



Approximation and ABox Segmentation

Alissa Kaplunova, Ralf Möller, Sebastian Wandelt, Michael Wessel

October 28, 2008

Abstract

In this technical report, we consider the problem of query answering over large ontologies. Traditional reasoning systems may have problems to deal with large amounts of expressive ontological data (terminological as well as assertional data) that usually must be kept in main memory. We propose to overcome this problem with a new so-called *filter and refine paradigm for ontology-based query answering*.

In the filter step, the terminological part of an ontology is approximated to a less expressive ontology language, e.g., the description logic *DL-Lite*, which allows for efficient and complete, but possibly unsound query answering. A query is first evaluated w.r.t. this *DL-Lite* ontology; thus, the set of ontology individuals is filtered with the help of a *DL-Lite* query. In the second step, this set of retrieved individuals is refined further, without having to perform reasoning over the whole ontology. For this, a partition-based approach is investigated; the aim is to compute partitions which are small enough for main memory reasoning systems. However, the first step can be performed on secondary memory, exploiting database technology.

The contribution of this report is twofold: (1) For both steps, novel algorithms are presented. (2) We evaluate our approach on real-world multimedia ontologies.

1 Introduction

Applying semantic web technologies to enable the semantic retrieval of multimedia documents is a hot research topic. This interest is testified by, for example, emerging core / foundational ontologies for multimedia such as COMM, the “Core Ontology on Multimedia” [6], or the “Athletics Event Ontology” (AEO) developed in the BOEMIE project¹. Here, multimedia “assets” (documents) are annotated with OWL ontologies in order to enable semantic retrieval of multimedia documents, e.g. texts, HTML pages, images, video and sound files. OWL DL is formalized in description logics (DLs); OWL DL corresponds to *SHOIN* (*D*). In this work we are focusing on the DL *SHI* (extensions to larger OWL fragments will be considered in future research). We believe that rather expressive DLs such as *SHI* are required in order to capture important domain constraints in an ontology (e.g., less expressive DLs may not provide the required expressivity for the modeling problems at hand).

Real world ontologies tend to be large, both in numbers of concepts as well as in numbers of individuals. Unfortunately, the data complexity of instance retrieval in *SHI* (and more expressive DLs) is EXPTIME-complete. Thus, from a computational perspective, instance retrieval with large ontologies (containing lots of instances) may be very hard. Although mature DL/ OWL reasoning systems such as RACERPRO exist [11], many reasoning systems for expressive DLs nowadays still work on main memory only. This obviously prevents their usage on very large ontologies, which may contain millions of “facts”, so query answering simply runs out of main memory, or even loading of the whole ontology is already impossible.

Recently, query answering in less expressive DLs received great attention. E.g., the QUONTO system [5] is able to perform query answering on secondary memory by taking advantage of (relational) database technology.

In this work, we propose a pragmatic method to combine the *high performance and data scalability* achieved by the techniques realized in the QUONTO architecture with the *high expressivity and expressivity scalability* realized by state-of-the-art DL reasoners such as RACERPRO. We propose a new *filter & refine strategy* for expressive multimedia ontologies in order to address the *data- and expressivity scalability problem* [12] as follows: In the *filter step*, the QUONTO techniques are utilized to retrieve a set of *candidate instances*, in principle, from secondary memory. This obviously requires that the ontology is a *DL-Lite* ontology (or some other less expressive DLs which allows for complete query answering on secondary memory by exploiting existing RDBMs or RDF Triple Stores). Thus, we present a framework to *automatically approximate, i.e. rewrite* a *SHI* ontology into a *DL-Lite* ontology, such that retrieval w.r.t. the *DL-Lite* ontology remains *complete* compared to retrieval w.r.t. the original ontology. Thus, the query answers computed in such a way constitute a *superset* of the set of *true* query answers (w.r.t. the original ontology). This set is also called the set of *candidate individuals* in the following, from which *false or wrong query answers have to be removed*. This is the purpose of the refine step.

In the *refine step*, the original ontology is used to refine this set of candidate individuals by means of expressive DL reasoning. Since, as already mentioned, the ontology may be too big to be loaded into a main memory-based reasoner all at once, a novel *automatic partitioning algorithm* [13] is applied, which partitions the extensional part of the ontology (also called the ABox) into so-called islands in such a way that, for query answering, it

¹<http://www.boemie.org/>

is sufficient to only load the islands relevant for the query at hand into main memory. In many cases, these islands are small enough to fit into main memory, so false candidate individuals can effectively be rejected using a main memory-based expressive DL reasoner such as RACERPRO.

This report is structured as follows. First, the basics of descriptions logics (as far as relevant for this report) are introduced; i.e., the DLs \mathcal{SHI} and $DL-Lite$, as well as basic inference problems. Then, we describe the novel approximation algorithm which reformulates \mathcal{SHI} ontologies as $DL-Lite$ ontologies for the filter step. We then apply a novel partitioning algorithm for the refine step and perform a preliminary evaluation of our framework applied to the AEO ontology. Open problems are discussed and provide motivation for future research. Finally comes the conclusion and some discussion of related and future work.

2 Basics and Guiding Example

In the following part we will define mathematical notions, which are relevant for the remaining paper.

The Description Logic \mathcal{SHI} We briefly recall syntax and semantics of the description logic \mathcal{SHI} (also called \mathcal{ALCHIR}_+). For the details, please refer to [7]. We assume a collection of disjoint sets: a set of *concept names* N_{CN} , a set of *role names* N_{RN} and a set of *individual names* N_I . The *set of roles* N_R is $N_{RN} \cup \{R^- \mid R \in N_{RN}\}$, where R^- denotes inverse roles. A distinguished subset N_T of roles is called a *set of transitive roles*. The set of \mathcal{SHI} -*concept descriptions* is given by the following grammar:

$$C, D ::= \top \mid \perp \mid A \mid \neg C \mid C \sqcap D \mid C \sqcup D \mid \forall R.C \mid \exists R.C$$

where $A \in N_{CN}$ and $R \in N_R$. We say that a concept description is *atomic*, if it is a concept name. With N_C we denote all atomic concepts. For defining the semantics of concept descriptions and roles we consider *interpretations* \mathcal{I} that consist of a non-empty set $\Delta^{\mathcal{I}}$, the domain, and an interpretation function $\cdot^{\mathcal{I}}$, which assigns to every atomic concept description A a set $A^{\mathcal{I}} \subseteq \Delta^{\mathcal{I}}$ and to every role R a set $R^{\mathcal{I}} \subseteq \Delta^{\mathcal{I}} \times \Delta^{\mathcal{I}}$. For complex concept descriptions the interpretation function is extended as shown in [7]. The semantics of description logics is based on the notion of satisfiability. An interpretation $\mathcal{I} = (\Delta^{\mathcal{I}}, \cdot^{\mathcal{I}})$ *satisfies* a concept description C if $C^{\mathcal{I}} \neq \emptyset$. In this case, \mathcal{I} is called a *model* for C . A concept description is in *negation normal form* if negation occurs only in front of concept names.

A *TBox* is a finite set of axioms of the form $C \sqsubseteq D$ (so-called *generalized concept inclusions*, GCIs) together with a finite set of axioms $R \sqsubseteq S$ (*role inclusions*). An interpretation \mathcal{I} *satisfies* a GCI $C \sqsubseteq D$ if $C^{\mathcal{I}} \subseteq D^{\mathcal{I}}$. An interpretation \mathcal{I} *satisfies* a role inclusion $R \sqsubseteq S$ if $R^{\mathcal{I}} \subseteq S^{\mathcal{I}}$. An interpretation is a *model* of a TBox \mathcal{T} if it satisfies all generalized concept inclusions and role inclusions in \mathcal{T} . An *ABox* is a finite set of so-called *concept and role assertions* $a : C$ and $R(a, b)$, where a and b are elements of $\Delta^{\mathcal{I}}$. An interpretation \mathcal{I} *satisfies* a concept assertion $a : C$ (resp. role assertion $R(a, b)$) if $a \in C^{\mathcal{I}}$ (resp. $(a, b) \in R^{\mathcal{I}}$).

A *ontology* \mathcal{O} is a pair $\langle \mathcal{T}, \mathcal{A} \rangle$, where \mathcal{T} is a TBox and \mathcal{A} is a ABox. We restrict the concept assertions in \mathcal{A} in such a way that each concept description is an atomic concept or a negated atomic concept. This is a common assumption, e.g., in [9], when dealing with large assertional datasets in ontologies. With $Ind(\mathcal{A})$ we denote the set of individuals occurring in \mathcal{A} . We say that \mathcal{O} is *inconsistent*, denoted with $INC(\mathcal{O})$, if there exists no

model for \mathcal{O} . We say that \mathcal{O} is *consistent*, denoted with $CON(\mathcal{O})$, if there exists at least one model for \mathcal{O} . Given an individual a and an atomic concept C , we have $\langle \mathcal{T}, \mathcal{A} \rangle \models a : C$ iff $INC(\langle \mathcal{T}, \mathcal{A} \cup \{a : \neg C\} \rangle)$.

By *instance retrieval for concept C* , we obtain all individuals $a \in Ind(\mathcal{A})$, s.t. we have $\langle \mathcal{T}, \mathcal{A} \rangle \models a : C$.

We denote the set of instances for a given concept C with $concept_instances(C, \mathcal{A}, \mathcal{T})$.

In the following we define some additional notions, which will be used throughout the remaining part of the paper. A \exists -*constraint* is a concept description of the shape $\exists R.C$, s.t. C is an arbitrary concept description. A \forall -*constraint* is a concept description of the shape $\forall R.C$, s.t. C is an arbitrary concept description.

The *subsumption hierarchy* (so-called *taxonomy*) of parents and children for each concept name can be obtained by classification. For \mathcal{SHI} ontologies it is possible to compute the subsumption hierarchy in advance given only the TBox \mathcal{T} , i.e. without the ABox \mathcal{A} . This is possible since \mathcal{SHI} does not allow the use of nominals. With $\sqsubseteq_{\mathcal{T}}: N_C \times N_C$ we denote the precomputed taxonomy obtained by classification, e.g., we have $\sqsubseteq_{\mathcal{T}}(C, D)$ iff $\mathcal{O} \models C \sqsubseteq D$ for atomic concepts C and D . The role hierarchy of a \mathcal{SHI} -ontology can be computed in advance given the TBox \mathcal{T} only as well. With $\sqsubseteq_{\mathcal{R}}: N_R \times N_R$ we denote the precomputed role hierarchy, e.g. we have $(R, S) \in \sqsubseteq_{\mathcal{R}}$ iff $\mathcal{O} \models R \sqsubseteq S$ for roles R and S .

An atomic concept D is a synonym for a concept description C if we have $\mathcal{T} \models C \sqsubseteq D$ and $\mathcal{T} \models D \sqsubseteq C$. With $synonyms(C, \mathcal{T})$ we denote the set of atomic concepts, which are synonyms for concept C with respect to \mathcal{T} . With $parents(C, \mathcal{T})$ ($children(C, \mathcal{T})$) we denote the set of atomic concepts which are more general (specific) than a given concept C .

The Description Logic $DL-Lite_{\mathcal{F}}$ $DL-Lite_{\mathcal{F}}$ ([8]) is a fragment of OWL DL with the set of concept and role descriptions defined by the following grammar: $B \rightarrow A \mid \exists R.T, C \rightarrow B \mid \neg B, R \rightarrow P \mid P^-, E \rightarrow R \mid \neg R$

where A is an *atomic concept*, P is an *atomic role*, and P^- is the *inverse* of the atomic role P . B denotes a *basic concept*, i.e., a concept that can be either an atomic concept or a concept of the form $\exists R.T$, and \top is the top concept, R denotes a *basic role*, i.e., a role that is either an atomic role or the inverse of an atomic role. $DL-Lite_{\mathcal{F}}$ TBox is a set of axioms of the form $B \sqsubseteq C$ (concept inclusions), where only basic concepts may occur on the left-hand side, and a set of global role functionality assertions of the form $funct(R)$. An ABox includes as usually concept and role assertions $a : C$ and $R(a, b)$.

In the following we define an example ontology, which is used throughout the remaining part of the paper. The ontology is a simplified version of the LUBM [10]. Let $\mathcal{O}_{EX} = \langle \mathcal{T}_{EX}, \mathcal{A}_{EX} \rangle$, s.t.

$$\begin{aligned} \mathcal{T}_{EX} = \{ & Chair \sqsubseteq \exists headOf.Department \sqcap Person, \\ & Professor \sqsubseteq Faculty, Book \sqsubseteq Publication, \\ & GraduateStudent \sqsubseteq Student, Student \sqsubseteq Person \sqcap \exists takesCourse.Course, \\ & \top \sqsubseteq \forall teacherOf.Course, \exists teacherOf.\top \sqsubseteq Faculty, Faculty \sqsubseteq Person, \\ & \top \sqsubseteq \forall publication.Author^-.(Book \sqcup ConferencePaper), \\ & headOf \sqsubseteq worksFor, worksFor \sqsubseteq memberOf, memberOf \sqsubseteq member^- \} \end{aligned}$$

The ABox \mathcal{A}_{EX} is shown in Fig. 1. Please note that individual $p2$ is a non-obvious instance (i.e. we need to perform reasoning) of concept *Chair*, since $p2$ has an outgoing *headOf*-edge to a *Department* and every *Professor* is necessarily a *Person*.

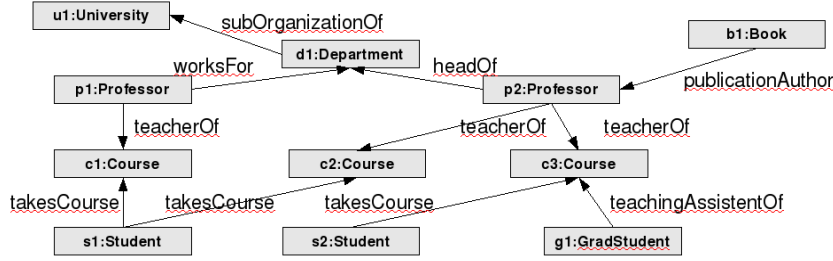


Figure 1: Guiding Example: ABox \mathcal{A}_{EX} for ontology \mathcal{O}_{EX}

3 Terminological Approximation - The Filter Step

Definition of Approximation Let us start with some basic definition. First we define the notion of an approximation of a TBox \mathcal{T} :

Definition 1 Let \mathcal{T}_1 be a TBox in some DL \mathcal{DL} . A \mathcal{T}_2 is called an approximation of \mathcal{T}_1 iff a) \mathcal{T}_2 is a \mathcal{DL}' TBox, with $\mathcal{DL}' \subseteq \mathcal{DL}$, and b) $\mathcal{T}_2 \models \mathcal{T}_1$ holds².

Definition 2 For two TBoxes \mathcal{T}_1 and \mathcal{T}_2 , $\mathcal{T}_2 \models \mathcal{T}_1$ iff all models of \mathcal{T}_2 are also models of \mathcal{T}_1 .

TBox entailment is decidable if \mathcal{DL} is decidable, since $\mathcal{T}_2 \models \mathcal{T}_1$ iff for all $C \sqsubseteq D \in \mathcal{T}_1$, $C \sqcap \neg D$ is unsatisfiable w.r.t. \mathcal{T}_2 . Note that this is well-defined, since we assume $\mathcal{DL}' \subseteq \mathcal{DL}$.

After all, our intention for this definition is that instance retrieval over \mathcal{A} w.r.t. \mathcal{T}_2 shall be complete, but possibly unsound compared with instance retrieval w.r.t. \mathcal{A} and \mathcal{T}_1 :

Proposition 1 Let \mathcal{A} be an ABox which contains only atomic concept assertions, i.e., for all $i : C \in \mathcal{A}$, C is an atomic concept: $C \in N_{CN}$. Let D be an atomic query concept, the concept whose instances shall be retrieved. Let \mathcal{T}_2 be an approximation of \mathcal{T}_1 . Then, the following holds: $\text{concept_instances}(D, \mathcal{A}, \mathcal{T}_1) \subseteq \text{concept_instances}(D, \mathcal{A}, \mathcal{T}_2)$.

Proof 1 (Sketch) Assume that $i \in \text{concept_instances}(D, \mathcal{A}, \mathcal{T}_1)$, but $i \notin \text{concept_instances}(D, \mathcal{A}, \mathcal{T}_2)$. Let Σ_i , $i \in \{1, 2\}$ denote the logical theory of \mathcal{T}_i and \mathcal{A} , where $\alpha \in \Sigma_i$ is either an ABox assertion or a TBox axiom in \mathcal{DL}' (not \mathcal{DL}). Then, $\Sigma_1 \subseteq \Sigma_2$, since \mathcal{T}_2 is an approximation of \mathcal{T}_1 , so $\mathcal{T}_2 \models \mathcal{T}_1$ by definition. Obviously, $i \in \text{concept_instances}(D, \mathcal{A}, \mathcal{T}_j)$ for $j \in \{1, 2\}$ iff $i : D \in \Sigma_j$ (note that $D \in \mathcal{DL}'$). But then, $i \in \text{concept_instances}(D, \mathcal{A}, \mathcal{T}_1)$ implies $i \in \text{concept_instances}(D, \mathcal{A}, \mathcal{T}_2)$, contradicting the assumption.

How to Compute Approximations Having given these definitions, the question arises, how to actually compute an approximation of \mathcal{T} . The idea of the *approximation algorithm* is quite simple. W.l.o.g. we assume that a TBox \mathcal{T} contains only implication axioms (axioms of the form $C \sqsubseteq D$; an axiom $C \equiv D$ is transformed into two axioms $C \sqsubseteq D$, $D \sqsubseteq C$). Please note that \mathcal{SHI} admits role inclusion axioms (for roles R, S) $R \sqsubseteq S$, which are valid in $DL\text{-Lite}$ as well. Regarding transitive roles, which are not allowed in $DL\text{-Lite}$, the following well-known “trick” from the modal logic realm can be applied:

²We are discussing the case where $\mathcal{DL} = \mathcal{SHI}$, and $\mathcal{DL}' = DL\text{-Lite}$.

Definition 3 Let R be a transitive role in \mathcal{T}^3 . Let $N_{CN}(\mathcal{T})$ denote the set of concept names appearing in \mathcal{T} , and $N_T(\mathcal{T})$ the set of transitively closed roles in \mathcal{T} . The K_4 closure of \mathcal{T} , \mathcal{T}^{K_4} , is defined as follows:

$$\mathcal{T}^{K_4} =_{def} \mathcal{T} \cup \{ \exists R. \exists R. C \dot{\sqsubseteq} \exists R. C, \forall R. C \dot{\sqsubseteq} \forall R. \forall R. C \mid C \in N_{CN}(\mathcal{T}), R \in N_T(\mathcal{T}) \}.$$

Moreover, we assume that R is an ordinary role in $\mathcal{T}^{K_4} =_{def}$ (not a transitively closed one).

Proposition 2 Let D be an atomic query concept, and \mathcal{A} an ABox in which all concept assertions refer to atomic concepts only.

Then, $\text{concept_instances}(D, \mathcal{A}, \mathcal{T}) = \text{concept_instances}(D, \mathcal{A}, \mathcal{T}^{K_4})$.

We assume a corresponding function `get_K4_closure` which computes the \mathcal{T}^{K_4} for a given \mathcal{T} . Please note that this proposition does not hold for arbitrary ABoxes and query concepts D (only for ABox containing atomic concept assertions, and atomic instance retrieval concepts).

Another preprocessing step is applied to remove nested occurrences of (sub) concepts of the form $\exists R.C$ and $\forall R.C$ from the axioms, so they can be better approximated to *DL-Lite* axioms. Thus, for each axiom $C \dot{\sqsubseteq} D$, and for each subconcept E in $\neg C \sqcup D$ with $E = \exists R.F$ or $E = \forall R.F$, and $F \notin N_{CN}$, we replace E with a new atomic concept C_E and add $\{C_E \dot{\sqsubseteq} E, E \dot{\sqsubseteq} C_E\}$ to \mathcal{T} . This process continues, until \mathcal{T} no longer contains such axioms (note that E itself might still contain such subconcepts as well). For example, $\{C \dot{\sqsubseteq} D \sqcap \exists R.(E \sqcap F)\}$ is rewritten into $\{C \dot{\sqsubseteq} D \sqcap C_{\exists R.(E \sqcap F)}, C_{\exists R.(E \sqcap F)} \dot{\sqsubseteq} E \sqcap F, E \sqcap F \dot{\sqsubseteq} C_{\exists R.(E \sqcap F)}\}$. Consequently, we assume a function `flatten_tbox` which applies this transformation to a TBox \mathcal{T} . Each model of `flatten_tbox`(\mathcal{T}) is trivially also a model of \mathcal{T} , and vice versa, each model of \mathcal{T} can uniquely be extended to a model of `flatten_tbox`(\mathcal{T}) (only the new atomic concepts must be interpreted correctly so that their axioms become satisfied).

For an *SHI* TBox \mathcal{T} we can now compute an approximated \mathcal{T}' by approximating each axiom. So, $C \dot{\sqsubseteq} D \in \mathcal{T}$ is replaced by a logically stronger axiom $C' \dot{\sqsubseteq} D'$, $\{C' \dot{\sqsubseteq} D'\} \models \{C \dot{\sqsubseteq} D\}$, which is a *DL-Lite* axiom. The algorithm is best understood as a non-deterministic algorithm which works as follows (the actual deterministic implementation is described briefly below):

Function `approximate`(\mathcal{T})

Parameter: *SHI* TBox \mathcal{T}

$\mathcal{T} := \text{flatten_tbox}(\text{get_K4_closure}(\mathcal{T}))$

$\mathcal{T}' := \{C \dot{\sqsubseteq} D \mid \mathcal{T} \models C \dot{\sqsubseteq} D, C, D \in N_{CN}\}$

while $\mathcal{T} \neq \emptyset$

$axiom := \text{select_axiom}(\mathcal{T})$

$\mathcal{T}' := \mathcal{T}' \cup \text{approximate_axiom}(axiom, \mathcal{T}')$

$\mathcal{T} := \mathcal{T} \setminus \{axiom\}$

end while

return \mathcal{T}'

The algorithm first syntactically transforms the input TBox \mathcal{T} as explained. Although `flatten_tbox` introduces new atomic concepts, no additional “K4” axioms need to be introduced for them by `get_K4_closure`. Then, the taxonomy of \mathcal{T} is made explicit by adding corresponding axioms to \mathcal{T}' ; these axioms are *DL-Lite* axioms. The reason for this addition to \mathcal{T}' is that the taxonomy of \mathcal{T} shall be available for `approximate_axiom` (see below). Both `select_axiom` and `approximate_axiom` are non-deterministic as well. Given an axiom $C \dot{\sqsubseteq} D$, the basic idea of `approximate_axiom` is to *generalize* the left-hand side C to C' , and to *specialize* the right-hand side D to D' . This ensures that the approximated axiom is stronger than the original axiom, since $C' \dot{\sqsubseteq} D' \models C \dot{\sqsubseteq} D$ iff $\neg C' \sqcup D' \models \neg C \sqcup D$

³*DL-Lite* does not offer transitive roles.

iff $(\neg C' \sqcup D') \sqcap \neg(\neg C \sqcup D)$ is unsatisfiable iff $(\neg C' \sqcup D') \sqcap C \sqcap \neg D$ is unsatisfiable iff both $\neg C' \sqcap C \sqcap \neg D$ and $D' \sqcap C \sqcap \neg D$ are unsatisfiable. Then, either $C \dot{\sqsubseteq} D$ (so this is a tautology, and thus the trivial case), or $C \dot{\sqsubseteq} C'$ (then $C \sqcap \neg C'$ is unsatisfiable) and $D' \dot{\sqsubseteq} D$, (so $D' \sqcap \neg D$ is unsatisfiable). In principle, it is of course also sufficient to find equivalent C', D' in *DL-Lite*. The concepts C' and D' are called *possible rewritings* of C resp. D , and $C' \dot{\sqsubseteq} D'$ is called a *possible rewriting* of $C \dot{\sqsubseteq} D$ in the following, or also a *candidate rewriting*.

For example, the axiom $C \dot{\sqsubseteq} D \sqcup E$ can be rewritten to $C \dot{\sqsubseteq} D$, or to $C \dot{\sqsubseteq} E$ (assuming that $C, D, E \in N_{CN}$). Moreover, $C \dot{\sqsubseteq} D \sqcup E$ can also be written as $\neg D \dot{\sqsubseteq} \neg C \sqcup E$, $\neg E \dot{\sqsubseteq} \neg C \sqcup D$, or even $\neg D \sqcap C \dot{\sqsubseteq} E$, and so on, yielding additional rewriting possibilities. Thus, re-arranging the left-hand sides of the axioms maximizes the number of rewriting possibilities. Even though these axioms are still equivalent to the original one, after rewriting into *DL-Lite* they no longer are. Perhaps for some reordering, no better approximations than $\top \dot{\sqsubseteq} \perp$ can be found. It is thus even more important to maximize the number of possible approximations in order to avoid bad approximations which are *too strong* (rendering the whole TBox unsatisfiable).

The `approximate_axiom` function considers the input axiom $C \dot{\sqsubseteq} D$ as a disjunction $\neg C \sqcup D$ which, in a first step, is brought into *disjunctive normal form (DNF)*. A concept is in DNF if it is in *negation normal form (NNF)*, and does not contain any (sub)concepts of the form $D \sqcap (E \sqcup F)$. Using simple boolean algebra, each concept can be brought into DNF. Note that the concepts are even simpler at this step in the processing chain, because complex qualification concepts have been removed in advance. In the following, we use the set notation for disjuncts of a concept in DNF: $\text{DNF}(C \sqcap (E \sqcup F)) = (C \sqcap E) \sqcup (C \sqcap F) = \{C \sqcap E, C \sqcap F\}$. The function `approximate_axiom` non-deterministically chooses a subset of $\text{DNF}(\neg C \sqcup D)$ as a possible left-hand side of the axiom, and uses the remaining disjuncts as right-hand side. Then, `approx_axiom` calls the non-deterministic functions `generalize` and `specialize`:

Function `approximate_axiom(axiom, T')`

Parameter: \mathcal{SHI} axiom $axiom = C \dot{\sqsubseteq} D$ and partial approximation T'

```

if  $T' \models axiom$  then return  $T'$ 
else if  $axiom$  is a DL-Lite axiom then return  $\{axiom\} \cup T'$ 
else
   $concept := \text{DNF}(\neg C \sqcup D)$ 
   $left\_side := \text{some\_subset\_of}(concept)$ 
   $right\_side := concept \setminus left\_side$ 
   $left\_side' := \text{generalize}(\neg left\_side, T')$ 
   $right\_side' := \text{specialize}(right\_side, T')$ 
  if  $left\_side' \neq \emptyset$  and  $right\_side' \neq \emptyset$  then
     $axiom' := left\_side' \dot{\sqsubseteq} right\_side'$ 
    if  $T' \not\models axiom'$  then return  $\{axiom'\} \cup T'$ 
return  $T'$ 

```

Both `specialize` and `generalize` first bring their argument concepts in DNF, and then specialize or generalize using a set of *non-deterministic rewriting rules* which are guided by the structure of the concept. The rules are applied exhaustively to the concept C until no more rule is applicable.

The rules make use of the helper function `syns_or_parents` which returns a non-empty result for *non-atomic concepts only*:

$$\text{syns_or_parents}(C, T') =_{def} \begin{cases} \text{synonyms}(C, T') & \text{if } \text{synonyms}(C, T') \neq \emptyset, C \notin N_{CN} \\ \text{parents}(C, T') & \text{if } \text{synonyms}(C, T') = \emptyset, C \notin N_{CN} \\ \emptyset & \text{otherwise} \end{cases}$$

Note that \mathcal{T}' is only partially available, but already contains the taxonomy axioms derived from \mathcal{T} (see `approximate`). Note that $C \in \text{synonyms}(C, \mathcal{T}, \mathcal{T}')$ for all $C \in N_{CN}$.

The function `generalize` uses the following non-deterministic *generalization rules*; $C \rightarrow_G C'$ means that C is generalized to C' :

- $C \rightarrow_G C'$, if C is a valid left-hand side of a *DL-Lite* axiom
- $\exists R.C \rightarrow_G C'$, $C' \in \{\exists R.\top\} \cup \text{syns_or_parents}(\exists R.C, \mathcal{T}, \mathcal{T}')$ (note: $C \in N_{CN}$)
- $C \sqcap D \rightarrow_G C'$, $C' \in \{C, D\} \cup \text{syns_or_parents}(C \sqcap D, \mathcal{T}, \mathcal{T}')$
- $C \sqcup D \rightarrow_G C'$, where $C' = C_1 \sqcup D_1$, with $C \rightarrow_G C_1$, $D \rightarrow_G D_1$,
or $C' \in \text{syns_or_parents}(C \sqcup D, \mathcal{T}, \mathcal{T}')$
- for all other concepts C : $C \rightarrow_G C'$, $C' \in \text{syns_or_parents}(C, \mathcal{T}, \mathcal{T}')$

To give an example, consider `generalize` is applied to $\exists R.C \sqcup (E \sqcap F)$. First, the DNF is computed: $(\exists R.C \sqcap E) \sqcup (\exists R.C \sqcap F)$. Then, a possible rewriting is: $(\exists R.C \sqcap E) \sqcup (\exists R.C \sqcap F) \rightarrow_G \exists R.\top \sqcup F$, since $(\exists R.C \sqcap E) \rightarrow_G \exists R.\top$ and $(\exists R.C \sqcap F) \rightarrow_G F$. There are many other different rewritings.

Please note that *DL-Lite* does not permit negation or conjunctions on the left-hand sides of axioms; thus, it is impossible to generalize conjunctions by generalizing the arguments analog to the \sqcup -case. Note that, from this definition, in most cases $\forall R.C \rightarrow_G \top$ unless `syns_or_parents` finds some parent for $\forall R.C$ in \mathcal{T}' . In principle, it is also possible to generalize a disjunction $C \sqcup D$ to something like $C \sqcup D \sqcup E$, for some $E \in N_{CN}$ (although this will result in a huge search space in the implementation). The rules are designed in such a way to avoid *over-generalization* in order to keep the number of unsound query answers small. That means, more specific rewriting alternatives shall be favored over less specific ones. For example, although $C \sqcap D \rightarrow_G C \sqcup D$ is conceivable, it doesn't make much sense under this premise, since both $C \sqcap D \rightarrow_G C$ as well as $C \sqcap D \rightarrow_G D$ are more specific.

The rules for concept specialization, `specialize`, exploit a similar function `syns_or_children` and follow the principle to avoid *over-specialization*, i.e., more general rewriting alternatives are preferred over more specific ones. In these rules, there is the possibility to rewrite a concept C to \emptyset . In case $C \rightarrow_S \emptyset$ for a conjunct C in $C \sqcap D$, then \emptyset is considered as \top . However, in case C is a disjunct, then \emptyset is considered as \perp . So, \emptyset serves as the neutral element w.r.t. the surrounding operation:

- $C \rightarrow_S C'$, if C is a valid right-hand side for a *DL-Lite* axiom
- $\neg C \rightarrow_S \neg C'$, where $C \rightarrow_G C'$ (i.e., C is generalized),
or $C \in \text{syns_or_children}(\exists R.C, \mathcal{T}, \mathcal{T}')$.
- $\exists R.C \rightarrow_S C'$, $C' = \exists R_C.\top$ with $\mathcal{T}' := \mathcal{T}' \cup \{R_C \dot{\sqsubseteq} R, \exists R_C.\top \dot{\sqsubseteq} C\}$,
or $C' = \exists R.\top$ with $\mathcal{T}' := \mathcal{T}' \cup \{\exists R.\top \dot{\sqsubseteq} C\}$,
or $C \in \text{syns_or_children}(\exists R.C, \mathcal{T}, \mathcal{T}')$ (note: $C \in N_{CN}$)
- $\forall R.C \rightarrow_S \emptyset$, $\mathcal{T}' := \mathcal{T}' \cup \{\exists R.\top \dot{\sqsubseteq} C'\}$, where $C \rightarrow_S C'$
- $C \sqcup D \rightarrow_S C'$, $C' \in \{C, D\} \cup \text{syns_or_children}(C \sqcup D, \mathcal{T}, \mathcal{T}')$
- $C \sqcap D \rightarrow_S C'$, where $C' = C_1 \sqcap D_1$, with $C \rightarrow_S C_1$, $D \rightarrow_S D_1$,
or $C' \in \text{syns_or_children}(C \sqcap D, \mathcal{T}, \mathcal{T}')$
- for all other concepts C : $C \rightarrow_S C'$, $C' \in \text{syns_or_children}(C, \mathcal{T}, \mathcal{T}')$

In principle, it is possible to use $C \sqcap D \rightarrow_S C \sqcap D \sqcap E$, for some $E \in N_{CN}$, but the same comments as given above (for $C \sqcup D$) apply. Please note that *DL-Lite* does not permit disjunctions on the right-hand sides of axioms. Moreover, `specialize` has a side-effect on \mathcal{T}' , since it may introduce additional axioms. For example, the $\exists R.C$ -rule introduces a new range restriction on R by adding $\exists R^-. \top \sqsubseteq C$ to \mathcal{T}' . So, $\exists R.C$ is in fact *generalized* to $\exists R. \top$; however, due to the introduced range restriction $\exists R^-. \top \sqsubseteq C$ we get $\exists R. \top \models \exists R.C$. In combination this is a specialization of $\exists R.C$, as required. Another possibility would be to introduce a subrole R_C , $R_C \sqsubseteq R$ with range C , and rewrite $\exists R.C$ to $\exists R_C. \top$, but this would require a modification of the ABox during instance retrieval.

The rule $\forall R.C \rightarrow_S \emptyset$ deserves an explanation. The idea here is to completely ignore this (sub)concept on the right-hand side, and instead put a new axiom into \mathcal{T}' (which is modified per side-effect): $\mathcal{T}' := \mathcal{T}' \cup \{\exists R^-. \top \sqsubseteq C'\}$. For example, consider the TBox $\{C \sqsubseteq (\forall R.D) \sqcap E\}$. Since the left-hand side is already acceptable, only the right-hand side is rewritten: $(\forall R.D) \sqcap E \rightarrow_S \top \sqcap E$, since $\forall R.D \rightarrow_S \emptyset$ and $E \rightarrow_S E$. However, also $\exists R^-. \top \sqsubseteq D$ has been added to \mathcal{T}' , thus the approximation is $\mathcal{T}' = \{C \sqsubseteq E, \exists R^-. \top \sqsubseteq D\}$. It is easy to see that $\mathcal{T}' \models \mathcal{T}$ holds. In case the input TBox is $\{C \sqsubseteq (\forall R.D) \sqcup E\}$, then the following rewriting is possible: $(\forall R.D) \sqcup E \rightarrow_S \forall R.D \rightarrow_S \emptyset$. Since `approximate_axiom` will reject axioms with $right_side' = \emptyset$, the approximation is simply $\mathcal{T}' = \{\exists R^-. \top \sqsubseteq D\}$. Another possibility is of course $\mathcal{T}' = \{C \sqsubseteq E\}$, according to the \sqcup -rule.

Proposition 3 *Let $\mathcal{T}' = \text{approximate}(\mathcal{T})$ for a SHI TBox \mathcal{T} . Then, \mathcal{T}' is a DL-Lite approximation of \mathcal{T} .*

Proof 2 (Sketch) *This can be shown by induction on \mathcal{T} .*

It is easy to check that $\text{approximate_axiom}(C \sqsubseteq D, \mathcal{T}') \cup \mathcal{T} \models \{C \sqsubseteq D\} \cup \mathcal{T}$ for every \mathcal{T}' by monotonicity and by definition of `approximate_axiom`.

An Example Approximation If `approximate` is applied to the example TBox \mathcal{T} , then after the preprocessing only the following non-*DL-Lite* axioms remain which thus have to be approximated:

$$\{ \top \sqsubseteq \forall publicationAuthor^-. (book \sqcup conferencePaper), \top \sqsubseteq \forall teacherOf.course, \\ person \sqcap \exists takesCourse.course \sqsubseteq student, student \sqsubseteq \exists takesCourse.course, \\ person \sqcap \exists headOf.department \sqsubseteq chair, chair \sqsubseteq \exists headOf.department \}$$

One possible *DL-Lite* approximation \mathcal{T}' is:

$$\{ chair \sqsubseteq \exists headOf. \top, \exists headOf. \top \sqsubseteq chair, chair \sqsubseteq person, \\ professor \sqsubseteq faculty, book \sqsubseteq publication, faculty \sqsubseteq person, \\ graduateStudent \sqsubseteq student, student \sqsubseteq \exists takesCourse. \top, student \sqsubseteq person, \\ \exists teacherOf^-. \top \sqsubseteq course, \exists takesCourse. \top \sqsubseteq student, \exists teacherOf. \top \sqsubseteq faculty, \\ \exists publicationAuthor. \top \sqsubseteq book, \\ \exists headOf^-. \top \sqsubseteq department, \exists takesCourse^-. \top \sqsubseteq course \}$$

If this \mathcal{T}' is used for retrieval on the example ABox, then no unsound answers to atomic instance retrieval queries are delivered. Thus, the approximation is *perfect* for every atomic concept of the ABox. Of course, this depends on the ABox. Note that $p2$ is a *chair* instance, as required, since $\exists headOf. \top \sqsubseteq chair \in \mathcal{T}'$.

However, there also exist imperfect approximations. On the one hand, there are many \mathcal{T}' 's containing incoherent concepts, or even unsatisfiable \mathcal{T}' 's. In the latter case, the ABox is definitely inconsistent (unsatisfiable), and in the former case, inconsistency of the ABox is likely (if the ABox contains instances of incoherent concepts). From a logical perspective, an instance retrieval query performed on an inconsistent ABox returns the set of *all* ABox individuals (since, from an inconsistent theory, everything follows). So, such a query answer is still complete.⁴ On the other

⁴Of course, a DL reasoner typically does not permit ABox retrieval on inconsistent ABoxes.

hand, an example for an unsound approximation is given by \mathcal{T}' , if $\top \sqsubseteq \forall publicationAuthor^{\neg} . (book \sqcup conferencePaper)$ is approximated to $\exists publicationAuthor . \top \sqsubseteq conferencePaper$ instead of $\exists publicationAuthor . \top \sqsubseteq book$. Then, $b1$ will be a false answer to the *conferencePaper* instance retrieval query.

Even for the simple example TBox we get 757 coherent approximations (there are a few hundred thousand consistent approximations containing incoherent concepts); w.r.t. retrieval, there is one perfect approximation (see above), and the worst coherent approximation has an average failure of 2.5 individuals which means that an atomic instance retrieval query, in the average, returns 2.5 false query answers.

An Implementation of the Approximation Algorithm We have eliminated the non-determinism in the `approximate` algorithm by implementing it in a depth-first (backtracking) search algorithm. Thus, for a given *axiom*, `approximate_axiom(axiom)` returns a set of *candidate axioms*, representing possible approximations of *axiom*. Each axiom thus represents a state in the search space, whose branching factor is given by the number of its candidate approximation axioms.

In principle, the number of possible approximations is truly astronomic for larger TBoxes. Consider a TBox with 500 axioms to approximate, in which each axiom can be approximated in three different ways – the number of atoms in the universe is $10^{80} \approx 3^{167.6722}$ and thus tiny compared to the 3^{500} nodes in this search space. Thus, clever heuristics are needed to guide the search. Since, in principle, one is only interested in coherent approximations (containing only one incoherent concept name, \perp), it is a good idea to prune a path in the search tree as soon as more than one incoherent concept is discovered in the partial \mathcal{T}' . Of course, this requires a TBox coherence check by the DL reasoner (RACERPRO) at each step. This would be a good use case for *incremental* reasoning. In order to filter out candidate axioms which are too specific, RACERPRO is used as well.

It may not be possible to compute a coherent approximation at all. In this case, an incoherent TBox is wanted which at least leaves the ABox satisfiable (but even this may be impossible), or contains only a minimal number of incoherent concepts.

Sometimes it is possible to compute more than one approximation. Even if each computed approximation is unsound for retrieval on an actual ABox, it is good to have a multitude of approximations available, since their retrieval results can be intersected. Even if no – w.r.t. an actual ABox – perfect approximation is among the computed approximations, this intersected answer set may be perfect for some concept C on this ABox.

4 Island-Based Instance Retrieval – The Refine Step

In the following section we discuss how to post-filter individuals, which were obtained by the Filter Step before. The original algorithm was proposed in [13] for the DL *ALCHI*.

The idea for the refine step is that only a subset of role and concept assertions is necessary/used to perform instance checking for a particular given individual a and a given concept C . The approach undertaken here is to identify role assertions which can be used during the application of a tableau algorithm for instance checking (note that $\langle \mathcal{T}, \mathcal{A} \rangle \models^? a : C$ can be reduced to checking whether $\langle \mathcal{T}, \mathcal{A} \cup \{a : \neg C\} \rangle$ is unsatisfiable via a tableau algorithm).

For *ALCHI*-ontologies the important tableau-rule is the \forall -rule, since it is the only rule, which makes use of existing role assertions in the ABox. We know that \forall -constraints cannot come directly from the ABox, since we only allow for concept assertions of the kind $a : C$, where C is an atomic concept or its negation. Thus, \forall -constraints can only be derived from the TBox. Unfortunately, the shape of TBox axioms, i.e. $C \sqsubseteq D$, does not allow to easily read off “possible” \forall -constraints.

This is due to the implicit presence of negation in the concept C . We propose some kind of normal form (*Shallow Normal Form*), which allows for extraction of a superset of \forall -constraints, which can possibly be used in a tableau algorithm:

Definition 4 A concept description C is in *Shallow Normal Form (SNF)*, if it has the shape $C = C_1 \sqcup C_2 \sqcup \dots \sqcup C_n$, s.t. each C_i is either

- an atomic concept,
- a negated atomic concept,
- an \exists -constraint $\exists R.D$, s.t. D is an arbitrary concept description in negation normal form
- a \forall -constraint $\forall R.D$, s.t. D is an arbitrary concept description in negation normal form

Please note that we do not enforce anything on concepts “hidden” behind \forall/\exists -constraints, but that they are in negation normal form. Thus the name *shallow* normal form.

Lemma 1 Each GCI $C \sqsubseteq D$ can be converted into a set S of equivalent concept descriptions in SNF. Here, equivalent means that $C \sqsubseteq D$ is unsatisfiable iff the conjunction of the formulas in S is unsatisfiable.

For the details of the transformation please refer to [13]. For completeness, we only provide an example for \mathcal{T}_{EX} from Example 1 in Shallow Normal Form. The TBox \mathcal{T}_{EX} in SNF is as follows: $Shallow(\mathcal{T}_{EX}) =$

$$\left\{ \begin{array}{l} \neg Chair \sqcup \exists headOf.Department, \neg Chair \sqcup Person, \\ \forall headOf.\neg Department \sqcup \neg Person \sqcup Chair, \\ \neg Professor \sqcup Faculty, \neg Book \sqcup Publication, \\ \neg GraduateStudent \sqcup Student, \neg Student \sqcup Person, \\ \neg Student \sqcup \exists takesCourse.Course, \\ \neg Person \sqcup \forall takesCourse.\neg Course \sqcup Student, \\ \forall teacherOf.Course, \forall teacherOf.\perp \sqcup Faculty, \neg Faculty \sqcup Person, \\ \forall publicationAuthor^-(Book \sqcup ConferencePaper) \end{array} \right\}$$

Having the TBox in Shallow Normal form, we can read off a superset of possible \forall -constraints from the terminology.

Managing \forall -constraints in *ALCHI*-ontologies.

Definition 5 A \forall -info structure for the conceptual part of TBox \mathcal{T} is a function $f_{\mathcal{T}}^{\forall} : N_R \rightarrow \mathcal{P}(N_C \cup \{\neg A \mid A \in N_C\} \cup \{\perp\}) \cup \{*\}$, s.t. N_C (N_R) is a set of atomic concepts (roles) used in \mathcal{T} .

Informally speaking, the function $f_{\mathcal{T}}^{\forall}$ is used to manage the \forall -constraints, i.e. the function assigns to each role name in N_R one of the following entries:

- \emptyset , if we know that there is no \forall -constraint for R in \mathcal{T}
- a subset S of $N_C \cup \{\neg A \mid A \in N_C\} \cup \{\perp\}$, s.t. there is no other concept but those in S , which occur \forall -bound on R in \mathcal{T}
- $*$, if there are arbitrary complex \forall -constraints on role R in \mathcal{T} , but we don't give additional information on the structure of these constraints.

Definition 6 A \forall -info structure for ontology $\mathcal{O} = \langle \mathcal{T}, \mathcal{A} \rangle$ is a function $f_{\mathcal{O}}^{\forall} : N_R \rightarrow \mathcal{P}(N_C \cup \{\neg A \mid A \in N_C\} \cup \{\perp\}) \cup \{*\}$, s.t.

$$f_{\mathcal{O}}^{\forall}(R) = \begin{cases} * & \text{if } \exists S \in N_R. \sqsubseteq_R(R, S) \wedge (f_{\mathcal{T}}^{\forall}(S) = *) \\ \bigcup_{R \sqsubseteq_{\mathcal{R}} S} f_{\mathcal{T}}^{\forall}(S) & \text{else} \end{cases}$$

Let us look at our ontology \mathcal{O}_{EX} again to give an example for a \forall -info structure. The \forall -info structure for ontology \mathcal{O}_{EX} is as follows:

$$f_{\mathcal{O}_{EX}}^{\forall}(R) = \begin{cases} \{\neg Department\} & \text{if } R = headOf \\ \{\neg Course\} & \text{if } R = takesCourse \\ \{\perp, Course\} & \text{if } R = teacherOf \\ * & \text{if } R = publication.Author^- \\ \{\} & \text{else} \end{cases}$$

Lemma 2 *Given an ontology $\mathcal{O} = \langle \mathcal{T}, \mathcal{A} \rangle$, the following holds: if $f_{\mathcal{O}}^{\forall}(R) \neq *$, then for each valid tableau proof P of \mathcal{O} and for each application of a \forall -rule (on R and subconcept C) in P , we have that $C \in f_{\mathcal{O}}^{\forall}(R)$.*

Thus, a \forall -info structure for an ontology \mathcal{O} helps us to determine which concepts are (worst-case) propagated over role assertions in an ABox.

Island identification. Given the notions in the previous part, in the following we will derive means to rewrite an ontology \mathcal{O} , s.t. inconsistency is preserved, i.e, $INC(\mathcal{O})$ if and only if $INC(\mathcal{O}_{rewritten})$. Inconsistency tests are important for instance checking, since it holds that $INC(\langle \mathcal{T}, \mathcal{A} \cup \{a : \neg C\} \rangle)$ if and only if $\langle \mathcal{T}, \mathcal{A} \rangle \models a : C$. The following definition of \mathcal{O} -separability is used to determine the importance of role assertions in a given ABox. Informally speaking, the idea is that \mathcal{O} -separable assertions will never be used to propagate “complex and new information” (see below) via role assertions.

Definition 7 *Given an ontology $\mathcal{O} = \langle \mathcal{T}, \mathcal{A} \rangle$, a role assertion $R(a, b)$ is called \mathcal{O} -separable, if we have $INC(\mathcal{O})$ iff $INC(\langle \mathcal{T}, \mathcal{A}_2 \rangle)$, where*

$$\mathcal{A}_2 = \mathcal{A} \setminus \{R(a, b)\} \cup \{R(a, i_1), R(i_2, b)\} \cup \{i_1 : C \mid b : C \in \mathcal{A}\} \cup \{i_2 : C \mid a : C \in \mathcal{A}\},$$

s.t. i_1 and i_2 are fresh individual names.

Lemma 3 *Given an ontology $\mathcal{O} = \langle \mathcal{T}, \mathcal{A} \rangle$ and a role assertion $R(a, b) \in \mathcal{A}$, it holds that $R(a, b)$ is \mathcal{O} -separable, if we have*

1. For each $C \in f_{\mathcal{O}}^{\forall}(R)$
 - (a) $C = \perp$ or
 - (b) we can find a concept description $D \in \{E \mid b : E \in \mathcal{A}\}$, s.t. we have $D \sqsubseteq_{\mathcal{T}} C$
 - (c) we have $b : nnf(\neg C) \in \mathcal{A}$
2. For each $C \in f_{\mathcal{O}}^{\forall}(R^-)$
 - (a) $C = \perp$ or
 - (b) we can find a concept description $D \in \{E \mid a : E \in \mathcal{A}\}$, s.t. we have $D \sqsubseteq_{\mathcal{T}} C$
 - (c) we have $a : nnf(\neg C) \in \mathcal{A}$

Let us consider an example for \mathcal{O} -separability w.r.t. Example 1:

For instance the ABox assertion $teacherOf(p2, c3)$ in Example 1 is \mathcal{O}_{EX} -separable, since we have

- $f_{\mathcal{O}_{EX}}^{\forall}(teacherOf) = \{\perp, Course\}$ and $c3 : Course \in \mathcal{A}_{EX}$
- $f_{\mathcal{O}_{EX}}^{\forall}(teacherOf^-) = \{\}$

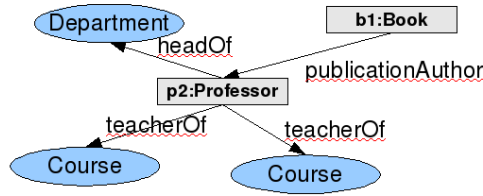


Figure 2: Example island for individual $b1$ in \mathcal{O}_{EX}

The extraction of islands for instance checking in ontology \mathcal{O} , given an individual a , is now straightforward. From an individual a one just follows each non- \mathcal{O} -separable role assertion in the original ABox, until at most \mathcal{O} -separable role assertions are left. For the details of this algorithm please refer to [13]. In Figure 2 we show the island computed for individual $b1$ from Example 1. Please recall that $b1$ potentially was a false answer to the *conferencePaper* instance retrieval query. It is easy to see that the computed island does not entail $b1 : conferencePaper$ and thus $b1$ can be eliminated from the set of candidates. In Figure 2 it is also easy to see that $p2$ can be verified to be an instance of *Chair*, by only taking into account the corresponding island, since the *headOf*-edge was preserved during the computation of the island.

Extension from DL *ALCHI* to DL *SHI*. Transitive roles can be easily read off from the TBox by additionally taking into account the role hierarchy. Then, whenever we want to compute the island for an individual w.r.t. DL *SHI*, then we have to additionally “follow” all transitive role assertions. This proposal for the extension to DL *SHI* is quite straight-forward and we do not prove it here.

5 Test Results

5.1 Evaluation of the Abox partitioning algorithm

First, we have evaluated our proposal for island reasoning. We have implemented the proposed partitioning algorithms in Java. For ontology access we used a Java interface and implementation for the Ontology Web Language, called OWLAPI⁵. Given the OWLAPI interface, implementing the above algorithms was straightforward.

Ontology	—Inclusions—	—Equivalences—	— N_{RN} —	Time for analysis (ms)
LUBM[10]	75	6	25	4
Cyc[1]	43541	2	4853	93
GODaily[3]	28997	0	1	103
Galen[2]	3388	699	413	55
Pizza[4]	57	2	7	3

Figure 3: Extraction of \forall -info structures for several ontologies

First, we investigate, whether it is feasible to convert the TBox of a given ontology into Shallow Normal Form and extract the \forall -info structure $f_{\mathcal{O}}^{\forall}$. For this purpose we have selected five arbitrary ontologies. For lack of space please refer to the given references for a detailed description of the ontologies. The results of our investigation are shown in Figure 3. We think the numbers indicate

⁵<http://owlapi.sourceforge.net/>

that extracting a \forall -info structure for most ontologies should be feasible. Even ontologies whose TBox is considered large for current description logic systems do perform quite well.

Second, we look at the average size of extracted islands. Our tests w.r.t. LUBM are quite encouraging. The average size of an island is 29 nodes. The actual island size depends on the chosen individual. Some preliminary statistics on example islands are given in Figure 4. Please note that the island size does not depend on the number of universities, i.e. for LUBM our approach promises quite good scalability.

The big island for *Dep0.Uni0/FullProfessor7* can be explained as follows. LUBM imposes several \forall -constraints on the role *headOf*. Due to the definitions of *Chair*, *Dean* and *Director*, the atomic concept descriptions \neg *Department*, \neg *College* and \neg *Program* can be propagated via role *headOf*. Since all individuals with an incoming *headOf*-edge are of type *Department*, we have to take them completely into account for building the islands. If there were further constraints in the TBox, e.g. disjointness of *Department*, *College* and *Program*, we could further reduce the size of islands for individuals of type *FullProfessor*.

Individual	Island size	Time for island computation
<i>Dep0.Uni0/GraduateStudent128</i>	9	0 ms
<i>Dep0.Uni0/Publication2</i>	4	1 ms
<i>Dep0.Uni0/FullProfessor7</i>	93	2 ms
<i>Dep0.Uni0/Course4</i>	37	0 ms

Figure 4: Statistics for islands in LUBM

5.2 Evaluation of the whole approach

We have performed an initial evaluation of the complete algorithms on a version of the “Athletics Event Ontology” (AEO) developed in the BOEMIE project. Using RACERPRO, we have transformed this OWL ontology into a DL ontology (= TBox, ABox). The utilized TBox DL is *ALCHf*. It contains 1061 axioms which are already in *DL-Lite*, and 499 axioms which have to be approximated to *DL-Lite*. AEO also contains some so-called number restrictions, which we simply approximate to functional roles in *DL-Lite_F* (since only $\leq_1 R$ concepts appear).

The ABox of the AEO version we used is rather small – it only contains 138 individuals (266 concept assertions plus 70 role assertions = 336 assertions). We have chosen this ABox since some interesting reasoning is required in order to retrieve the instances of the concept *HighJump* (similar to the *chair* example, but over 2 role fillers).

As illustrated previously, it is very demanding to approximate a TBox with 499 axioms. Unfortunately, we were not successful to compute a coherent approximation of this AEO ontology in reasonable time. Better heuristics are needed here. The reason for this is a massive number of *disjointness axioms*; e.g., axioms of the form $A \sqsubseteq \neg B$, $A \sqsubseteq \neg C$, \dots . Additionally, `get_K4_closure` introduces another 1110 additional axioms.

We have thus simplified AEO substantially by removing all disjointness axioms and ignoring transitivity (so `get_K4_closure` adds no axioms). With this version, a coherent approximation could be computed within 5 minutes. These simplifications do not affect retrieval. In the average, it returns 0,984 false instances for a concept (w.r.t. to the original AEO). The original AEO contains one instance of *HighJump*, and no instances of *SprintCompetition*. The approximated version is perfect for *HighJump*, but delivers 7 wrong *SprintCompetitions*. The *HighJump* instance is in fact also a (false) *SprintCompetition* here. Thus, 7 islands were computed by the partitioning method, ranging in size from 7 to 45 assertions. The average island contains 33.75 assertions. So, in the average, only one tenth of the assertions from the original ABox have to be loaded in order to verify or falsify the candidates for *HighJump* and *SprintCompetition*. Each instance test

requires ≈ 180 msec per candidate, thus, after $\approx 1,440$ seconds the candidate individuals have been refined. Some additional time is needed to compute the islands. Computation of an islands needs milliseconds only (for such small islands).

Since this ABox was rather small, we expanded the ABox artificially by a factor of 500. We thus created an ABox which contained 500 copies of the original ABox, simply by prefixing the individuals with numbers (0 to 499). We then connected these 500 separated ABox parts using some new artificial role assertions, resulting in a connected ABox, containing 415330 assertions. The ABox still fits into main memory, because otherwise we could not have performed this experiment (the secondary memory-access is not yet realized). The ABox consistency check already needs 3,5 minutes on this ABox now; retrieval requires ≈ 34 seconds (for each concepts). As expected, from the approximated version of the TBox, we got 500 *HighJump* instances, and 3500 *SprintCompetitions*. As expected, the newly introduced artificial role assertions connecting the 500 copies have no influence on the size of the islands. Thus, the average islands size is still 33.75; that this is only 0.0008126068 % of the whole ABox. However, now 4000 candidate tests have to be performed, which will require ≈ 14 minutes of RACERPRO reasoning time. Additional time for accessing and loading from secondary memory etc. (since also the island partitioning has to work on secondary memory in the future) is taken⁶.

Acknowledgments. This work has also been partially supported by the EU-funded project BOEMIE (Bootstrapping Ontology Evolution with Multimedia Information Extraction, IST-FP6-027538).

Finally, we would like to thank the BOEMIE team at STS, and especially Irma Sofia Espinosa Peraldi for providing us with the AEO ontology and the *HighJump* example.

References

- [1] CYC. <http://www.opencyc.org>.
- [2] Galen. <http://www.co-ode.org/galen/>.
- [3] GODaily. <http://www.geneontology.org>.
- [4] Pizza ontology. <http://www.co-ode.org/resources/tutorials/>.
- [5] A. Acciarri, D. Calvanese, G. De Giacomo, D. Lembo, M. Lenzerini, M. Palmieri, and R. Rosati. Quonto: Querying ontologies. In *Proc. of the 20th Nat. Conf. on Artificial Intelligence (AAAI 2005)*, 2005.
- [6] R. Arndt, R. Troncy, S. Staab, L. Hardman, and M. Vacura. Comm: Designing a well-founded multimedia ontology for the web. In *In 6th International Semantic Web Conference (ISWC'2007)*, 2007.
- [7] F. Baader, D. Calvanese, D. L. McGuinness, D. Nardi, and P. F. Patel-Schneider. *The Description Logic Handbook*. Cambridge University Press, New York, NY, USA, 2007.
- [8] D. Calvanese, G. De Giacomo, D. Lembo, M. Lenzerini, and R. Rosati. Tractable Reasoning and Efficient Query Answering in Description Logics: The DL-Lite Family. 2007.

⁶The test files and results can be downloaded from <http://www.sts.tu-harburg.de/~mi.wessel/download/boemie-experiment.zip>.

- [9] Yuanbo Guo and Jeff Heflin. A Scalable Approach for Partitioning OWL Knowledge Bases. In *SSWS 2006*, Athens, GA, USA, November 2006.
- [10] Yuanbo Guo, Zhengxiang Pan, and Jeff Heflin. Lubm: A benchmark for owl knowledge base systems. *J. Web Sem.*, 3(2-3):158–182, 2005. <http://swat.cse.lehigh.edu/projects/lubm/index.htm>.
- [11] V. Haarslev, R. Möller, and M. Wessel. RacerPro User’s Guide and Reference Manual Version 1.9.1, May 2007.
- [12] R. Möller, V. Haarslev, and M. Wessel. On the scalability of description logic instance retrieval. In Chr. Freksa and M. Kohlhase, editors, *29. Deutsche Jahrestagung für Künstliche Intelligenz*, Lecture Notes in Artificial Intelligence. Springer Verlag, 2006.
- [13] Sebastian Wandelt and Ralf Moeller. Island reasoning for alchi ontologies. In *Proceedings of the 5th International Conference on Formal Ontology in Information Systems (FOIS-04)*. IOS Press, 2008.