for some pre-defined queries and store them in some additional *edb* relations.

Next, in Section 4.2, we show that replacing the pre-materialized body atoms in the original rules with atoms using the introduced edbs is harmless (after the pre-materialization algorithm is computed), since this modified program computes the same answers as the original one. This last argument finally settles the correctness of our approach, since it ensures that, after the pre-materialization is computed, no derivation is missed.

4.1. Pre-Materialization

Let \mathcal{J} be a database and P the program with a set L of atomic queries that are selected for pre-materialization. The pre-materialization is performed by Algorithm 2. The reason why we do not simply introduce auxiliary relations named S_Q to \mathcal{J} for each $Q \in L$ and populate these by setting $S_Q^{\mathcal{J}} := \text{main}(Q)$ (for *main* cf. Algorithm 1) is that the efficiency of Algorithm 1 hinges upon the fact that as many body atoms as possible are not unfoldable, but are edbs for which merely look-ups have to be performed during backward chaining.

We explain now Algorithm 2 in detail and discuss its completeness: In a first step (lines 1–3), the database is extended with auxiliary relations named S_Q for $Q \in L$. Each rule of the program P is rewritten (lines 5–12) by replacing every body atom $R_i(\bar{t}_i)$ with the atom $S_Q(\bar{t}_i)$ if $R_i(\bar{t}_i) \sqsubseteq Q$ (cf. page 4), i.e. if the “answers” to $R_i(\bar{t}_i)$ are also yielded by Q .

Clearly, the result of the rewriting need not to be deterministic in case there are two or more atomic queries $Q_0, Q_1 \in L$ with $R_i(\bar{t}_i) \sqsubseteq Q_0, Q_1$. However, we shall show that either rewriting is good enough. The new rule thus obtained is stored in a new program P' . In case the rule p contains no body atoms that need to be replaced, p is stored in P' as well.

In each repeat-loop pass (cf. lines 15–24), \mathcal{J} is extended in an external step (lines 17–19) with all answers for $Q \in L$, which are copied into the auxiliary relation $S_Q^{\mathcal{J}}$. Since this is repeated between each derivation until no new answers for any $Q \in L$ are

yielded, this is equivalent to adding $S_Q(\bar{t}) \leftarrow R(\bar{t})$ for each $Q \in L$ with $Q = R(\bar{t})$ to P' directly.⁴ One can derive from this argument that Algorithm 2 terminates and is sound in the sense that $S_Q^{\mathcal{J}_0} \subseteq Q^{P(\mathcal{J})}$, where \mathcal{J}_0 the database \mathcal{J} after Algorithm 2 has terminated. As we shall show in Proposition 2, Algorithm 2 is complete in the sense that, after termination of this algorithm, $S_Q^{\mathcal{J}_0}$ contains all answers for Q in $P(\mathcal{J})$.

Example 6. Take the altered program from Example 1 and add the appropriate $S_Q(x, \text{SPO}, y) \leftarrow Q$ with $Q = T(x, \text{SPO}, y)$ to it. In this case we obtain

$$\begin{aligned} T(a, p1, b) &\leftarrow S_Q(p, \text{SPO}, p1) \wedge T(a, p, b) \\ T(x, \text{SPO}, y) &\leftarrow S_Q(x, \text{SPO}, w) \wedge S_Q(w, \text{SPO}, y) \\ S_Q(x, \text{SPO}, y) &\leftarrow T(x, \text{SPO}, y) \end{aligned}$$

It is trivially clear, that this program yields for every Database exactly the same results for $T(x, \text{SPO}, y)$ as the original program

$$\begin{aligned} T(a, p1, b) &\leftarrow T(p, \text{SPO}, p1) \wedge T(a, p, b) \\ T(x, \text{SPO}, y) &\leftarrow T(x, \text{SPO}, w) \wedge T(w, \text{SPO}, y) \end{aligned}$$

□

Proposition 2. *Algorithm 2 is complete in the sense that for the database \mathcal{J}_0 which we obtain after Algorithm 2 has terminated $S_Q^{\mathcal{J}_0} \supseteq Q^{P(\mathcal{J})}$ for all $Q \in L$, i.e. every answer that could be derived from Q under P in \mathcal{J} is contained in $S_Q^{\mathcal{J}_0}$.*

Proof. Let P' be the rewritten program P , defined as *NewRuleset* in the pseudocode of Algorithm 2. We have $Q^{P(\mathcal{J})} \subseteq Q^{P(\mathcal{J}_0)}$ (monotonicity) and $Q^{P'(\mathcal{J}_0)} \subseteq S_Q^{\mathcal{J}_0}$ (line 18). In order to show $Q^{P(\mathcal{J})} \subseteq S_Q^{\mathcal{J}_0}$ we show $Q^{P(\mathcal{J}_0)} \subseteq Q^{P'(\mathcal{J}_0)}$.

Assume no new element could be derived from \mathcal{J}_0 using the program P' but for some $Q \in L$, $\text{main}(Q)$ could derive another yet unknown ground atom from \mathcal{J}_0 using the original program P . Let hence $\Delta := \bigcup_{Q \in L} Q^{P(\mathcal{J}_0)} \setminus Q^{P'(\mathcal{J}_0)}$ and let $R(\bar{a}) \in \Delta$ such that it has a Datalog proof-tree (G, ℓ) in $P(\mathcal{J}_0)$ of height m , where $m < \omega$ is for all atoms from Δ the least height of Datalog proof-trees in $P(\mathcal{J}_0)$.

We change the tree (G, ℓ) recursively as follows: If the root v is a leaf, the tree stays unaltered. Otherwise, there is a rule $r \in P$ and an assignment β

⁴By definition, this would render $S_Q(\bar{t})$ into an idb, which we thus only propose for the sake of explaining correctness.

such that there is a bijection ι between the successor set $E(v)$ and the set of body atoms of r such that $\ell(v') = \beta \circ \iota(v')$ for all $v' \in E(v)$. As we only exchanged predicate names in the rewriting of $r \in P$ to $r' \in P'$, there is a bijection ι' between $E(v)$ and the body atoms of r' and we set $\ell'(v') := \beta \circ \iota'(v')$. For all $v' \in E(v)$ we have $\ell(v') = \ell'(v')$ if the body atom was not replaced during rewriting. If for any $v' \in E(v)$ $\ell'(v')$ has predicate name S_Q for some $Q \in L$, prune its subtree but keep v' . Otherwise recursively continue to change the subtree $\langle v' \rangle$.

Line 18 guarantees that $S_Q^{\mathcal{J}_0} = Q^{\mathcal{J}_0}$ for all $Q \in L$ and since the (G, ℓ) was chosen minimal, every leaf of (G, ℓ') is labelled with an atom in \mathcal{J}_0 . Hence (G, ℓ') is a Datalog proof-tree in $P'(\mathcal{J}_0)$ for $R(\bar{a})$ and so $R(\bar{a}) \in \bigcup_{Q \in L} Q^{P'(\mathcal{J}_0)} \setminus Q^{P(\mathcal{J}_0)}$. A contradiction! \square

Algorithm 2 Overall algorithm of the precomputation procedure: L is a set containing all queries that were selected for pre-materialization, $RuleSet$ is a constant containing a program P and $Database$ represents \mathcal{J} .

```

1  for every  $Q \in L$ 
2    introduce a new predicate symbol  $S_Q$  to
      Database
3  end for
4
5  for every rule  $p : R_0(\bar{t}_0) \leftarrow R_1(\bar{t}_1) \wedge \dots \wedge R_n(\bar{t}_n)$  in
      RuleSet
6    for every  $Q \in L$ 
7      if  $R_i(\bar{t}_i) \sqsubseteq Q$  then
8        replace  $R_i(\bar{t}_i)$  in  $p$  with  $S_Q(\bar{t}_i)$ 
9      end if
10   end for
11   add this (altered) rule to  $NewRuleSet$ 
12 end for
13
14 Derivation :=  $\emptyset$ 
15 repeat
16   Database := Database  $\cup$  Derivation
17   for every  $R(\bar{t}) \in L$ 
18     Perform  $S_{R(\bar{t})}(\bar{t}) \leftarrow R(\bar{t})$  on Database
19   end for
20
21   for every  $Q$  in  $L$ 
22     Derivation := Derivation  $\cup$  main(Q) using
      NewRuleSet as program on Database
23   end for
24 until Derivation  $\subseteq$  Database

```

4.2. Reasoning with Pre-Materialized Predicates

We show now that replacing body atoms with auxiliary predicates that contain the full materialization of the body atom w.r.t. a given database, yields the same full materialization of the database as under the origi-

nal program. We will formally define what conditions must be fulfilled and prove that if they hold, the two programs will produce the same derivation. Finally, we point out that the output of Algorithm 2 satisfies these condition and therefore guarantees the correctness of the approach.

We start by letting P be an arbitrary Datalog program and \mathcal{J} a database. We assume that \mathcal{J} has already been enriched with the results of the pre-materialization.

As an example, assume the binary relation $S^{\mathcal{J}}$ contains all answer tuples of the query

$$query(x, y) \leftarrow T(x, SPO, y).$$

under the program P .

From an abstract point of view we can define S as an *extensional database predicate* (edb) of P , i.e. it is not altered by P so that the *interpretation* $S^{\mathcal{J}}$ of S under \mathcal{J} equals the interpretation $S^{P(\mathcal{J})}$ of S under the *least fix-point model* $P(\mathcal{J})$ of P over \mathcal{J} .

Since S is an edb and therefore does not appear in the head of any rule of P , S cannot be unfolded and so the evaluation of S during the backward chaining process is reduced to a mere look-up in the database.

Such a replacement is harmless only if \mathcal{J} has been adequately enriched. Thus, the question arises which abstract conditions must be satisfied to allow such a replacement: In essence, we want that a rule fires under “almost the same” variable assignment as its replacement, which we formalize in the following two paragraphs.

Let $R(\bar{t}_0) \leftarrow R_1(\bar{t}_1) \wedge \dots \wedge R_n(\bar{t}_n)$ be a rule in P . We define two queries, one being the body of the rule and one being the body of the rule where for some $i \in \{1, \dots, n\}$ the body atom $R_i(\bar{t}_i)$ is replaced by $S(\bar{t})$ where \bar{t} is some arbitrary tuple having the arity of S . Let $\bar{z} := \bar{t}_0 \cdot \bar{t}_1 \cdots \bar{t}_{i-1} \cdot \bar{t}_{i+1} \cdots \bar{t}_n$, i.e. the concatenation of all tuples except \bar{t}_i and

$$\begin{aligned} q_0(\bar{z}) &\leftarrow R_1(\bar{t}_1) \wedge \dots \wedge R_i(\bar{t}_i) \wedge \dots \wedge R_n(\bar{t}_n) \\ q_1(\bar{z}) &\leftarrow R_1(\bar{t}_1) \wedge \dots \wedge S(\bar{t}) \wedge \dots \wedge R_n(\bar{t}_n) \quad (*) \end{aligned}$$

Now, the rule and its replacement fire under *almost the same variable assignment* iff $q_0(\bar{z})^{P(\mathcal{J})} = q_1(\bar{z})^{P(\mathcal{J})}$, i.e. q_0 and q_1 yield the same answers under P in \mathcal{J} . We see, that it is “almost the same” variable assignment, as we do not require variable assignments to coincide on \bar{t}_i and \bar{t} . In this way we do not require, e.g., $S^{\mathcal{J}} = R_i^{P(\mathcal{J})}$. S is merely required to contain the *necessary* information. This is important, if we want

to apply the substitution to rdf-triples, where we lack distinguished predicate names:

Example 7. Since there is only one generic predicate symbol T , requiring $S^{\mathcal{J}} = T^{P(\mathcal{J})}$ would mean that S contains the complete materialization of \mathcal{J} under P which would render our approach obsolete. \square

Also note that it is not sufficient to merely require $q_0(\bar{t}_0)^{P(\mathcal{J})} = q_1(\bar{t}_0)^{P(\mathcal{J})}$, i.e. that both queries yield the same answer tuples \bar{t}_0 under $P(\mathcal{J})$, as the following example shows.

Example 8. Let the program P which computes the transitive closure of R_0 in R_1 consist of the two rules:

$$\begin{aligned} R_1(x, z) &\leftarrow R_1(x, y) \wedge R_0(y, z) \\ R_1(x, y) &\leftarrow R_0(x, y) \end{aligned}$$

Consider a database \mathcal{J} with $R_0^{\mathcal{J}} := \{(a, b), (b, c), (b, b), (c, c)\}$. In the least fix-point model $P(\mathcal{J})$ of P we expect $R_1^{P(\mathcal{J})} = \{(a, b), (b, c), (a, c), (b, b), (c, c)\}$. Let S have the interpretation $S^{\mathcal{J}} = \{(b, b), (c, c)\}$. Since $R_1^{P(\mathcal{J})}$ is the transitive closure, the following two queries deliver the same answer tuples under $P(\mathcal{J})$, i.e.

$$\begin{aligned} q_0(x, z) &\leftarrow R_1(x, y) \wedge R_0(y, z) \\ q_1(x, z) &\leftarrow R_1(x, y) \wedge S(y, z) \end{aligned}$$

Yet the program P'

$$\begin{aligned} R_1(x, z) &\leftarrow R_1(x, y) \wedge S(y, z) \\ R_1(x, y) &\leftarrow R_0(x, y) \end{aligned}$$

will *not* compute the transitive closure of R_0 in R_1 , as $R_1^{P'(\mathcal{J})} = \{(a, b), (b, c), (b, b), (c, c)\}$. \square

We show now that substituting a body atom $R_i(\bar{t}_i)$ by $S(\bar{t})$ under the condition that the queries in (*) yield the same answer tuples under $P(\mathcal{J})$, generates the same least fix-point:

Proposition 3. *Let P' be the program P where the rule*

$R_0(\bar{t}_0) \leftarrow R_1(\bar{t}_1) \wedge \dots \wedge R_i(\bar{t}_i) \wedge \dots \wedge R_n(\bar{t}_n) \in P$ has, for some tuple \bar{t} and edb S , been replaced by

$R_0(\bar{t}_0) \leftarrow R_1(\bar{t}_1) \wedge \dots \wedge S(\bar{t}) \wedge \dots \wedge R_n(\bar{t}_n)$.

Let q_0 and q_1 be defined as in ().*

If $q_0(\bar{z})^{P(\mathcal{J})} = q_1(\bar{z})^{P(\mathcal{J})}$ then $P(\mathcal{J}) = P'(\mathcal{J})$.

Proof. In order to show the implication we assume $q_0(\bar{z})^{P(\mathcal{J})} = q_1(\bar{z})^{P(\mathcal{J})}$. Let T_P and $T_{P'}$ be the im-

mediate consequence operators (mentioned on page 5) for each program. We show for all $k < \omega$ and every ground atom $R(\bar{a})$ that if $R(\bar{a}) \in T_P^k(\mathcal{J})$ then there is an $\ell < \omega$ such that $R(\bar{a}) \in T_{P'}^\ell(\mathcal{J})$ and vice versa. Since we start out from the same database \mathcal{J} we have $T_P^0(\mathcal{J}) = T_{P'}^0(\mathcal{J})$ which settles the base case.

Let $R(\bar{a}) \in T_P^{k+1}(\mathcal{J})$ then either $R(\bar{a}) \in T_P^0(\mathcal{J})$ and we are done or there is some rule $R(\bar{t}_0) \leftarrow R_1(\bar{t}_1) \wedge \dots \wedge R_n(\bar{t}_n)$ and some variable assignment β such that $\beta(\bar{t}_0) = \bar{a}$ and $R_j(\beta(\bar{t}_j)) \in T_P^k(\mathcal{J})$ for all $j \in \{1, \dots, n\}$.

If $R(\bar{t}_0) \leftarrow R_1(\bar{t}_1) \wedge \dots \wedge R_n(\bar{t}_n) \in P'$, i.e. none of its body atoms were substituted, the induction hypothesis shows for each $j \in \{1, \dots, n\}$ that we can find $\ell_j < \omega$ such that $R_j(\beta(\bar{t}_j)) \in T_{P'}^{\ell_j}(\mathcal{J})$. Let $\ell_0 := \max(\{0\} \cup \{\ell_j \mid 1 \leq j \leq n\})$. Note that we add $\{0\}$ for the case where the rule body is empty. In any case, we have $R_j(\beta(\bar{t}_j)) \in T_{P'}^{\ell_0}(\mathcal{J})$ for all $j \in \{1, \dots, n\}$. Since all premises of this rule are satisfied, there is some $\ell := \ell_0 + 1$ such that $R_0(\beta(\bar{t}_0)) \in T_{P'}^\ell(\mathcal{J})$.

If $R(\bar{t}_0) \leftarrow R_1(\bar{t}_1) \wedge \dots \wedge R_n(\bar{t}_n) \notin P'$ it is the rule where $R_i(\bar{t}_i)$ has been substituted with $S(\bar{t})$. For the assignment β we now know $\beta(\bar{t}_0 \cdot \bar{t}_1 \cdot \dots \cdot \bar{t}_{i-1} \cdot \bar{t}_{i+1} \cdot \dots \cdot \bar{t}_n) \in q_0(\bar{z})^{P(\mathcal{J})}$. Since $q_0(\bar{z})^{P(\mathcal{J})} = q_1(\bar{z})^{P(\mathcal{J})}$ we know that there is some assignment β' , which coincides with β on $(\bar{t}_0 \cdot \bar{t}_1 \cdot \dots \cdot \bar{t}_{i-1} \cdot \bar{t}_{i+1} \cdot \dots \cdot \bar{t}_n)$ such that $\beta'(\bar{t}) \in S^{P(\mathcal{J})}$.

Hence $R_j(\beta'(\bar{t}_j)) \in T_P^k(\mathcal{J})$ for all $j \in \{1, \dots, n\} \setminus \{i\}$ and $S(\beta'(\bar{t})) \in T_P^0(\mathcal{J})$ since S is an edb predicate. The induction hypothesis yields some $\ell_j < \omega$ for each $j \in \{1, \dots, n\} \setminus \{i\}$ such that $R_j(\beta'(\bar{t}_j)) \in T_{P'}^{\ell_j}(\mathcal{J})$. Let $\ell_0 := \max(\{0\} \cup \{\ell_j \mid 1 \leq j \leq n \text{ and } j \neq i\})$, then $R_j(\beta'(\bar{t}_j)) \in T_{P'}^{\ell_0}(\mathcal{J})$ for all $j \in \{1, \dots, n\} \setminus \{i\}$ and $S(\beta'(\bar{t})) \in T_{P'}^0(\mathcal{J})$. Since all premises of this rule are satisfied, there is some $\ell := \ell_0 + 1$ such that $R_0(\beta'(\bar{t}_0)) \in T_{P'}^\ell(\mathcal{J})$. As β coincides with β' also on \bar{t}_0 , i.e. $\beta'(\bar{t}_0) = \bar{a}$, we have in particular $R(\bar{a}) \in T_{P'}^\ell(\mathcal{J})$.

This shows $R^{P(\mathcal{J})} \subseteq R^{P'(\mathcal{J})}$ for all predicate names $R \in \text{SIG}$. For the converse, we merely show the case of the substituted rule: Assume $R(\bar{a}) \in T_{P'}^{k+1}(\mathcal{J})$ and there is an assignment β' such that $\beta'(\bar{t}_0) = \bar{a}$ and $R_j(\beta'(\bar{t}_j)) \in T_{P'}^k(\mathcal{J})$ for all $j \in \{1, \dots, n\} \setminus \{i\}$ as well as $\beta'(\bar{t}) \in S^{P'(\mathcal{J})}$.

The induction hypothesis yields for each $j \in \{1, \dots, n\} \setminus \{i\}$ some $\ell_j < \omega$ with $R_j(\beta'(\bar{t}_j)) \in T_{P'}^{\ell_j}(\mathcal{J})$. Since S is an edb predicate for P , we have $S(\beta'(\bar{t})) \in T_P^0(\mathcal{J})$. Hence for $\ell_0 := \max(\{0\} \cup \{\ell_j \mid$

$1 \leq j \leq n$ and $j \neq i$) we have $R_j(\beta'(\bar{t}_j)) \in T_P^{\ell_0}(\mathcal{J})$ for all $j \in \{1, \dots, n\} \setminus \{i\}$ and $S(\beta'(\bar{t})) \in T_P^0(\mathcal{J})$.

This implies $\beta'(\bar{t}_0 \cdot \bar{t}_1 \cdots \bar{t}_{i-1} \cdot \bar{t}_{i+1} \cdots \bar{t}_n) \in q_1(\bar{z})^{P(\mathcal{J})}$ and since $q_0(\bar{z})^{P(\mathcal{J})} = q_1(\bar{z})^{P(\mathcal{J})}$ there is an assignment β coinciding on $(\bar{t}_0 \cdot \bar{t}_1 \cdots \bar{t}_{i-1} \cdot \bar{t}_{i+1} \cdots \bar{t}_n)$ with β' such that $R_i(\beta(\bar{t}_i)) \in T_P^{j_0}(\mathcal{J})$ for some $j_0 < \omega$. Let $\ell_1 := \max\{\ell_0, j_0\}$ then $R_j(\beta'(\bar{t}_j)) \in T_P^{\ell_1}(\mathcal{J})$ for all $j \in \{1, \dots, n\}$. Since all premises of the rule $R_0(\bar{t}_0) \leftarrow R_1(\bar{t}_1) \wedge \dots \wedge R_n(\bar{t}_n)$ are satisfied, there is some $\ell := \ell_1 + 1$ such that $R_0(\beta(\bar{t}_0)) \in T_P^\ell(\mathcal{J})$, which shows, as β coincides on t_0 with β' that $R(\bar{a}) \in T_P^\ell(\mathcal{J})$.

Together with $R^{P(\mathcal{J})} \subseteq R^{P'(\mathcal{J})}$ this shows $R^{P(\mathcal{J})} = R^{P'(\mathcal{J})}$ for all predicate names $R \in \text{SIG}$ and hence that $P(\mathcal{J}) = P'(\mathcal{J})$. \square

It now becomes clear, how Algorithm 2 and Proposition 3 fit together: For a given database \mathcal{J} and a set of atomic queries L , Algorithm 2 computes for each $Q \in L$ the query answers under the program P , which are stored in the relation $S_Q^{\mathcal{J}_0}$, where \mathcal{J}_0 is \mathcal{J} after Algorithm 2 has finished. These S_Q are edbs for P .

Let now $r : R_0(\bar{t}_0) \leftarrow R_1(\bar{t}_1) \wedge \dots \wedge R_n(\bar{t}_n)$ be a rule in this program and $Q \in L$ an atomic query s.t. $R_i(\bar{t}_i) \sqsubseteq Q$, then $Q = R_i(\bar{t})$ such that $\bar{t}_i \sqsubseteq \bar{t}$ by definition of \sqsubseteq . Correctness of Algorithm 2 yields $R(\bar{t}_i)^{P(\mathcal{J}_0)} = S_Q(\bar{t}_i)^{\mathcal{J}_0}$ and hence that $q_0(\bar{z})^{P(\mathcal{J}_0)} = q_1(\bar{z})^{P(\mathcal{J}_0)}$ where $\bar{z} = \bar{t}_0 \cdots \bar{t}_n$ and

$$\begin{aligned} q_0(\bar{z}) &\leftarrow R_1(\bar{t}_1) \wedge \dots \wedge R_i(\bar{t}_i) \wedge \dots \wedge R_n(\bar{t}_n) \\ q_1(\bar{z}) &\leftarrow R_1(\bar{t}_1) \wedge \dots \wedge S_Q(\bar{t}_i) \wedge \dots \wedge R_n(\bar{t}_n) \end{aligned}$$

Proposition 3 guarantees that the substitution of $R_i(\bar{t}_i)$ by $S_Q(\bar{t}_i)$ in rule r is harmless w.r.t. \mathcal{J} . By applying this argument iteratively, one eventually obtains a program P' in which all precomputed atoms have been replaced and which yields the same materialization for \mathcal{J}_0 as P .

In the following section, we apply this rewriting to the OWL RL rule set.

5. Hybrid reasoning for OWL RL

In the previous sections, we described the two main components of our method which consist of the backward-chaining algorithm used to retrieve the inference and the pre-materialization procedure which ensures the completeness of our approach.

We now discuss the implementation of the OWL RL rules using our approach. In fact, while our method can

OWL RL	RDFS	(?X rdfs:subPropertyOf ?Y) (?X rdfs:subClassOf ?Y) (?X rdfs:domain ?Y) (?X rdfs:range ?Y)
	PID*	(?P rdf:type owl:FunctionalProperty) (?X owl:allValuesFrom ?Y) (?P rdf:type owl:InverseFunctionalProperty) (?X owl:inverseOf ?Y) (?P rdf:type owl:TransitiveProperty) (?X rdf:type owl:Class) (?P rdf:type owl:SymmetricProperty) (?X rdf:type owl:Property) (?X owl:equivalentClass ?Y) (?X owl:onProperty ?Y) (?X owl:hasValue ?Y) (?X owl:equivalentProperty ?Y) (?X owl:someValuesFrom ?Y)
		(?X owl:propertyChainAxiom ?Y) (?X owl:hasKey ?Y) (?X owl:intersectionOf ?Y) (?X owl:unionOf ?Y) (?X owl:oneOf ?Y) (?X owl:maxCardinality 1) (?X owl:maxQualifiedCardinality 1) (?X owl:onClass ?Y) (?X rdf:type owl:Class) (?X rdf:type owl:DatatypeProp) (?X rdf:type owl:ObjectTypeProp)

Table 2

Triple patterns that are precalculated considering the OWL RL rules (1 is an abbreviation for the URI "1"^^xsd:nonNegativeInteger).

in principle be applied to any Datalog program, our implementation heavily relies on the fact that our program consists of inference rules on RDF data and thus, our prototype is unable to execute generic Datalog programs.

The official OWL RL rule set contains 78 rules, for which the reader is referred to the official document overview [13]. With some selected examples from [13] we illustrate some key features of our algorithm.

Initial assumptions. First of all, we exclude some rules from our discussion and implementation for various reasons. These are:

- All the rules whose purpose is to derive an inconsistency, i.e. rules with predicate *false* in the head of the rule. We do not consider them because our purpose is to derive new triples rather than identify an inconsistency;
- All the rules which have an empty body. These rules cannot be triggered during the unfolding process of backward-chaining. These rules include those that encode the semantics of datatypes,

therefore our implementation does not support datatypes;

- The rules that exploit the *owl:sameAs* transitivity and symmetry.⁵ These rules require a computation that is too expensive to perform at query time since they can be virtually applied to every single term of the triples. These rules can be implemented by computing the *sameAs* closure and maintaining a consolidation table.⁶

In our approach, we decided to pre-materialize all triple patterns that are used to retrieve “schema” triples, also referred to as the terminological triples. Table 2 reports the set of the patterns that are pre-materialized using our method described in Section 4.

Singling out exactly those triple patterns from Table 2 is motivated by the empirical observation that:

- they appear in many of the OWL rules;
- their answer sets are very small compared to the entire input;
- their answer sets are not as frequently updated as the rest of the data.

These characteristics make the set of inferred schema triples the ideal candidate to be pre-materialized. All rules which have a pre-materialized pattern amongst their body atoms are substituted replacing the pre-materialized pattern with its corresponding auxiliary relation as justified by Proposition 3. This affects 25 rules out of 78 and hence reduces reasoning considerably.

After the pre-materialization procedure is completed, each rule which has a pre-materialized pattern in its head can be reduced to a mere look-up:

Example 9. Consider (scm-sco) from Table 9 in [13]:

$$T(x, \text{SCO}, z) \leftarrow T(x, \text{SCO}, y) \wedge T(y, \text{SCO}, z)$$

can be replaced according to Proposition 3 by r' :

$$T(x, \text{SCO}, z) \leftarrow S_{\text{sco}}(x, \text{SCO}, y) \wedge S_{\text{sco}}(y, \text{SCO}, z)$$

where the answers to $T(x, \text{SCO}, y)$ under \mathcal{I} are contained in $S_{\text{sco}}^{\mathcal{I}}$. By adding the rule

$$r'': T(x, \text{SCO}, z) \leftarrow S_{\text{sco}}(x, \text{SCO}, z)$$

every rule whose head atom unifies with $T(x, \text{SCO}, y)$ can be deleted. Adding r'' is harmless; we can render r' into r'' using Proposition 3: Since S_{sco} is transitive

we may replace an imaginary conjunction with \top (the atom which is always true) in r' with $S_{\text{sco}}(x, \text{SCO}, z)$. Using Proposition 3 twice more, we can successively replace $S_{\text{sco}}(x, \text{SCO}, y)$ and $S_{\text{sco}}(y, \text{SCO}, z)$ each by \top obtaining r'' . \square

This removes a further 30 rules from unfolding.

On the implementation of RDF lists. Some of the inference rules in the OWL RL rule set use RDF lists to define a variable number of antecedents. The RDF lists cannot be represented in Datalog in a straightforward way since they rely on *rdf:first* and *rdf:next* triples to represent the elements of the list. Therefore, they need to be processed differently.

In our implementation, at every step of the pre-materialization procedure, we launch two additional queries to retrieve all the (inferred and explicit) *rdf:first* and *rdf:rest* triples with the purpose of constructing such lists. Once we have collected them, we perform a join with the other schema triples, and determine the sequence of elements by repeatedly joining the *rdf:first* and *rdf:rest* triples. After this operation is completed, from the point of view of the Datalog program, RDF lists appear as simple list of elements and are used according to the various rule logics. For example, if the rule requires a matching of all the elements of the list with the other antecedents, then the rule executor will first use the first element to perform the join, then it will use the second, and so on, until all the elements are considered.

5.1. Detecting duplicate derivation in OWL RL

Since the OWL RL fragment consists of a large number of rules, there is a high possibility that the proof tree contains branches that lead to the same derivation. Detecting and avoiding the execution of these branches is essential in order to reduce the computation.

After empirically analyzing some example queries, we determined that there are two types of sources in the generation of duplicates. The first comes from the nature of the rule set. The second comes from the input data.

First type of duplicates source. The most prominent example of generation of duplicates of the first type is represented by the *symmetric* rules which have the same structure but have the variables positioned at different locations. We refer with rule names and tables in the following list to the OWL RL rule set in [13]:

⁵The list of these rules is reported in Table 4 of [13].

⁶This procedure is explained in detail in [19].

prp-eqp1 and prp-eqp2 from Table 5
 cax-eqc1 and cax-eqc2 from Table 7
 prp-inv1 and prp-inv2 from Table 5.

We analyze each of these three cases below.

Let $S_{eqp}^{\mathcal{J}}$ be the pre-materialization of the triple pattern $T(?x, EQP, ?y)$. The rules scm-eqp1 and scm-eqp2 render $S_{eqp}^{\mathcal{J}}$ symmetric. Hence

$$q(x, p_2, y) \leftarrow T(x, p_1, y) \wedge S_{eqp}(p_1, EQP, p_2)$$

yields the same results under $P(\mathcal{J})$ as

$$q(x, p_1, y) \leftarrow T(x, p_2, y) \wedge S_{eqp}(p_2, EQP, p_1)$$

and Proposition 3 yields that scm-eqp2 can be replaced by scm-eqp1, effectively deleting scm-eqp2 from the rule set.

Similarly, rules scm-eqc1 and scm-eqc2 render the results of the pre-materialized query $T(?x, EQC, ?y)$, symmetric.

In contrast, the pre-materialized query $T(?x, INV, ?y)$ is not symmetric. However, we can first observe that Proposition 3 allows us to replace the rules by

$$\begin{aligned} T(x, p, y) &\leftarrow T(x, q, y) \wedge S_{inv}(p, INV, q) \\ T(x, p, y) &\leftarrow T(x, q, y) \wedge S_{inv}(q, INV, p) \end{aligned}$$

which defuses the idb atom $T(q, INV, p)$ into the harmless edb S_{inv} , for which $S_{inv}^{\mathcal{J}}$ contains all answers to the query $T(?x, INV, ?y)$. Let this new program be called P' . Further, let P'' be the program where both rules have been replaced by

$$T(x, p, y) \leftarrow T(x, q, y) \wedge S'_{inv}(p, INV, q)$$

with $S'_{inv}^{\mathcal{J}}$ being the symmetric closure of $S_{inv}^{\mathcal{J}}$. It is now not difficult to see, that every model of P' is a model of P'' and vice versa. In particular the least fix-point model of $P(\mathcal{J})$ is equal to the least fix-point model $P''(\mathcal{J})$. Hence we can replace prp-inv1 and prp-inv2 by one rule under the condition that we pre-materialize the symmetric closure of $T(?x, INV, ?y)$.

Second type of duplicates source. The second type of duplicate generations comes from the input data which might contain some triples that make the application of two different rules perfectly equivalent.

We have identified an example of such a case in the Linked Life Dataset, that is one realistic dataset that we used to evaluate our approach. In this dataset there is the triple $T(SCO, TYPE, TRANS)$ which states that the *rdfs:subClassOf* predicate is transitive.

In this case, during the precomputation phase the query $T(a, SCO, b)$ will be launched several times, and

each time the reasoner will trigger the application of both the rules scm-sco and prp-trp.

However, since the application of these two rules will lead to the same derivation, the computation is redundant and inefficient. To detect such cases, we can apply a special algorithm when the system is starting up and initializing the rule set. A complete description of this algorithm is outside the scope of this paper and we will simply illustrate the main idea behind it.

Basically, this algorithm compares each rule with all the others in order to identify under which conditions the two will produce the same output to a given query. For example, the rules scm-sco and prp-trp will produce the same derivation if (i) the input contains the triple $T(SCO, TYPE, TRANS)$ and if (ii) there is a query with SCO as a predicate.

In order to verify this is the case, the algorithm checks whether the triple $T(SCO, TYPE, TRANS)$ exists in the input and there is a matching on the position of the variables in the two rules (if one rule contains more variables than the other, then the algorithm will substitute the corresponding terms). If such a matching exists, then the two rules are equivalent. In our example, the algorithm will find out that the rule prp-trp is equivalent to scm-sco if we replace $?p$ with SCO. Therefore, if there is an input query with SCO as predicate, the system will execute only one of the two rules, avoiding in this way a duplicated derivation.

6. Evaluation

We have implemented our approach in a Java prototype that we called *QueryPIE* and we evaluated the performance using one machine of the DAS-4 cluster,⁷ which is equipped with a dual Intel E5620 quad core CPU of 2.4 GHz, 24 GB of memory and 2 hard disks of 1 TB each, configured in RAID-0 mode.

We used two datasets as input. LUBM [9], which is one of the most popular benchmarks for OWL reasoning and LLD (Linked Life Data),⁸ which is a curated collection of real-world datasets in the bioinformatics domain.

LUBM allows us to generate datasets of different sizes. For our experiments, we generated a dataset of 10 billion triples (which corresponds to the generation of 80000 LUBM universities). The Linked Life Data dataset consists of about 5 billion triples. Both

⁷<http://www.cs.vu.nl/das4>

⁸<http://www.linkedlifedata.com/>

Dataset	Reasoning time		N. iterations	N. derived triples
	Our approach	Full materialization		
LUBM	1.0s	4d4h16m	4	390
LLD	16m	5d10h45m	7	10 millions

Table 3

Execution time of the pre-materialization algorithm compared to a full closure.

datasets were compressed using the procedure described in [20].

The QueryPIE prototype uses six indices stored alongside with the triple permutations on disk using an optimized B-Tree data structure. During the reasoning process, the inferred triples are stored and cached in the main memory. Furthermore, the content of the pre-materialization is indexed at every iteration to facilitate the retrieval of the *lookup* function.

We organized this section as follows. First, in Section 6.1 we report a set of experiments to evaluate the performance of the pre-materialization phase. Next, in Section 6.2 we focus on the performance of the backward-chaining approach and analyze its performance on some example queries. Finally, in Section 6.3, we present a more general discussion on the results that we obtained.

6.1. Performance of the pre-materialization algorithm

We launched the pre-materialization algorithm on the two datasets to measure the reasoning time necessary to perform the partial closure. The results are reported in the second column of Table 3. Our prototype performs joins between the pre-materialized patterns when it loads the rules in memory, therefore, we have also included the startup time along with the query runtimes to provide a fair estimate of the time requested for the reasoning.

From the results, we notice that the pre-materialization is about three orders of magnitude faster for the LUBM dataset than for LLD. The cause for this difference is that the ontology of LUBM requires much less reasoning than the one of LLD in order to be pre-materialized. In fact, in the first case the pre-materialization algorithm has derived 390 triples and needing four iterations to reach a fix point. In the other case, the pre-materialization required 7 iterations and returned about 10 million triples.

We intend to compare the cost of performing the partial closure against the cost of a full materialization, which is currently considered as state of the art in the field of large scale OWL reasoning. However, to the best of our knowledge there is no approach described in literature which supports the OWL RL fragment and that is able to scale to the input size that we consider.

The closest approach we can use for a comparison is WebPIE [19], which has demonstrated OWL reasoning up to the pD^* fragment to a hundred billion triples. Since WebPIE uses the MapReduce programming model, an execution on a single machine would be suboptimal. Therefore, we launched it using eight machines and multiplied the execution time accordingly to estimate the runtime on one machine (the estimation is in line with the performance of WebPIE which has shown linear scalability in [19]).

The runtime of the complete materialization performed with this method is reported in the third column of Table 3. Note that in both cases a complete materialization requires between four and five days against the seconds or minutes required for our method. This comparison clearly illustrates the advantage of our approach in terms of pre-materialization cost. However, this advantage comes at a price: although after a complete materialization reasoning is no longer needed, our backward-chaining algorithm still has to perform some inference at query time. The impact of this operation on the query-time performance is analyzed in the next section.

6.2. Performance of the reasoning at query time

In order to analyze the performance of reasoning at query time, we launched some example queries after we computed the closure using our backward-chaining algorithm to retrieve the results. For this purpose, we selected six example queries for both the LUBM and LLD datasets and reported them in Table 4.

While LUBM provides an official set of queries for benchmarking, there is unfortunately no official set of queries that can be used for benchmarking the performance on the LLD dataset. Therefore, we took some queries that are reported on the official web page of the LLD dataset and modified them so that they could trigger different types of reasoning.

These queries were selected according to the following criteria:

- *Number of results*: We selected queries that return a number of results that varies from no results to a large set of triples;

ID	Dataset	Query
Pattern 1	LUBM	?x ?y <http://www.Department0.University0.edu/GraduateCourse0>
Pattern 2	LUBM	?x <lubm:subOrganizationOf> <http://www.University0.edu>
Pattern 3	LUBM	<http://.../GraduateStudent124> <lubm:degreeFrom> <http://www.University114.edu>
Pattern 4	LUBM	?x ?y <http://www.Department0.University0.edu/AssistantProfessor0>
Pattern 5	LUBM	?x <lubm:memberOf> <http://www.Department0.University0.edu>
Pattern 6	LUBM	?x <rdf:type> <lubm:Department>
Pattern 7	LLD	?x ?y <lifeskim:mentions>
Pattern 8	LLD	?x <lifeskim:mentions> <http://linkedlifedata.com/resource/umls/id/C0439994>
Pattern 9	LLD	<http://.../resource/pubmed/id/15964627> ?x ?y
Pattern 10	LLD	?x ?y <http://purl.uniprot.org/go/0006952>
Pattern 11	LLD	?x ?y <http://linkedlifedata.com/resource/umls/id/C0439994>
Pattern 12	LLD	?x <http://www.biopax.org/release/biopax-level2.owl#NAME> ?y

Table 4

List of example queries

Q.	Runtime (ms)		Derived Triples		I/O access	
	Cold	Warm	Total	Output	#	MB
1	60.43	6.39	5	5	43	8
2	1099.28	129.31	463	239	12	205
3	49.18	6	3	1	18	5
4	73.06	11.17	37	29	86	8
5	118.71	13.97	1480	719	18	8
6	4026.27	2590.27	1599987	1599987	2	12
7	228.26	214.57	0	0	670	23
8	23.74	6.29	4466	4466	1	4
9	7064.04	609.4	140	128	3540	105
10	2535.38	1103.48	28446	26860	14372	337
11	2613.37	1883.14	8546	4504	15128	64
12	2334.70	2059.20	1187944	1187944	1	10

Table 5

Runtime of the queries in Table 4 on the LUBM and LLD datasets

- *Reasoning complexity*: Some queries in our example set require no reasoning to be answered, in contrast other queries generate a very large proof-tree;
- *Amount of data processed*: In order to answer a query, the system might need to access and process a large set of data. We selected queries that read and process a variable amount of data to verify the impact of I/O on the overall performance.

We performed a number of experiments to analyze three aspects of the performance of our algorithm during query time: the *absolute response time*, the *reduction of the proof tree*, and the *overhead induced by rea-*

soning during query-time. Each of these aspects is analyzed below.

6.2.1. Absolute response time

We report in Table 5 the execution times obtained launching the selected example queries in Table 4. In the second and third columns we report both the cold and warm runtime.⁹ With cold runtime, we identify the runtime that is obtained by launching the query right after the system has started. Since the data is stored on disk, we also measure with the cold runtime the time to read the data from disk. On the other side, the warm runtime measures the average response time of launching the same query thirty more times. During this execution the data is already cached in memory and the Java VM has already initialized the internal data structures, so the warm runtime is significantly faster than the cold one.

The fourth and fifth column, respectively, report the total number of derivations that were inferred during the execution of the query, and the number of triples returned to the user.

The sixth and seventh column report the number of data lookups required to answer the query and the amount of data that is read from disk. These two numbers are important for estimating the impact of reasoning at query time. Whilst querying without reasoning only requires a data lookup, the backward-chaining algorithm might require to access the database multiple times. For example, in order to answer query 11 the

⁹Please note that the reported runtime does not include the time required to compress/decompress the numerical terms to their string counterpart.

program had access to the data indices about 15000 times.

Using the results reported in Table 5, we can make several observations. First, we notice that the cold runtime is in general significantly higher than the warm runtime between one and two orders of magnitude. This is primarily due to the time consuming I/O access to disk especially because reasoning requires to read in different locations of the data indices, and therefore the system is required to read several blocks of the B-Tree from the disk. For most of the queries, the I/O access dominates the execution time. The worst case presents query 10 where the program reads from disk about 337 MB of data. We conclude from these results that the performance of the program is essentially I/O bounded, if the data is stored on disk. After the data is loaded in memory, the execution time drops by about one order of magnitude on average and the performance becomes CPU bounded.

Another factor that impacts the performance is the number of the inferred triples that are calculated during the execution of the query. In fact, we notice that absolute performance is lower in case a large number of triples is either inferred or retrieved from the database. The behavior is due to the fact that the algorithm needs to temporarily store these triples as it must consider them in each repeat-loop pass until the closure is reached. This means that these triples must be stored and indexed to be retrieved during the following iterations and the response time consequently increases.

Summarizing our analysis, we make the following conclusions: (i) the runtime is influenced by several factors among which the most prominent is the amount of I/O access that is requested to answer the query (this number is proportional to the size of the proof tree) and the number of derivations produced. (ii) There is a large difference in the runtime observed in our experiments. In the worst case the absolute runtime is in the range of few seconds, while in the best cases the performance is in the order of dozens of milliseconds. However, even in the worst case the system allows an interactive usage since few seconds are acceptable in most scenarios.

In Section 6.2.3 we compare the response times to those without reasoning in order to have a better overview of the overall performance and understand what the overhead induced by reasoning at query time is.

Query	# Leaves proof tree		Reduction ratio
	Without precomp.	Our approach	
1	16	4	4.0
2	2	1	2.0
3	12	3	4.0
5	26	7	3.7

Table 6

Estimation of the reduction of the proof tree caused by the pre-materialization algorithm.

6.2.2. Reduction of the proof tree

The backward-chaining algorithm and more in general our approach relies on the pre-materialization of some selected queries which serve a variety of purposes such as performing efficient sideways information passing or excluding rules that derive duplicates. Another advantage of performing the pre-materialization is that it reduces the size of the proof tree during query-time.

In this section, we evaluate the effective reduction in terms of the size of the proof tree obtained by avoiding performing inference on the pre-materialized patterns.

However, since the method presented in this paper is embedded in the implementation of our prototype, and since the optimizations introduced are crucial to its execution, we cannot disable them. To overcome this problem, we have manually analyzed the execution of the LUBM queries with our prototype on a much smaller dataset (of about 100 thousand triples) and manually constructed the proof tree without pre-materialization (note that we excluded queries 4 and 6 since in these cases reasoning did not contribute to derive new answers). In principle, we identified for each query the rules which produce their answers and for each pre-materialized body atom, we added the corresponding branch that was generated when that query was calculated during the pre-materialization phase.

We report the results of this analysis in Table 6. The last column reports the obtained reduction ratio and shows that the number of leaves shrinks between two and four times due to our pre-materialization.

The results delivered by this method of evaluation must be seen as an underestimate, because we could not deactivate all the optimizations, and therefore in reality the gain is even higher than the one calculated. Nevertheless, this shows that our pre-calculation is indeed effective. For a very small cost in both data space and upfront computation time, we substantially reduce the proof tree. Apparently, the pre-materialization pre-

Q.	Only Lookup		RDFS		pD*		OWL RL	
	No Ins.	Ins.	No Ins.	Ins.	No Ins.	Ins.	No Ins.	Ins.
1	0.81	0.83	1.88	1.79	5.4	6.13	6.39	5.89
2	0.82	1.51	1.56	2.83	128.78	131.05	129.31	138.53
3	0.82	0.83	3.55	2.72	5.50	4.51	6	4.83
4	0.88	0.94	2.01	2.32	10.06	9.48	11.17	10.63
5	1.5	1.61	7.01	4.95	13.58	10.52	13.97	10.8
6	405.42	418.38	2605.68	2630.08	2608.20	2619.17	2590.27	2618.66
7	0.77	0.79	176.19	1.26	203.23	17.93	214.57	16.78
8	1.96	1.89	6.23	6.34	6.39	6.46	6.29	6.36
9	0.84	0.90	262.7	46.53	590.34	277.55	609.4	277.02
10	7.90	7.29	212.57	115.16	903.31	814.95	1103.48	1053.33
11	1.85	1.93	200.55	8.35	1695.73	1468	1883.14	1529.64
12	338.14	337.41	2129.49	2044.34	2055.02	2077.55	2059.2	2062.65

Table 7

Runtime (in ms.) of the example queries changing the rule set.

cisely captures small amounts of inferences that contribute substantially to the reasoning costs because they are being used very often.

6.2.3. Overhead of reasoning during query-time

While we are able to significantly reduce the size of the proof tree and apply other optimizations to further reduce the computation, we still have to perform some reasoning during the execution of a query. It is important to evaluate what the cost for the remaining reasoning is when we compare our approach to a full-materialization approach (which is currently the de-facto technique for large scale reasoning), where a large pre-materialization is performed so that during query time reasoning is avoided altogether.

To this end, we launched a number of experiments activating different types of reasoning at query time and report the warm runtime in Table 7.

We proceeded as follows: we first launch the queries, deactivating all rules at query time, and state their execution time in the first column of the table (the title “No Ins.” indicates no insertion). We then reissued the queries activating only the RDFS rules (in the third column), then the pD^* rules and finally the OWL RL ones.

The results reported under the “Ins.” columns (“Ins.” means insertion) were calculated differently. In fact, in the previous experiments the number of retrieved results for a specific query might differ because we changed the rule set and this can influence the general performance. To maintain the number of results constant, we have repeated the same experiment adding to the knowledge base all the possible results so that even

if reasoning is not activated the same number of results is retrieved .

From the results presented in the table, we notice that in both cases (“Ins” and “No Ins.”) the response time progressively increases as we include more rules. Such a behavior is expected since more computation must be performed as we add new rules. However, in some cases (like query 12) there is a significant difference even if the query does not require the application of any rule. The difference is due to the cost of storing the results into main memory during the query execution to ensure the completeness of the backward-chaining algorithm. This operation is clearly a non-negligible contributor to the overall performance.

We can compare the response times reported in the third column with the ones of the penultimate column to compare the performance of the reasoning at query time of our approach against traditional full materialization. In fact, because the input data already contains the whole derivation, a single lookup can be used to estimate the cost that we would have to pay if all the inferences were pre-materialized beforehand. From the results we notice that on average the response time is between one and three order of magnitude slower. In case the query needs to process and/or return many triples, the difference is certainly significant. However, the response time is still in the order of the hundreds of milliseconds, from the user perspective, the difference is less noticeable and more easily tolerated especially considering that a large precomputation phase is no longer needed.

6.3. Discussion

In this paper, we chose to evaluate our method using the most common and large-scale datasets currently available in order to evaluate how hybrid reasoning would perform on *current* data and *realistic* queries.

The measurements that we report have shown that large scale OWL RL reasoning is indeed possible, even on a relatively modest computer architecture. However, we must point out that these datasets do not (yet) use all the features introduced with the OWL 2 language and, to the best of our knowledge, there is no large-scale benchmark that extensively uses these new features.

Therefore, there is a remaining open question on what the performance would be on an input that exploits all the features of the OWL 2 language. While this problem is beyond the scope of this article, we can make some considerations by looking at the experiments here presented.

First of all, our approach would most likely not be able to guarantee the same response time in the worst-case scenario. This is not particularly due to a limitation of our method, but rather to the high computational complexity intrinsically required by reasoning.

However, even without looking at the complete worst-case scenario (which is very unlikely to happen in practice), there can be other cases where the performance could be significantly worse. First of all, during the pre-materialization phase the backward-chaining algorithm returns only the answers to the query and discards the intermediate triples. In some cases, this can become particularly inefficient since the same intermediate triples will be re-derived multiple times in answering multiple queries.

Furthermore, we noticed in our experiments that as the size of the proof tree increases, so does the potential derivation of duplicates due to the potential higher number of combinations. In Section 5.1, we tackled this problem by proposing some initial algorithms to limit the number of duplicates. However, our work in this respect is still initial and further research on this particular aspect might become necessary in order to scale not only in terms of input size but also in terms of reasoning complexity.

Summarizing, we observe in our evaluation that fairly complex reasoning can be performed rather quickly (in a matter of few seconds in the worst case) on a small set of realistic queries and on large data. However, the reader should keep in mind that there could be worst-case scenarios (which do not seem to

appear on current data) where the performance is significantly worse, and this is mainly due to the theoretical high worst-case complexity that is inherently present in the reasoning process.

7. Related Work

Applying rules with a top-down method like backward-chaining is a well-known technique in rule-based languages like Datalog [6]. Datalog is a generic and powerful language that can handle an arbitrary number of rules using predicates of any arity. In our work, we optimized the implementation to exploit the characteristics of RDF data and to execute a standard rule set, and ignored features of the language that are not necessary to execute the OWL rules. Therefore, and since there is, to the best of our knowledge, no Datalog engine that can load billions of triples, a direct comparison between our work and similar Datalog engines such as IRIS [5] or Jena [12] is not possible.

The backward-chaining algorithm that we present in Section 3 is inspired by the QSQ and the semi-naive evaluation algorithms which are well-known techniques in logic programming. A similar termination condition to ours is employed also in the RQA/FQI algorithm [14].

In our approach, we exploit the availability of the precomputation using a *sideway information passing* (SIP) technique [3] during the execution of the rules. This technique is used in other approaches like in the magic set rewriting algorithm [2]. However, while the magic set algorithm uses SIP at compile-time to construct rules which are later used in a bottom-up fashion, we employ this technique at runtime to execute queries in a top-down manner. Also, SIP strategies are similarly used in generic query processing to prune irrelevant results. In [10] the authors propose two adaptive SIP strategies where information is passed adaptively between operators that are executed in parallel.

A similar technique to our method is memoing [26]. Memoing is a technique where queries that are frequently requested are cached to improved the performance. The difference between memoing and our work is that in memoing caching is dynamic, whilst in our approach caching is static: we manually select which queries are to be cached.

Some RDF Stores support various types of inference. 4store [17] applies the RDFS rules with backward-chaining. Virtuoso [8] supports the execution of few (but not all) OWL rules. BigOWLIM [4]

is an RDF store that supports the OWL 2 RL rule set by performing a full materialization when the data is being loaded. Another database system that performs OWL RL reasoning in a similar way is Oracle: In [11] the authors describe their approach reporting the performance of the inference over up to seven billion triples. An approach in which the OWL RL rules are used is presented in [18] where the authors have encoded OWL RL reasoning in the context of embedded devices, and therefore optimizing the computation for devices with limited resources.

Some work has been presented to distribute the reasoning process using supercomputers or clusters of machines. In our previous work we used the MapReduce programming model to improve the scalability [19]. In [27], the authors implement RDFS reasoning using the Opteron blade cluster. To the best of our knowledge, none of these approaches supports the OWL RL rules.

Implicit information can be derived not only with rule-based techniques. In [15], the authors focus on ontology based query answering using the OWL 2 QL profile [13] and present a series of techniques based on query rewriting to improve the performance. While we demonstrate inference over a much larger scale, a direct comparison of our technique with this work is difficult since both the language and reasoning techniques are substantially different.

A series of work has been done on reasoning using the OWL EL profile. This language is targeted to domains in which there are ontologies with a very large number of properties and/or classes. [7] presented an extensive survey of the performance of OWL EL reasoners analyzing tasks like classification or consistency checking. Again, the different reasoning tasks and considered language makes a direct comparison difficult for our approach.

8. Conclusions

Until now, all inference engines that can handle reasonably expressive logics over very large triple stores (in the orders of billion of triples) have deployed full materialization. In the current paper we have broken with this mold, showing that it is indeed possible to do efficient backward-chaining over large and reasonably expressive knowledge bases.

The key to our approach is to precompute a small number of inferences which appear very frequently in the proof tree. This of course re-introduces some

amount of preprocessing, but this computation is measured in terms of minutes, instead of the hours needed for the full closure computation.

By pre-materializing part of the inferences upfront instead of during query-time, we are able to introduce a number of optimizations which exploit the pre-computation to improve the performance during query-time. To this end, we adapted a standard backward-chaining algorithm like QSQ to our use case exploiting the parallelization of current architectures.

Since our approach deviates from standard practice in the field, we have formalized the computation using the theory of deductive databases and extensively analyzed and proved its correctness.

We have implemented our method in a proof-of-concept Java prototype and analyzed the performance over both real and artificial datasets of five and ten billion triples using most of the OWL RL rules. The performance analysis shows that the query response-time for our approach is in the low number of milliseconds in the best cases, and increasing up to few seconds as the query increases in its complexity. The loss of response time is offset by the great gain in not having to perform a very expensive computation of many hours before being able to answer the first query.

Obvious next steps in future work would be to investigate how our approach can further scale in terms of data size and reasoning complexity and to understand the properties of the knowledge base that influence both the cost of the limited forward computation and the size of the inference tree. Also, it is worth to explore whether related techniques such as ad-hoc query-rewriting like the one presented in [15] can be exploited to further improve the performance.

To the best of our knowledge, this is the first time that complex backward-chaining reasoning over realistic OWL knowledge bases of a ten billion triples has been realized. Our results show that this approach is feasible, opening the door to reasoning over much more dynamically changing datasets than was possible until now.

Acknowledgments: We would like to thank Stefan Schlobach, Barry Bishop, Boris Motik, and Ian Horrocks for providing useful comments and advice. We also would like to thank Cerial Jacobs for the support in developing and debugging the prototype used in the evaluation.

A sincere thanks goes to the reviewers of this article, Aidan Hogan, Matthias Knorr, Raghava Mutharaju, and Peter Patel-Schneider who have significantly im-

proved the quality of this article with a careful review of the work.

This work was partly supported by the LarKC project (EU FP7-215535), the COMMIT project and by the SEALS Project.

References

- [1] S. Abiteboul, R. Hull, and V. Vianu. *Foundations of Databases*. Addison-Wesley, 1995.
- [2] F. Bancilhon, D. Maier, Y. Sagiv, and J. Ullman. Magic sets and other strange ways to implement logic programs. In *Proceedings of the fifth ACM SIGACT-SIGMOD symposium on Principles of database systems*, pages 1–15. ACM, 1985.
- [3] C. Beeri and R. Ramakrishnan. On the power of magic. *The journal of logic programming*, 10(3):255–299, 1991.
- [4] B. Bishop and S. Bojanov. Implementing owl 2 rl and owl 2 ql rule-sets for owlim. In *Pro. of the 8th International Workshop OWL: Experiences and Directions*, 2011.
- [5] B. Bishop and F. Fischer. Iris-integrated rule inference system. In *International Workshop on Advancing Reasoning on the Web: Scalability and Commonsense (ARea 2008)*, 2008.
- [6] S. Ceri, G. Gottlob, and L. Tanca. *Logic Programming and Databases*. Springer, 1990.
- [7] K. Dentler, R. Cornet, A. ten Teije, and N. de Keizer. Comparison of reasoners for large ontologies in the owl 2 el profile. *Semantic Web*, 2(2):71–87, 2011.
- [8] O. Erling and I. Mikhailov. Sparql and scalable inference on demand. Available at http://www.openlinksw.co.uk/virtuoso/Whitepapers/pdf/Scalable_Inference.pdf.
- [9] Y. Guo, Z. Pan, and J. Hefflin. LUBM: A benchmark for OWL knowledge base systems. *Journal of Web Semantics*, 3:158–182, 2005.
- [10] Z. Ives and N. Taylor. Sideways information passing for push-style query processing. In *Data Engineering, 2008. ICDE 2008. IEEE 24th International Conference on*, pages 774–783. IEEE, 2008.
- [11] V. Kolovski, Z. Wu, and G. Eadon. Optimizing enterprise-scale owl 2 rl reasoning in a relational database system. In *International Semantic Web Conference (1)*, pages 436–452, 2010.
- [12] B. McBride. Jena: A semantic web toolkit. *Internet Computing, IEEE*, 6(6):55–59, 2002.
- [13] B. Motik, B. Grau, I. Horrocks, Z. Wu, A. Fokoue, and C. Lutz. OWL2 web ontology language: Profiles. *W3C Recommendation*, 2009.
- [14] W. Nejdl. Recursive strategies for answering recursive queries—the rqa/fqi strategy. In *Proc. International Conference on Very Large Data Bases (VLDB)*, pages 43–50, 1987.
- [15] H. Pérez-Urbina, I. Horrocks, and B. Motik. Efficient query answering for owl 2. *The Semantic Web-ISWC 2009*, pages 489–504, 2009.
- [16] E. Prud’hommeaux and A. Seaborne. SPARQL Query Language for RDF. W3C Recommendation, 2008.
- [17] M. Salvadores, G. Correndo, S. Harris, N. Gibbins, and N. Shadbolt. The design and implementation of minimal RDFS backward reasoning in 4store. In *ESWC (2)*, pages 139–153, 2011.
- [18] C. Seitz and R. Schönfelder. Rule-based owl reasoning for specific embedded devices. In *International Semantic Web Conference (2)*, pages 237–252, 2011.
- [19] J. Urbani, S. Kotoulas, J. Maassen, F. V. Harmelen, and H. Bal. Webpie: A web-scale parallel inference engine using mapreduce. *Web Semantics: Science, Services and Agents on the World Wide Web*, 10(0):59 – 75, 2012.
- [20] J. Urbani, J. Maassen, N. Drost, F. Seinstra, and H. Bal. Scalable rdf data compression with mapreduce. *Concurrency and Computation: Practice and Experience*, 2012.
- [21] J. Urbani, F. van Harmelen, S. Schlobach, and H. E. Bal. Querypie: Backward reasoning for owl horst over very large knowledge bases. In *International Semantic Web Conference (1)*, pages 730–745, 2011.
- [22] M. H. van Emden and R. A. Kowalski. The semantics of predicate logic as a programming language. *J. ACM*, 23(4):733–742, 1976.
- [23] L. Vieille. Recursive axioms in deductive databases: The query/subquery approach. In *Expert Database Conf.*, pages 253–267, 1986.
- [24] L. Vieille. A database-complete proof procedure based on sld-resolution. In *ICLP*, pages 74–103, 1987.
- [25] L. Vieille. Recursive query processing: The power of logic. *Theor. Comput. Sci.*, 69(1):1–53, 1989.
- [26] D. Warren. Memoing for logic programs. *Communications of the ACM*, 35(3):93–111, 1992.
- [27] J. Weaver and J. Hendler. Parallel materialization of the finite rdfs closure for hundreds of millions of triples. In *8th International Semantic Web Conference (ISWC2009)*, October 2009.