

# Efficient reasoning on large SHIN Aboxes in relational databases

Julian Dolby<sup>1</sup>, Achille Fokoue<sup>1</sup>, Aditya Kalyanpur<sup>1</sup>, Li Ma<sup>2</sup>, Chintan Patel<sup>3</sup>,  
Edith Schonberg<sup>1</sup>, Kavitha Srinivas<sup>1</sup>, and Xingzhi Sun<sup>2</sup>

<sup>1</sup> IBM Watson Research Center, P.O. Box 704, Yorktown Heights, NY 10598, USA  
dolby, achille, adityakal, ediths, ksrinivs@us.ibm.com

<sup>2</sup> IBM China Research Lab, Beijing 100094, China  
malli, sunxingz@cn.ibm.com

<sup>3</sup> Columbia University Medical Center  
chintan.patel@dbmi.columbia.edu

**Abstract.** As applications based on semantic web technologies enter the mainstream, there is a need to provide highly efficient ontology reasoning over large Aboxes. However, achieving sufficient scalability is still a challenge, especially for expressive ontologies. In this paper, we present a hybrid approach which combines a fast, incomplete reasoning algorithm with a slower complete reasoning algorithm to handle the more expressive features of DL. Our approach works for *SHLN*. We demonstrate the effectiveness of this approach on large datasets (30-60 million assertions), including a clinical-trial patient matching application, where we show significant performance gains (an average of 15 mins per query compared to 100 mins) without sacrificing completeness or expressivity.  
**keywords:** Reasoning, Description Logic, Ontology.

## 1 Introduction

As applications based on semantic web technologies enter the mainstream, there is a need to provide highly efficient ontology reasoning over large Aboxes. However, achieving sufficient scalability is still a challenge. DL reasoning is intractable in the worst case. In [1], we reported on the use of expressive reasoning for matching patient records to clinical trial criteria. While the system was able to successfully reason on 240,269 patient records, a knowledge base with 59 million Abox and 33,561 Tbox assertions, the execution time was prohibitive. In some cases, the system took hours to respond.

The expressivity of the patient knowledge base was *ALCH*, so expensive reasoning was needed to be complete. However, most typical queries were simple, and could have been answered faster with a less expensive reasoner. A high cost was paid by all queries to support rarer complex queries. In this paper, we present a hybrid approach, that combines a fast, incomplete reasoning algorithm with a slower complete reasoning algorithm to handle the more expressive features of DL. In this way, we were able to dramatically lower the cost of typical simple queries, without losing the ability to answer more complex queries.

An interesting feature of our technique is that *any sound and incomplete algorithm* may be used in the first phase to quickly find as many solutions as possible to the query. The key novelty in the approach is a mechanism to incorporate these solutions into a slower, complete reasoning algorithm for *SHLN*, providing much better performance characteristics overall, without sacrificing completeness or expressivity. This approach can be described as self-adjusting, since the reasoner dynamically defaults to the expensive complete algorithm only when deeper inferencing is actually required. On large datasets (30-60 million assertions), this hybrid approach provides significant performance gains (an average of 15 mins per query on the 60 million dataset compared to 100 mins) without sacrificing completeness or expressivity.

At its core, this hybrid approach builds on the summarization and refinement techniques we described earlier to perform sound and complete reasoning on large Aboxes in relational databases [2] [3]. Briefly, this technique applies a standard tableaux algorithm on a *summary Abox*  $\mathcal{A}'$  rather than the original Abox  $\mathcal{A}$  to answer queries. A summary Abox is created by aggregating individuals which are members of the same concepts, so when any given individual is tested in the summary Abox, all individuals mapped to the summary individual are effectively tested at the same time. For a tested individual  $s$  in  $\mathcal{A}'$ , if the summary is found to be consistent, then we know that all individuals mapped to that summary individual  $s$  are not solutions. But if the summary is found to be inconsistent, it is possible that either (a) a subset of individuals mapped to the summarized individual  $s$  are instances of the query or (b) the inconsistency is a spurious effect of the summarization. We determine the answer through *refinement*, which selectively expands the summary Abox to make it more precise. Refinement is an iterative process that partitions the set of individuals mapped to a single summary individual based on the common edges they have in the original Abox, and remaps each partition to a new summary individual. The iteration ends when either the expanded summary is consistent, or it can be shown that all individuals mapped to the tested summary individual are solutions. Significantly, convergence on the solution is based only on the structure of the refined summary, without testing individuals in  $\mathcal{A}$ . In practice, the scalability of this algorithm is limited by the number of refinement steps that are needed. Refinement is performed by database join operations, which become expensive when the database is large.

The key insight of our hybrid approach is that the solutions from the sound and incomplete reasoner can be used as a partitioning function for refinement instead of partitioning based on common edges, as described in our earlier work. This effectively removes the obvious solutions from the summary Abox. If the sound and incomplete reasoning algorithm finds all solutions, there will be no solutions left in the summary Abox after this first refinement, so the algorithm will converge very quickly. Any remaining inconsistencies are spurious, and can be resolved in one or a few refinement steps. If the sound and incomplete algorithm finds only some of the solutions, then the refinement process will find the rest of the solutions with fewer refinement steps.

Our key contributions in this paper are as follows: (a) we develop a fast, sound but incomplete algorithm based on query expansion, and describe how to incorporate solutions from this and other such techniques into a sound and complete hybrid algorithm for reasoning over large expressive Aboxes, and (b) we demonstrate its effectiveness in providing performance gains (from 100 minutes per query to 15 minutes per query) on expressive Aboxes with 60 million assertions.

## 2 Related Work

There have been efforts in the semantic web community to define less expressive subsets of OWL-DL for which reasoning is tractable. The EL-family of languages [4] is one such example, for which classification can be done in polynomial time. To take advantage of this fact, various query answering algorithms for EL have been proposed (e.g. [5]). Another example is the DL-Lite family [6], for which conjunctive query answering is expressible as a first-order logic formula (and hence an SQL query) over the Abox stored in a relational database. The QuOnto algorithm [6] is a sound and complete query expansion algorithm for DL-Lite.

Our query expansion algorithm described in Section 6 is not significantly novel. It is similar in spirit to the EL and DL-Lite query expansion approaches, with some differences, namely: (i) instead of using an EL reasoner to compute additional subclasses during the normalization process (as in [5]), we use a sound and complete OWL-DL reasoner (Pellet) which enables us to discover more entailments outside of EL; (ii) we use a datalog reasoner to compute *same-as*-individual inferences (considering functional properties) and transitive closure for transitive properties that exist in the ABox.

Furthermore, a key point is that *any* query answering algorithm for a subset of OWL can be plugged into our sound and complete hybrid OWL-DL reasoning system. When it is known that the optimization is complete based on the underlying logic of the KB<sup>4</sup> and the manner in which it is implemented, fallback to our refinement strategy is not necessary. Otherwise, the refinement process will find any remaining solutions.

## 3 Background

Query answering in expressive DLs can be reduced to consistency detection. For instance, assume that we want to find all instances of the concept  $C$ . To answer this query, each individual  $a$  is tested by adding the assertion  $a : \neg C$  to the Abox, and checking the new Abox for consistency. If the Abox is inconsistent, then  $a$  is an instance of  $C$ . For large Aboxes, this approach will clearly not scale. Therefore, in our previous work [3], [7], we modify this approach to perform tableau reasoning on a summarized version of the Abox rather than the original

---

<sup>4</sup> Checking whether the logic falls in EL or DL-Lite is a matter of syntactic checking of the KB axioms which can be done easily

Abox. Formally, an Abox  $\mathcal{A}'$  is a summary Abox of a  $\mathcal{SHIN}$  Abox  $\mathcal{A}$  if there is a mapping function  $\mathbf{f}$  that satisfies the following constraints<sup>5</sup>:

- (1) if  $a : C \in \mathcal{A}$  then  $\mathbf{f}(a) : C \in \mathcal{A}'$
- (2) if  $R(a, b) \in \mathcal{A}$  then  $R(\mathbf{f}(a), \mathbf{f}(b)) \in \mathcal{A}'$
- (3) if  $a \neq b \in \mathcal{A}$  then  $\mathbf{f}(a) \neq \mathbf{f}(b) \in \mathcal{A}'$

If the summary Abox  $\mathcal{A}'$  obtained by applying the mapping function  $\mathbf{f}$  to  $\mathcal{A}$  is consistent w.r.t. a given Tbox  $\mathcal{T}$  and a Rbox  $\mathcal{R}$ , then  $\mathcal{A}$  is consistent w.r.t.  $\mathcal{T}$  and  $\mathcal{R}$ . However, the converse does not hold. In the case of an inconsistent summary, we use a process of iterative refinement to make the summary more precise, to the point where we can conclude that an inconsistent summary  $\mathcal{A}'$  reflects a real inconsistency in the actual Abox  $\mathcal{A}$ . Refinement is a process by which only the part of the summary that gives rise to the inconsistency is made more precise, while preserving the summary Abox properties (1)-(3). To pinpoint the portion of the summary that gives rise to the inconsistency, we focus on the *justification* for the inconsistency, where a justification is a minimal set of assertions which, when taken together, imply a logical contradiction.

We define refinement for a summary individual  $s$  in a justification  $\mathcal{J}$  as a partition where individuals mapped to  $s$  are partitioned based on which edges in  $\mathcal{J}$  each individual actually has. More specifically:

$$key(a, \mathcal{J}) \equiv \left\{ \begin{array}{l} \mathbf{f}(a) = s \wedge \\ R(t, s) \in \mathcal{J} \wedge \\ \exists b \text{ in } \mathcal{A} \text{ s.t.} \\ R(b, a) \in \mathcal{A} \wedge \\ \mathbf{f}(b) = t \end{array} \right\} \cup \left\{ \begin{array}{l} \mathbf{f}(a) = s \wedge \\ R(s, t) \in \mathcal{J} \wedge \\ \exists b \text{ in } \mathcal{A} \text{ s.t.} \\ R(a, b) \in \mathcal{A} \wedge \\ \mathbf{f}(b) = t \end{array} \right\}$$

Since an individual may be mapped to a summary individual that is in multiple overlapping justifications, we define:

$$key^*(a) = \bigcup_{\{\mathcal{J} | a \in \mathcal{J}\}} key(a, \mathcal{J})$$

In a *refinement step* that refines  $s$  in  $\mathcal{A}'$ , new individuals  $s_1 \dots s_k$  replace  $s$  in  $\mathcal{A}'$ , where there are  $k$  unique key sets  $key^*(a)$ , for all  $a$  in  $\mathcal{A}$  such that  $\mathbf{f}(a) = s$ . Individuals  $a$  and  $b$  in  $\mathcal{A}$  mapped to  $s$  in  $\mathcal{A}'$  are partitioned correspondingly, that is,  $\mathbf{f}(a) = \mathbf{f}(b)$  after the refinement step iff  $key^*(a) = key^*(b)$  before the refinement step.

In principle, in the presence of many justifications involving overlapping sets of nodes, the union of the keys could become very large. In practice, we have not observed this across the various knowledge bases we have evaluated, even for ones that do contain overlapping justifications.

If all individuals in  $\mathcal{A}$  mapped to a summary individual  $s$  have the same key w.r.t.  $\mathcal{J}$ , then it must be the case that they have all the edges in the justification

<sup>5</sup> We assume without loss of generality that  $\mathcal{A}$  does not contain an assertion of the form  $a \doteq b$

and hence  $s$  is precise w.r.t.  $\mathcal{J}$ . If a justification is precise, we can conclude that all individuals in  $\mathcal{A}$  mapped to the tested individual in the justification are solutions to the query. In the worst case, iterative refinement can expand a summary Abox into the original Abox, but in practice, we conclude on precise justifications with many individuals mapped to each summary node in the justification.

Our implementation of summarization and refinement in a system called SHER is in terms of RDBMS operations to allow the system to scale to large data sets. However, the iterative process of summarization and refinement is expensive, because (a) it requires expensive join operations on all role assertions in the Abox  $\mathcal{A}$  to define the  $key(a)$ , as well as expensive join operations of role assertions with type assertions to rebuild the summary, and (b) it requires several consistency checks to find the many sources of inconsistencies for each summary that gets built. For large knowledge bases with multiple ways in which one can derive a solution to the query, this becomes a serious performance bottleneck.

## 4 A Sample Knowledge Base

We illustrate our techniques with the sample knowledge base (Tbox  $\mathcal{T}$ , the Rbox  $\mathcal{R}$  and the Abox  $\mathcal{A}$ ) in Figures 1 and 2. This example is a small subset of the UOBM [8] benchmark that we use in our evaluation. To form the summary Abox for Figure 2, the individuals  $a$  and  $b$  are mapped to a single summary individual  $w$  with a concept set of *Woman*, and the individuals  $f$ ,  $g$  and  $j$  are mapped to another summary individual  $p$  with a concept set of *Person*. The summary Abox is shown in the Figure 3.

$\mathcal{T}$  assertions:

- (1)  $WomanCollege \sqsubseteq \forall hasStudent.Woman$
- (2)  $\top \sqsubseteq \leq 1 isTaughtBy$

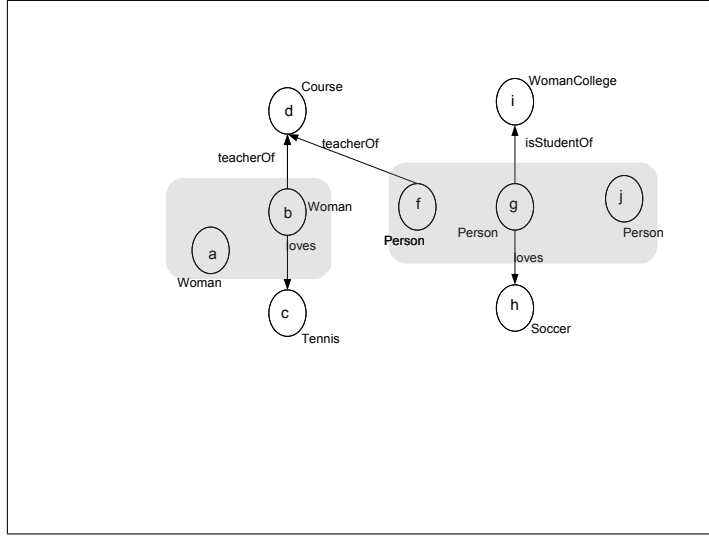
$\mathcal{R}$  assertions:

- (1)  $loves \sqsubseteq likes$
- (2)  $isStudentOf$  is inverse of  $hasStudent$
- (3)  $teacherOf$  is inverse of  $isTaughtBy$

**Fig. 1.** Example  $\mathcal{T}$ ,  $\mathcal{R}$

Consider the query *WomanWithHobby*, which is defined as  $Woman \sqcap \geq 1 likes$ . There are three solutions. The individual  $b$  is a solution because  $loves \sqsubseteq likes$ . The individual  $f$  is a solution because the course  $d$  can be taught by only one *Person*, and so  $f$  and  $b$  will be identified with each other during reasoning. Finally,  $g$  is a solution, since  $isStudentOf(g, WomenCollege)$  implies that  $g$  is a *Woman*.

Figure 3 shows the entire refinement process for answering this query:

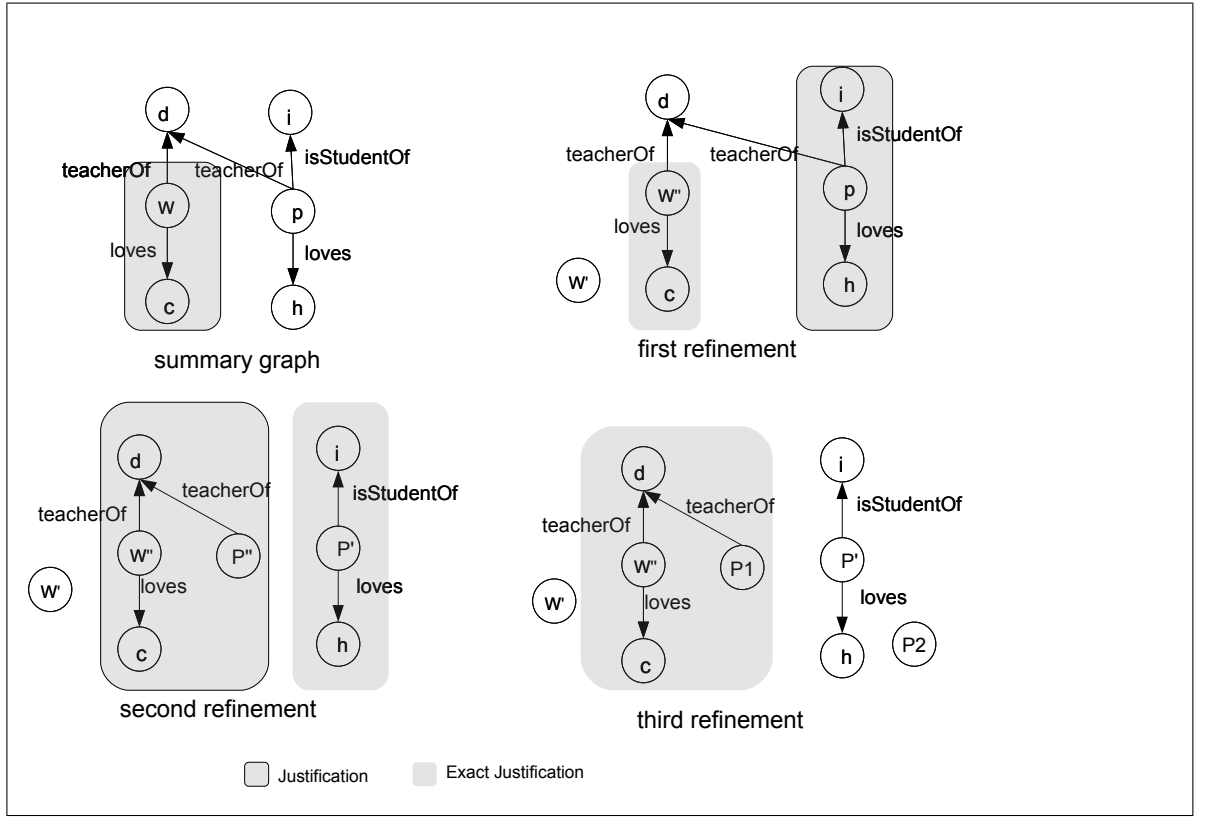


**Fig. 2.** Example  $\mathcal{A}$

- (1) Refine  $w$  by splitting it into two nodes  $w'$  which has  $a$  mapped to it, and  $w''$  which has  $b$  mapped to it.
- (2) Refine  $p$  by splitting it into two nodes  $p'$  which has  $g$  mapped to it, and  $p''$  which has  $f$  and  $j$  mapped to it.
- (3) Refine  $p''$  further, by splitting it into nodes  $p_1$  which has  $f$  mapped to it, and  $p_2$  which has  $j$  mapped to it.

We explain these steps in more detail. First,  $\neg WomanWithHobby$  is added to a tested summary individual  $w$ . The resulting Abox is inconsistent, and a justification  $\mathcal{J}$  contains the assertions:  $w : Woman$ ,  $loves \sqsubseteq likes$ , and  $loves(w, c)$ . For refinement, we target the summary individuals in  $\mathcal{J}$ , which are  $w$  and  $c$ . Refinement makes a justification  $\mathcal{J}$  *precise*, that is, it partitions the individuals mapped to the summary node  $w$  into a new set of summary nodes to reflect the fact that not all individuals in  $\mathcal{A}$  mapped to  $w$  have the  $loves(w, c)$  in  $\mathcal{J}$ . The summary individual  $w$  is therefore split into two new summary nodes,  $w'$  that has individuals with no  $loves(w, c)$  mapped to it (e.g.,  $a$ ), and  $w''$  that has individuals with  $loves(w, c)$  mapped to it (e.g.,  $b$ ). This new refined Abox is still inconsistent, with a new justification  $\mathcal{J}$  which contains the individuals  $w''$  and  $c$ . Refinement of  $w''$  or  $c$  however is no longer possible, because every individual in  $\mathcal{A}$  that is mapped to  $w''$  also has the  $loves(c, \cdot)$  and every individual mapped to  $c$  has the same edge (here  $c$  is the same as the summary node  $c$ ). At this point, the justification  $\mathcal{J}$  is precise, in that it cannot be refined further, and we conclude that all individuals in  $\mathcal{A}$  mapped to  $w''$  are solutions to the query.

For the second step,  $\neg WomanWithHobby$  is added to a tested summary individual  $p$ . The resulting Abox is inconsistent, and this time there is the jus-



**Fig. 3.** Refinement Steps for Example

tification:  $isStudentOf(p, i)$ ,  $loves \sqsubseteq likes$ , and  $loves(p, h)$ , combined with the axiom  $WomanCollege \sqsubseteq \forall hasStudent.Woman$ . The result of the second refinement is shown in Figure 3. After this refinement, the subgraph containing  $p'$  is still inconsistent, and  $p'$  is not refinable. Therefore, all individuals in  $\mathcal{A}$  mapped to  $p'$ , namely  $g$ , are solutions.

There is one final justification which is refinable:  $teacherOf(p'', d)$ ,  $teacherOf(w'', d)$ ,  $w'' : Woman$ ,  $loves \sqsubseteq likes$ ,  $loves(w'', c)$ , and  $\top \sqsubseteq \leq 1isTaughtBy$ . After the third refinement step, we conclude that  $f$  mapped to  $P1$  is a solution.

On large knowledge bases, the cost of each additional refinement is significant, so it is critical to reduce the number of refinements. We show in the next sections how our hybrid reasoning approach can reduce the number of refinements for this example.

## 5 Hybrid Algorithm

The key idea to reducing refinement iterations is to (a) quickly find solutions to the query, (b) refine the summary to isolate these solutions into new summary individuals, and (c) ignore these individuals for the rest of the refinement process. We find solutions quickly by using a sound and incomplete reasoning algorithm which does a form of query expansion described in Section 6. We point out that other reasoner implementations (such as QuOnto) for less expressive logics may also be plugged into this technique.

To illustrate the overall idea in terms of our example in Figure 2, we expand our query *WomanWithHobby* into the query  $WomanWithHobby(x) \sqcup (Woman(x) \sqcap likes(x, y)) \sqcup (Woman(x) \sqcap loves(x, y))$ . This query matches all pairs of individuals in the Abox bound to both  $x$  and  $y$ , namely the pair  $(b, c)$ , and this constitutes our set of known bindings. Our next step is to refine the summary Abox, so that the individuals in the solution, namely  $b$  and  $c$ , are mapped to distinct new summary individuals. We do this by refining the summary Abox in a manner similar to that described in Section 3; the only difference is that we now partition the Abox individuals according to whether they were bound to any variable in the query or not, rather than according to key sets. That is,  $\mathbf{f}(a) = \mathbf{f}(b)$  after the refinement step iff  $a$  and  $b$  are mapped to the same summary node before the refinement step and either both or neither  $a$  and  $b$  are individuals in the set of known bindings. Our algorithm keeps track of the subset of known bindings that actually are answers to the query, which is just  $b$  in this case. Next, consistency checking is applied to this refined summary, and any remaining inconsistencies are resolved using the standard iterative refinement and summarization process described in [3].

This approach has a nice property: in cases where the incomplete step actually does find all solutions and the summary itself is consistent, the complete reasoning step may simply be a single consistency check on the refined summary. Since there are no more solutions to be found, the only possible causes of inconsistency are spurious inconsistencies, which are the result of our summarization technique. In practice, we find that the incomplete step captures all solutions on most complex queries on most realistic datasets. This optimization therefore significantly reduces the number of refinements and makes query answering practical for large Aboxes.

One non-obvious part of the hybrid algorithm is that it is important to partition out *all* individuals that are bound to any variable in the query, and not just the individuals that are actual solutions to the query. To illustrate why this is the case, consider a simple Abox shown in Figure 4 with 3 patients ( $q, r, s$ ) who each have an associated lab event  $(l, m, n)$ , and each event indicates a presence of organisms of different types, where  $x, y$ , and  $z$  indicate individuals with organisms of type  $X, Y$  and  $Z$ , respectively. The summary Abox, as shown in the Figure will contain one patient individual  $p$ , which has  $q, r$  and  $s$  mapped to it, one lab event individual  $e$  which has  $l, m$  and  $n$  mapped to it, and 3 individual nodes for organisms  $x, y$ , and  $z$ . Consider a realistic query, which is to find all patients who have a laboratory event which shows the presence of



the organism X. As shown in the Figure 4, if a summary is built with only the solution individual  $q$  partitioned out, then it will contain spurious inconsistencies which will cause unnecessary refinement. To avoid this issue, we should not only partition out the solution individual  $q$  from  $p$ , but also other individuals bound to other variables in the query, which in our example would be  $l$  and  $x$ .

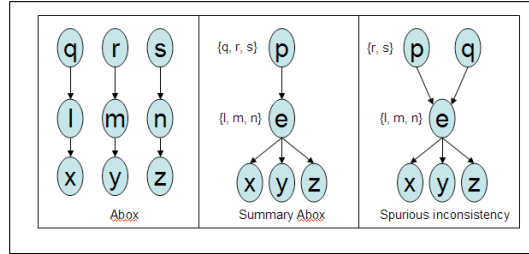


Fig. 4. Partitioning Complexity

The pseudo-code for our overall algorithm is shown in the function `ConjunctiveQuery` in Figure 5.

## 6 Query Expansion

Our sound but incomplete reasoning algorithm is based on the well-known recursive query expansion technique suggested in the EL [5] and DL-Lite [6] solutions. As discussed earlier, our approach differs in the following ways: (a) we refer to an OWL-DL reasoner (Pellet) for computing subclasses of a concept when performing the expansion, (b) we have an ABox pre-processing step that uses a datalog reasoner to compute transitive relations in the ABox and *same-as* inferences between ABox individuals due to functional property assertions. The *same* individuals are used to expand query solutions, i.e, if individual  $a$  is found to be a solution to the SQL query generated by query expansion, and  $sameAs(a,b)$  is inferred by the datalog reasoner, we add  $b$  to the solution set.

For any given query  $a : C$ , we recursively traverse the definitions and subclasses of the concept  $C$ . For our sample query  $x : WomanWithHobby$ , we first generate a union of SQL select statements which signify all the possible ways in which this query can be expanded. The first disjunct in the union matches individuals of *WomanWithHobby* directly,  $rd\!f : type(x, WomanWithHobby)$ . In this case, however, the *WomanWithHobby* type does not appear in the ABox, and so we drop this disjunct. Next we would generate disjuncts to match individuals that are in subclasses of *WomanWithHobby*, but in this case there are no subclasses (checked by calling a standard DL reasoner). We then add any complex subclasses of *WomanWithHobby* which can be inferred syntactically. In our example, we have one such obvious subclass because *WomanWithHobby*

```

Function:ConjunctiveQuery
Input: Conjunctive Query  $CQ: C_i(x) \wedge ..R_j(x, y)$ 
/* Get incomplete answers from sound but incomplete algorithm, which
   can be translated to SQL */
sqlQuery  $\leftarrow$  BuildQuery( $CQ$ );
/* Get the bindings for all variables in the expanded query, both
   distinguished and non-distinguished variables */
result  $\leftarrow$  execute(sqlQuery);
/* Build filtered summary for query answering, which is the basic
   summary Abox */
sum  $\leftarrow$  BuildSummary( $A, CQ$ );
/* Separate the bindings for distinguished variables  $x_{dist}$  from
   bindings for existentially quantified variables */
sqlsolutions  $\leftarrow$  getBindings(result,  $x_{dist}$ );
others  $\leftarrow$   $\bigcup_{v \in vars(result) - x_{dist}}$  getBindings(result,  $v$ );
/* Refine summary based on solutions found from SQL */
sum  $\leftarrow$  refineSummaryFromSolutions(sum, sqlsolutions  $\cup$  others) ;
/* Find all summary nodes in new summary which have sqlSolutions
   mapped to them */
sumSolutions  $\leftarrow$  getSummaryNodesForSQLSolutions(sum, sqlSolutions);
/* complete query answering, using refined summary */
restsolutions  $\leftarrow$  solveQuery(sum, allnodes - sumSolutions);
return sqlsolutions  $\cup$  restsolutions

```

**Fig. 5.** Overall optimized complete query algorithm

is defined as equivalent to  $Woman \sqcap \geq 1likes$ . The expansion process now recursively continues and we expand this complex concept into a select statement which is a disjunction of conjuncts; i.e., the selection must satisfy the two conditions  $rdf : type(x, Woman)$  and  $likes(x, y)$ , or alternatively, satisfy the two conditions  $rdf : type(x, Woman)$  and  $loves(x, y)$ , since  $likes$  has a subproperty  $loves$ . These queries are applied against an Abox that has been processed to include all edges materialized from the application of all deterministic merger and transitivity rules.

One technical challenge in query expansion in general is keeping the query relatively simple, especially when given very large Tboxes with deep subclass and subproperty hierarchies. Our approach to this problem was to eliminate forms of query expansion if the concept or role did not appear in the ABox. We therefore maintained a simple cache of all roles and concepts that appeared in the ABox, and limited our expansion to only these concepts and roles.

## 7 Evaluation

We evaluated our technique on two knowledge bases: the first is a real-world knowledge base, and real queries of clinical data that we had used in previ-

ous work[1], and the second is the UOBM benchmark[8]. Our experiments were conducted on a 2-way 2.4GHz AMD Dual Core Opteron system with 16GB of memory running Linux, and we used IBM DB2 V9.1 as our database. Our Java processes were given a maximum heap size of 8GB for clinical data, and 4GB for UOBM.

### 7.1 Clinical trials dataset

In prior work [1], we reported on the use of expressive reasoning for matching of patient records on clinical trials. The 1 year anonymized patient dataset we used contained electronic medical records from Columbia University for 240,269 patients with 22,561 Tbox subclass assertions, 26 million type assertions, and 33 million role assertions. The 22,561 Tbox subclass assertions are a subset of the a larger Tbox which combines SNOMED with Columbia’s local taxonomy called MED for a total of 523,368 concepts. For details of the partitioning algorithm used to define the subset see [1]. Although the expressivity of the SNOMED version we used falls in the EL fragment of DL, the expressivity needed to reason on the knowledge base is  $\mathcal{ALCH}$ . This is because we have type assertions in the Abox which includes assertions of the type  $\forall R.-C$ , where the concept  $C$  is itself defined in terms of a subclass or equivalence axiom. As a concrete example, for a given patient, and a specific radiology episode for the patient, the presence of *ColonNeoplasm* may be ruled out. *ColonNeoplasm* has complex definitions in SNOMED (e.g.,  $ColonNeoplasm \equiv \exists AssociatedMorphology.Neoplasm \sqcap \exists FindingSite.Colon \sqcap ColonDisorder$ ). We selected the 9 clinical trials we evaluated in our earlier work which are shown Table 1. Table 2 shows the DL version of the queries, in the order shown in Table 1. For query *NCT00001162*, the results shown are for the union of 7 different disorders, only 4 of which are illustrated in Table 2.

ClinicalTrials.gov ID	Description
<i>NCT00084266</i>	Patients with MRSA
<i>NCT00288808</i>	Patients on warfarin
<i>NCT00393341</i>	Patients with breast neoplasm
<i>NCT00419978</i>	Patients with colon neoplasm
<i>NCT00304382</i>	Patients with pneumococcal pneumonia where source specimen is blood or sputum
<i>NCT00304889</i>	Patients on metronidazole
<i>NCT00001162</i>	Patients with acute amebiasis, giardiasis, cyclosporiasis or strongloides...
<i>NCT00298870</i>	Patients on steroids or cyclosporine
<i>NCT00419068</i>	Patients on corticosteroid or cytotoxic agent

**Table 1.** Clinical Trial Requirements Evaluated

Table 3 shows the queries, the number of patients matched to the queries, the time to process the queries in minutes, the time in minutes for our hybrid

DL Query
$\exists \text{associatedObservation.MRSA}$
$\exists \text{associatedObservation.}$ $\exists \text{roleGroup.}$ $\exists \text{administeredSubstance.}$ $\exists \text{roleGroup.}\exists \text{hasActiveIngredient.Warfarin}$
$\exists \text{associatedObservation.BreastNeoplasm}$
$\exists \text{associatedObservation.ColonNeoplasm}$
$\exists \text{associatedObservation.}$ $\left( \begin{array}{l} \text{PneumococcalPneumonia} \\ \sqcap \\ \exists \text{hasSpecimenSource.Blood} \sqcup \text{Sputum} \end{array} \right)$
$\exists \text{associatedObservation.}$ $\exists \text{roleGroup.}$ $\exists \text{administeredSubstance.}$ $\exists \text{roleGroup.}\exists \text{hasActiveIngredient.Metronidazole}$
$\exists \text{associatedObservation.}$ $\left( \begin{array}{l} \text{acuteamebiasis} \sqcup \\ \text{giardiasis} \sqcup \\ \text{cyclosporiasis} \sqcup \\ \text{strongloides} \sqcup \\ \dots \end{array} \right)$
$\exists \text{associatedObservation.}$ $\exists \text{roleGroup.}$ $\exists \text{administeredSubstance.}$ $\exists \text{roleGroup.}\exists \text{hasActiveIngredient.cyclosporine} \sqcup \text{steroids}$
$\exists \text{associatedObservation.}$ $\exists \text{roleGroup.}$ $\exists \text{administeredSubstance.}$ $\exists \text{roleGroup.}\exists \text{hasActiveIngredient.corticosteroid} \sqcup \text{cytotoxicAgent}$

**Table 2.** DL Queries for Evaluated Clinical Trials

approach (HTime), the time in minutes for our previous approach (Time), the number of refinements with our hybrid approach (HRefinements) and the number of refinements with our previous approach (Refinements). As can be seen from the table, the hybrid approach reduced the number of refinements to 1 in all cases, which reflects the refinement needed to check that there are no additional solutions after the incomplete algorithm has completed (The one case where 0 refinements occurred was because for that specific query, our expressivity checker decided that no refinement was needed given the specific filtered Abox that was built for the query and the Tbox.) The hybrid approach improved our overall query times from 100.4 mins on average with a standard deviation of 113.7, to 15.6, with a standard deviation of 3.5. This is not surprising, given that the entire variability in query answering in our previous approach was due to the number of refinements.

Query	Matched Patients	Time (m)	HTime (m)	Refinements	HRefinements
<i>NCT00084266</i>	1052	68.9	17.8	6	1
<i>NCT00288808</i>	3127	63.8	11.6	5	0
<i>NCT00393341</i>	74	26.4	12.1	2	1
<i>NCT00419978</i>	164	31.8	12.4	3	1
<i>NCT00304382</i>	107	56.4	15.1	8	1
<i>NCT00304889</i>	2	61.4	20.7	3	1
<i>NCT00001162</i>	1357	370.8	13.5	58	1
<i>NCT00298870</i>	5555	145.5	19.3	8	1
<i>NCT00419068</i>	4794	78.8	17.5	5	1

**Table 3.** Patient Matches for Trial DL Queries for 240,269 Patients

## 7.2 UOBM

We evaluated our approach on the UOBM benchmark, modified to *SHIN* expressivity. This was done by adding a new concept to correspond to each of the nominals in the dataset (e.g. *SwimmingClass* for *Swimming*), adding a type assertion for each nominal (e.g., *Swimming : SwimmingClass*), and changing any of the references to nominals in the Tbox to point to the class. Currently, we have evaluated membership query answering, and we tested one membership query for each concept in the benchmark<sup>6</sup>, comparing the hybrid approach with our prior techniques. We report results for UOBM size 100—with roughly 7.8 million type assertions and 22.4 million role assertions—and UOBM size 150—with about 11.7 million type assertions and 33.5 million role assertions. The queries naturally fall into three categories:

**empty** Concepts that have no instances in the Abox.

**simple** Concepts that have only simple solutions (i.e. reasoning does not require iterative refinement because the justification viewed as a graph does not have path lengths greater than 1).

**complex** Concepts that have complex solutions (i.e. reasoning requires iterative refinement because the justification viewed as a graph has path lengths greater than 1).

We expect the hybrid approach to benefit only the third category of queries. One complication is that the summary Abox for the UOBM benchmark has a spurious inconsistency induced by the summarization process, so all membership query answering require 2 passes of refinement in order to make the summary consistent.

Table 4 shows results for the 3 query categories for UOBM sizes 100 and 150. The first three columns list the UOBM dataset size, the category of query, and how many such queries there are. For both sizes and each query category, we report the average and standard deviation for the query time and the number of passes of refinement. For both datasets, we timed out queries that took

<sup>6</sup> That is, all classes in the original benchmark. The extra classes introduced by our transformation to *SHIN* are ignored.

Size	Category	Count	Time (seconds)				Refinement			
			Original		Hybrid		Original		Hybrid	
			Average	Stdev	Average	Stdev	Average	Stdev	Average	Stdev
100	empty	11	214	37	214	19	2	0	2	0
100	simple	43	255	83	265	47	2	0	2	0
100	complex	14	891*	386*	377	105	14*	11*	3	.3
150	empty	11	301	35	347	45	2	0	2	0
150	simple	43	340	88	416	85	2	0	2	0
150	complex	14	1368*	508*	647	198	14*	11*	3	.3

**Table 4.** Results for UOBM Membership Queries for sizes 100 and 150

longer than 30 minutes to complete; the timeouts occurred on both the 100 size (1 timeout) and the 150 size (6 timeouts) for the original approach. Hence, those averages and standard deviations are significant underestimates, and so are marked with a \* in the table.

As one might expect, there is some overhead for executing the incomplete query, and so the simpler queries actually show some slowdown in the hybrid approach. However, the results do indicate that our hybrid approach greatly reduces the time for the complex queries, which were the most expensive ones with our previous approach. In fact, for all but one query, the incomplete reasoning algorithm found all the solutions. The one query which was the outlier, `GraduateCourse`, required propagation from a universal restriction for reasoning, which was not accounted for by our incomplete algorithm. In this case, we proceeded to find the answer through our prior complete reasoning algorithm.

## 8 Conclusion and Future Work

We have developed an efficient, scalable query answering system for large expressive ABoxes. The hybrid approach proposed in this paper combines our novel summarization and refinement technology to do sound and complete OWL-DL reasoning with any incomplete reasoning implementation (possibly for a subset of OWL).

We have used our hybrid solution to build a web-based semantic search engine for biomedical literature, known as *Anatomy Lens*, details of which can be found in [9]. Anatomy Lens has indexed 300 million RDF triples dealing with PubMed data, and utilizes ontological information from three large biomedical ontologies (Gene ontology, Foundational Model of Anatomy, and MeSH), doing query answering in a few seconds. Performing web-time reasoning for such a large expressive dataset would not have been possible without our approach.

We plan to further optimize our query expansion algorithm by pruning irrelevant queries considering the summary ABox, and to continue to explore the use of SHER in real world semantic web applications.

## References

1. C.Patel, J.Cimino, J.Dolby, A.Fokoue, A.Kershenbaum, L.Ma, E.Schonberg, K.Srinivas: Matching patient records to clinical trials. Proc. of the Int. Semantic Web Conf. (ISWC 2007) (2007)
2. A.Fokoue, A.Kershenbaum, L.Ma, E.Schonberg, K.Srinivas: The summary abox: Cutting ontologies down to size. Proc. of the Int. Semantic Web Conf. (ISWC 2006) (2006) 136–145
3. Dolby, J., A.Fokoue, Kalyanpur, A., A.Kershenbaum, L.Ma, E.Schonberg, K.Srinivas: Scalable semantic retrieval through summarization and refinement. Proc. of the 22nd Conf. on Artificial Intelligence (AAAI 2007) (2007)
4. Baader, F., Brandt, S., Lutz, C.: Pushing the  $\mathcal{EL}$  envelope. In: Proceedings of the Nineteenth International Joint Conference on Artificial Intelligence IJCAI-05, Edinburgh, UK, Morgan-Kaufmann Publishers (2005)
5. Rosati, R.: On conjunctive query answering in EL, CEUR Electronic Workshop Proceedings (2007)
6. Calvanese, D., Giacomo, G.D., Lembo, D., Lenzerini, M., Rosati, R.: DL-lite: Tractable description logics for ontologies. Proc. of AAAI (2005)
7. Dolby, J., Fokoue, A., Kalyanpur, A., Ma, L., Schonberg, E., Srinivas, K., Sun, X.: Scalable grounded conjunctive query evaluation over large and expressive knowledge bases. In: ISWC '08: Proceedings of the 7th International Conference on The Semantic Web, Berlin, Heidelberg, Springer-Verlag (2008) 403–418
8. Ma, L., Yang, Y., Qiu, Z., Xie, G., Pan, Y.: Towards a complete owl ontology benchmark. In: Proc. of the third European Semantic Web Conf.(ESWC 2006). (2006) 124–139
9. Dolby, J., Fokoue, A., Kalyanpur, A., Schonberg, E., Srinivas, K.: Scalable highly expressive reasoner (sher). In: Journal of Web Semantics, (accepted), <http://dx.doi.org/10.1016/j.websem.2009.05.002> (2009)