

Classifying and querying very large taxonomies with bit-vector encoding

Hassan Ait-Kaci & Samir Amir

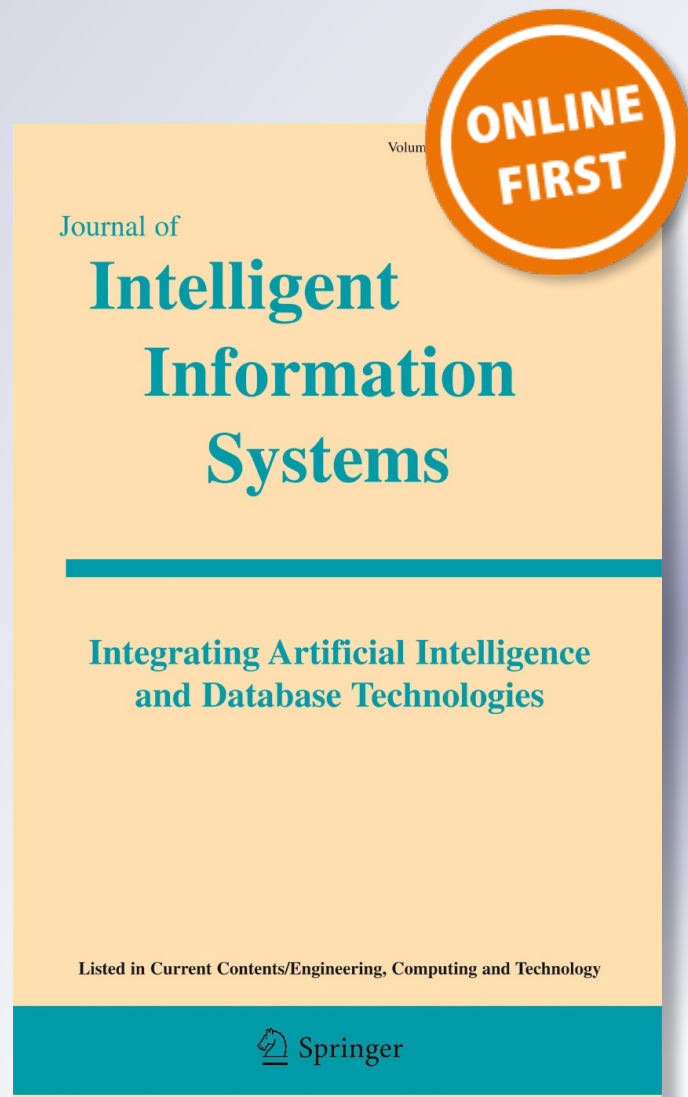
Journal of Intelligent Information Systems

Integrating Artificial Intelligence and Database Technologies

ISSN 0925-9902

J Intell Inf Syst

DOI 10.1007/s10844-015-0383-2



Your article is protected by copyright and all rights are held exclusively by Springer Science +Business Media New York. This e-offprint is for personal use only and shall not be self-archived in electronic repositories. If you wish to self-archive your article, please use the accepted manuscript version for posting on your own website. You may further deposit the accepted manuscript version in any repository, provided it is only made publicly available 12 months after official publication or later and provided acknowledgement is given to the original source of publication and a link is inserted to the published article on Springer's website. The link must be accompanied by the following text: "The final publication is available at link.springer.com".

Classifying and querying very large taxonomies with bit-vector encoding

Hassan Aït-Kaci¹ · Samir Amir¹

Received: 24 February 2015 / Revised: 13 July 2015 / Accepted: 31 August 2015
© Springer Science+Business Media New York 2015

Abstract In this article, we address the question of how efficiently Semantic Web (SW) reasoners perform in processing (classifying and querying) taxonomies of enormous size and whether it is possible to improve on existing implementations. We use a bit-vector encoding technique to implement taxonomic concept classification and Boolean-query answering. We describe the technique we have used, which achieves high performance, and discuss implementation issues. We compare the performance of our implementation with those of the best existing SW reasoning systems over several very large taxonomies under the exact same conditions for so-called TBox reasoning. The results show that our system is among the best for concept classification and several orders-of-magnitude more efficient in terms of response time for query answering. We present these results in detail and comment them. We also discuss pragmatic issues such as cycle detection and decoding.

Keywords Binary encoding · Taxonomic reasoning · Query optimization · Semantic web

1 Introduction

Let us first describe the context, research problem, challenges, and added-value contribution of the work reported in this article.

✉ Hassan Aït-Kaci
hassan.ait-kaci@univ-lyon1.fr
Samir Amir
samir.amir@univ-lyon1.fr

¹ LIRIS—Département Informatique, Université Claude Bernard Lyon 1, Villeurbanne, France

There exist many and varied formal systems purporting to express and use knowledge for inference. Despite this variety, they all share a common trait: knowledge is always organized into an “is-a” conceptual taxonomy where a concept denotes the set of all its instances. In such a taxonomy, a concept C_1 is deemed a subconcept of a concept C_2 whenever the set denoted by C_1 is a subset of the set denoted by C_2 . This is what “ C_1 is-a C_2 ” means. For example, “*employee is-a person*” means that all instances of the concept *employee* are also instances of the concept *person*. Hence, reasoning performance in all such systems will need to rely in a crucial way upon the performance of the underlying taxonomic reasoning. In particular, it will be necessary to identify all the maximal concepts that are subconcepts of two given concepts, those that are minimal superconcepts of two given concepts, or those that are incompatible with a given concept. These operations correspond respectively to concept conjunction (*i.e.*, intersection of the denoted sets), concept disjunction (*i.e.*, union of the denoted sets), and concept negation (*i.e.*, complement of the denoted set). Therefore, basic Propositional Algebra is central to *all* taxonomic reasoners.

The issue, then, is to enable the efficient evaluation of propositional expressions involving propositional symbols representing partially-ordered concepts making up these conceptual taxonomies. However, conceptual taxonomies can be of enormous size, making such evaluation a challenge to perform efficiently.

In this context, we address two topics: (1) robust and scalable taxonomic reasoning; and, (2) evaluation of semantic web technologies for such reasoning. Regarding the first topic, we demonstrate how a method for taxonomic reasoning based on bit-vector encoding we have implemented is both robust and scalable on very large taxonomies derived from real-life ontologies. As for the second topic, we measured the performance of our system and compared it with those of the best existing SW reasoning systems over several very large taxonomies under the exact same conditions for so-called TBox reasoning.¹ The results show that our system is among the best for concept classification and several orders-of-magnitude more efficient in terms of response time for query answering. We present these results in detail and comment them.

The rest of this paper is organized as follows. In Section 2, we overview the taxonomic reasoners used in our experiments: Section 2.1 is a brief survey of existing systems, and Section 2.2 describes our own method. Section 3 contains the main contribution of this article: Section 3.1 describes our experiments’ setup; Section 3.2 presents the results of these experiments; in Section 3.3, we discuss these results. To make this report complete, Section 4 describes two important pragmatic issues concerning encoding and decoding: detecting potential cycles in a poset formed by declaring a taxonomy (Section 4.1), and how to go from binary vectors back to sorts (Section 4.2). In Section 5, further work is discussed. We conclude in Section 6 with a summary of our contribution and perspectives. We added an appendix giving a formal specification of a space-efficient representation for extremely large binary codes.

¹In Semantic Web lingo, a Knowledge Base (KB) is defined as a formal ontology consisting of two parts (or “boxes”): (1) a *Terminological Box* (abbreviated as *TBox*); and, (2) an *Assertional Box* (abbreviated as *ABox*). The TBox contains the formal axioms that define the structure and semantic properties of the actual instance data; which instance data constitute the ABox. In Database lingo, the TBox corresponds to the schema and the ABox to the actual data.

2 Semantic-web reasoning

2.1 The state of the art

In this section, we give a brief description of the SW reasoners that we have used for our comparative experiments. Note that we have limited our selection to systems that are full-fledged *reasoners*, and not just *classifiers*. This is because our interest goes beyond concept classification and includes Boolean query answering as well. This rules out systems such as [ELK](#)² (Kazakov et al. 2011), [CEL](#)³ (Baader et al. 2006), [CB](#)⁴ (Kazakov 2009), *etc.*, that do not support query answering.

We retrieved and installed the following SW reasoners: [FaCT++](#),⁵ [HerMiT](#),⁶ [Pellet](#),⁷ [TrOWL](#),⁸ [Racerpro](#),⁹ and, [SnoRocket](#).¹⁰

[FaCT++](#) (Fast Classification of Terminologies) is a reasoner developed at the University of Manchester (Tsarkov and Horrocks 2006). It is based on the Description Logic fragment *SHOIQ* (Horrocks and Sattler 2007). It is implemented in C++ as a deductive tableau Manna and Waldinger (1991) adapted to the specifics of this logic. It is claimed to use a wide range of heuristic optimizations. [FaCT++](#) provides TBox reasoning (subsumption, satisfiability, classification) and partial support for ABox processing (retrieval).

[HerMiT](#) is also a reasoner for a (slight extension) of the Description Logic fragment *SHOIQ* (called *SHOIQ+*) (Shearer et al. 2008). It is based upon hypertableau reasoning, an optimized version of tableau reasoning (Motik et al. 2009). It purports to provide a faster process for classifying ontologies. The main optimization of hypertableau vs. tableau that it tries to minimize nondeterminism in the treatment of disjunctions and is more memory-efficient. [HerMiT](#) provides TBox reasoning, with the ability of checking the consistency of an ontology and inferring implicit relationships between concepts.

[Pellet](#) is a free open-source Java-based reasoner (Sirin et al. 2007). It, too, is based on the tableau algorithm and supports the Description Logic fragment *SHOIN(D)*. It provides TBox reasoning (subsumption, satisfiability, and classification) and ABox reasoning (retrieval, conjunctive query answering). It uses many optimization techniques and supports entailment checks and ABox querying through its interface.

[TrOWL](#) (Tractable reasoning infrastructure for OWL 2) was developed at the University of Aberdeen (Thomas et al. 2010). This is a system that starts by transforming an ontology from [OWL-DL](#) to [OWL-QL](#) Fikes et al. (2004) in order to classify it in polynomial time. Under this transformation, conjunctive query answering and consistency checking remain the same as for [OWL-DL](#). In addition, [TrOWL](#) can generate a database schema for storing normalized representations of [OWL-QL](#) ontologies.

²www.cs.ox.ac.uk/isg/tools/ELK/

³code.google.com/p/cel/

⁴code.google.com/p/cb-reasoner/

⁵owl.cs.manchester.ac.uk/fact++/

⁶www.hermit-reasoner.com/

⁷clarkparsia.com/pellet/

⁸trowl.eu/

⁹www.racer-systems.com/products/racerpro/

¹⁰research.ict.csiro.au/software/snorocket

Fig. 1 Java code for “in-place”
Warshall’s algorithm

```

int n = SORTS.size();
for (int k = 0; k < n; k++)
  for (int i = 0; i < n; i++)
    for (int j = 0; j < n; j++)
      if (!SORTS[i].code.get(j))
        SORTS[i].code.set(j,
          SORTS[i].code.get(k)
          &&
          SORTS[k].code.get(j));

```

Racerpro is a commercial version of **RACER** (Renamed ABoxes and Concept Expression Reasoner) (Haarslev and Möller 2001; Haarslev et al. 2011). It implements a reasoner for the description logic *SHIQ*. RACER provides both TBox and ABox reasoning. It supports all the optimizations of FaCT++ as well as new techniques for dealing with number restrictions and ABoxes.

Snorocket Lawley and Bousquet (2010) was proposed as a high-performance implementation of a polynomial-time classification algorithm for the lightweight Description Logic \mathcal{EL} (Baader et al. 2005).¹¹ It was primarily meant to be optimized for classifying **SNOMED CT**. It can process only conjunctive queries.

2.2 Our method

In this section, we give a self-contained summary of the method we have implemented in order to measure its performance for classification of bare taxonomies and query answering of Boolean queries.

Our method is an implementation in Java of a technique described in (Aït-Kaci et al. 1989). It consists in representing the elements of a taxonomy (*i.e.*, an arbitrary poset) as bit vectors. Thus, each element has a code (a bit vector) carrying a “1” in the position corresponding to the index of any other elements that it subsumes. In this manner, the three Boolean operations on sorts are readily and efficiently performed as their corresponding operations on bit-vectors. However, for this to be possible, these bit vectors must be encoded as the reflexive transitive closure of the “is-a” relation obtained from subsort declarations.

How to compute such a closure has been well-known—*e.g.*, the Warshall-Strassen method using clever matrix multiplication tricks (Strassen 1969; Fischer and Meyer 1971; Coppersmith and Winograd 1990). However, for a poset of n elements, this method has quite a large time complexity—even with the best known algorithm to date, it is $\mathcal{O}(n^{2.23727})$ (Stothers 2010; Williams 2011).¹² In fact, in practice, the straightforward $\mathcal{O}(n^3)$ in-place multiplication method known as **Warshall’s Algorithm** (Warshall 1962; Warren 1975) is used in most cases.

Figure 1 gives Java code for Warshall’s algorithm performed “in-place” on the binary codes of sorts stored in an array `SORTS`. The array `SORTS` contains the set of sorts as bit-set objects. A bit-set object has a field named “code” which is a bit vector. The class of bit vectors is endowed with a method `get(int)` that returns the value of its bit in the

¹¹Description Logics in the \mathcal{EL} -family are weaker versions that provide existential roles ($\exists r.C$) but no universal roles ($\forall r.C$).

¹²To the best of our knowledge, this is the latest best bound as of 2011. However, these algorithms are not implementable due to prohibitive size of constants. For more recent work on parallelizing Strassen’s algorithm, see (Ballard et al. 2014). This, however, requires special hardware (GPGPUs).

specified position as a `boolean`, and a method `set(int, boolean)` that sets its bit in the specified position to the specified Boolean value.

Now, while Warshall's algorithm may be viable for relatively small posets, it simply becomes unusable for posets of the size of the taxonomies we are considering.

Note, however, that transitive-closure methods need pay such a high performance cost only due to the fact that they are devised for arbitrary graphs. But concept taxonomies are *not* arbitrary graphs. Namely, a necessary condition for a set of partially-ordered concepts to be semantically consistent is that its graph must be *acyclic*. Thus, a consistent taxonomy must be a directed-acyclic graph (or *dag*) with a least element (\perp) and a highest element (\top). In Ait-Kaci et al. (1989), it is shown that for such a dag, an $\mathcal{O}(n)$ transitive-closure algorithm exists and is proven correct. This method is described as Algorithm 1.

Algorithm 1 Taxonomy Classification Algorithm

```

1: procedure CLASSIFY
2:    $L \leftarrow \text{Parents}(\perp)$ ;
3:   while  $L \neq \emptyset$  do
4:     for all  $x \in L$  do
5:        $x.\text{code} \leftarrow 2^{x.\text{index}} \vee \bigvee_{y \in \text{Children}(x)} y.\text{code}$ ;
6:        $x.\text{coded} \leftarrow \text{true}$ ;
7:     end for
8:      $L \leftarrow \bigcup_{x \in L} \text{Parents}(x)$ ;
9:     for all  $x \in L$  do
10:      if  $\exists y \in \text{Children}(x)$  and  $\neg y.\text{coded}$  then
11:         $L \leftarrow L - \{x\}$ ;
12:      end if
13:    end for
14:  end while
15: end procedure

```

The procedure CLASSIFY assumes that each sort s is provided with a set denoted as “*Parents(s)*” and a set denoted as “*Children(s)*” (*viz.*, from the “is-a” declarations of members of the taxonomy). Each sort is an object that has a field (called “*code*”), which is its bit vector representation (initialized to be all zeroes). A sort also has an integer field (called “*index*”) which is unique per sort (*viz.*, its index in the taxonomy). Finally, a sort also has a Boolean field (called “*coded*”), which denotes whether or not this sort has been encoded yet (it is initially set to **false**).

Algorithm 1 computes the reflexive-transitive closure of a the “is-a” relation on the set of sorts comprising a taxonomy (Ait-Kaci et al. 1989).¹³ This algorithm can be explained as follows. It proceeds layer by layer, starting with the parents of \perp (*i.e.*, the minimal sorts in the taxonomy) [Line 2], assigning a code to each element in the current layer to be the bitwise `or` of its children and also setting the bit in its index position [Line 5]. Each time an element is encoded, it is marked to be so by setting its *coded* flag to **true** [Line 6]. Then, a new layer is computed from the current one as the union of all its parents [Line 8] from which any sort that has at least one child not encoded is removed [Line 11]. Indeed, by construction, such sorts can always be reached later. This proceeds until an empty layer is obtained [Line 3]—which is when all sorts have been encoded.

¹³*op. cit.*, pages 125–126.

This algorithm's main loop [Lines 3–14] clearly visits each sort exactly once, and is thus linear. The auxiliary computation of the next layer [Lines 8–13] has comparatively marginal cost as it can be made efficient using constant access-time data structures for the sets of parents and children, making set operations on them negligible. Also proceeding bottom up has a clear performance advantage for dags such as most concept taxonomies, where sorts tend to have many less parents than they have children.

It is this method that we have implemented and tested on very large taxonomies, and compared it with the best SW reasoners we could retrieve. We extracted such taxonomies from existing publicly accessible ontologies of enormous size. The real bonus of this method is, of course, that all three Boolean operations on sorts stay virtually $\mathcal{O}(1)$ irrespective of the size of the taxonomy nor that of the number of concepts in the query. The result of any such query is the set of sorts with codes in the set of maximal common lower-bounds of the computed code.¹⁴ All the above claims are clearly demonstrated on all the performance graphs we are reporting in the next section.

As for incrementality, removing a sort amounts simply to erasing its index position in all codes that have it set. Adding a sort (through a new “is-a” declaration $s_1 < s_2$) is done by restarting the bottom up propagation of the loop [Lines 3–14] starting from the parents of lower of the two sorts, after having reset all its ancestors to the uncoded status. In case of several new “is-a” declarations, the same procedure is applied but starting from the union of the parents of the minimal new sorts after resetting all their ancestors to the uncoded status.

3 Experimental work

3.1 Experiment setup

We extracted the bare concept taxonomies gathered from four very large ontologies. Here they are, listed in order of increasing sizes:

1. [Wikipedia](#)—derived from the Wikipedia online database (size: 111,599 sorts);¹⁵
2. [BioModels](#)—various biological models (size: 182,651 sorts);¹⁶
3. [MeSH](#)—Medical Subject Headings of the National Library of Medicine (size: 286,381 sorts);¹⁷
4. [NCBI](#)—National Center for Biotechnology Information's for all known living organisms (size: 903,617 sorts).¹⁸ This taxonomy (the largest today) is a “bare” taxonomy—*i.e.*, it contains only partially “is-a” ordered classes symbols (concepts), but no properties (*i.e.*, no roles, nor other kinds of attributes).

We focused only on bare conceptual taxonomic reasoning.¹⁹ That is, we considered no attributes (roles or features), just sorts (*i.e.*, concepts). Thus, reasoning on such sorts amounts to pure propositional logic. In other words, this boils down to computing Boolean

¹⁴See Section 4.2.

¹⁵<http://www.h-its.org/en/research/nlp/wikitaxonomy/>

¹⁶biportal.bioontology.org/ontologies/3022

¹⁷www.nlm.nih.gov/mesh/meshhome.html

¹⁸www.ncbi.nlm.nih.gov/Taxonomy/taxonomyhome.html/

¹⁹See Section 3.3 for a discussion concerning this point.

expressions consisting of sorts and the three operations: *and*, *or*, *not*. Seen another (equivalent) way, these operations applied to set-denoting expressions are interpreted respectively as set intersection, union, and complementation. The topmost sort \top denotes the set of all things, and the bottommost sort \perp denotes the empty set—*i.e.*, the set of no thing.

Section 3.2 reports the results of our comparative experiments of our Java implementation of our method with the state-of-the-art reasoners on bare conceptual taxonomies using only propositional-logic queries.

3.2 “Just the facts, Ma’am!”

This section reports the results of our experiments with the six reasoners that we retrieved (FaCT++, HerMiT, Pellet, Racerpro, TrOWL, SnoRocket) and ours (*CE_DAR*). All runs for all reasoners were carried out for exactly the same queries under the *exact* same conditions (on an Intel Core Duo CPU, 2.20 Ghz, 64-bit processor, running Windows 64, with 250 GB SATA HDD, and 16 GB main memory).

3.2.1 Classification

Figure 2 shows the comparative classification time performances for each reasoner on each of the large taxonomies we have selected. This makes six out of the seven reasoners. We did not consider SnoRocket in these classification-performance graphs because we realized that

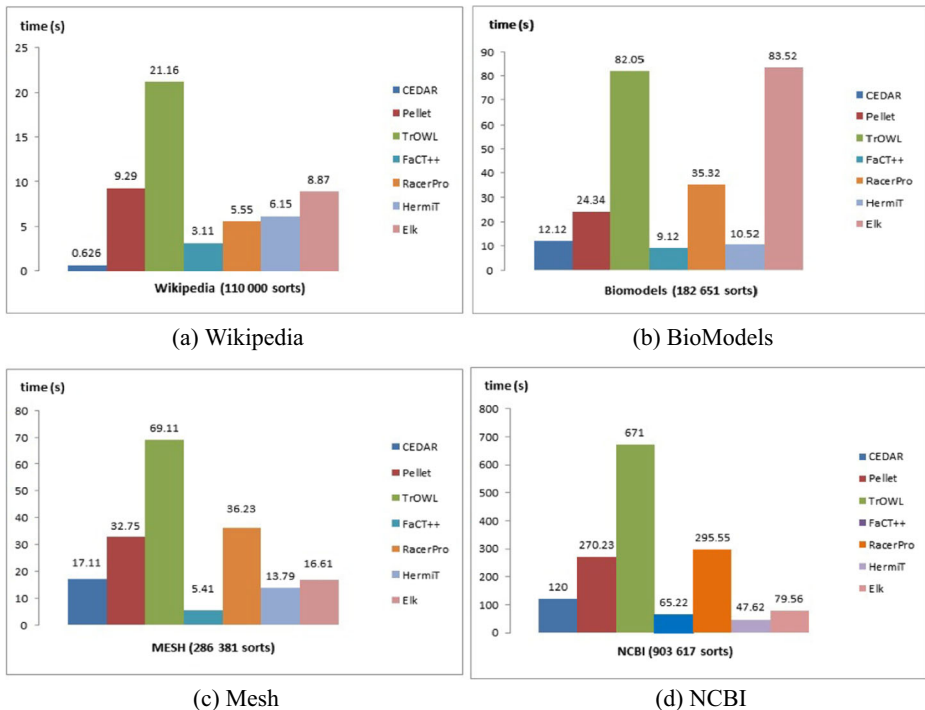


Fig. 2 Classification time per reasoner

that system does not actually perform any preliminary classification, but does so on demand at query time.

While our system's classification time (*CEDAR*—the leftmost performance on all graphs) is not always the best, it is always among the three best out of six, the worst being systematically *TrOWL*. This latter point may be due to the fact that it involves a preliminary compilation from DL to QL. Be that as it may, one can accept a longer classification time if it means faster query answering. In this regard, as illustrated next, our system definitely keeps a huge margin over all the others. It is to be noted also that *TrOWL* is faster at query answering relatively to the others (although still quite worse than our system on all tested taxonomies). This, again, may be justified by the longer classification time. On the other hand, *Hermit*, that is among the better classifiers, is always the worst for query answering. This is shown next.

3.2.2 Querying

For each of the seven reasoners, Fig. 3 shows the comparative query response time performances for two kinds of queries. If some reasoners are missing on some of these graphs, it is because they could not provide an answer before a time-out period that we set to 30 minutes.

The first series of graphs (leftmost column) are for mixed conjunctive and disjunctive queries, of the form: $s_1 \ \& \ \dots \ \& \ s_{n/2} \ \& \ (s_{n/2+1} \ | \ \dots \ | \ s_n)$, for $n = 10, 20, \dots, 100$.²⁰ We made one exception for *SnoRocket*, for which the queries were all conjunctive since the latest system available does not support disjunctive queries.

The second series of graphs (rightmost column) are for purely disjunctive queries of the form: $s_1 \ | \ \dots \ | \ s_n$, for $n = 10, 20, \dots, 100$. We did not include *SnoRocket* in this series of tests since, again, the system we retrieved does not support disjunctive queries.

In both series of graphs, it is clear that our system (*CEDAR*) systematically achieves the best performance. Moreover, it does so by several orders of magnitude (recall that the scale of time is logarithmic).

3.3 Discussion

Relative performance Although the graphs reported in Fig. 3 speak for themselves, it is interesting to get an appreciation of the relative performances for query answering of all the reasoners we have tested. In order to do so, Tables 1 and 2 sum up the facts displayed in the graphs by taking the average over all query sizes (*viz.*, from 10 to 100 concepts), giving the maximum of these averages the value 100, and showing all the other averages as percent values.

Empty cells mean that the reasoner was never able to provide an answer within our time-out limit (which, again, was set to 30 minutes). On the “BioModels” taxonomy, *Pellet* stumbles into a Java runtime error for some unknown reason.

Table 1 shows these figures for the mixed conjunctive and disjunctive queries. Table 2 shows these figures for purely disjunctive queries. Again, *SnoRocket* does not appear in the latter because it could not be tested on disjunctive queries (hence the N/A entries).

²⁰We use “&” to denote “and,” and “|” to denote “or.”

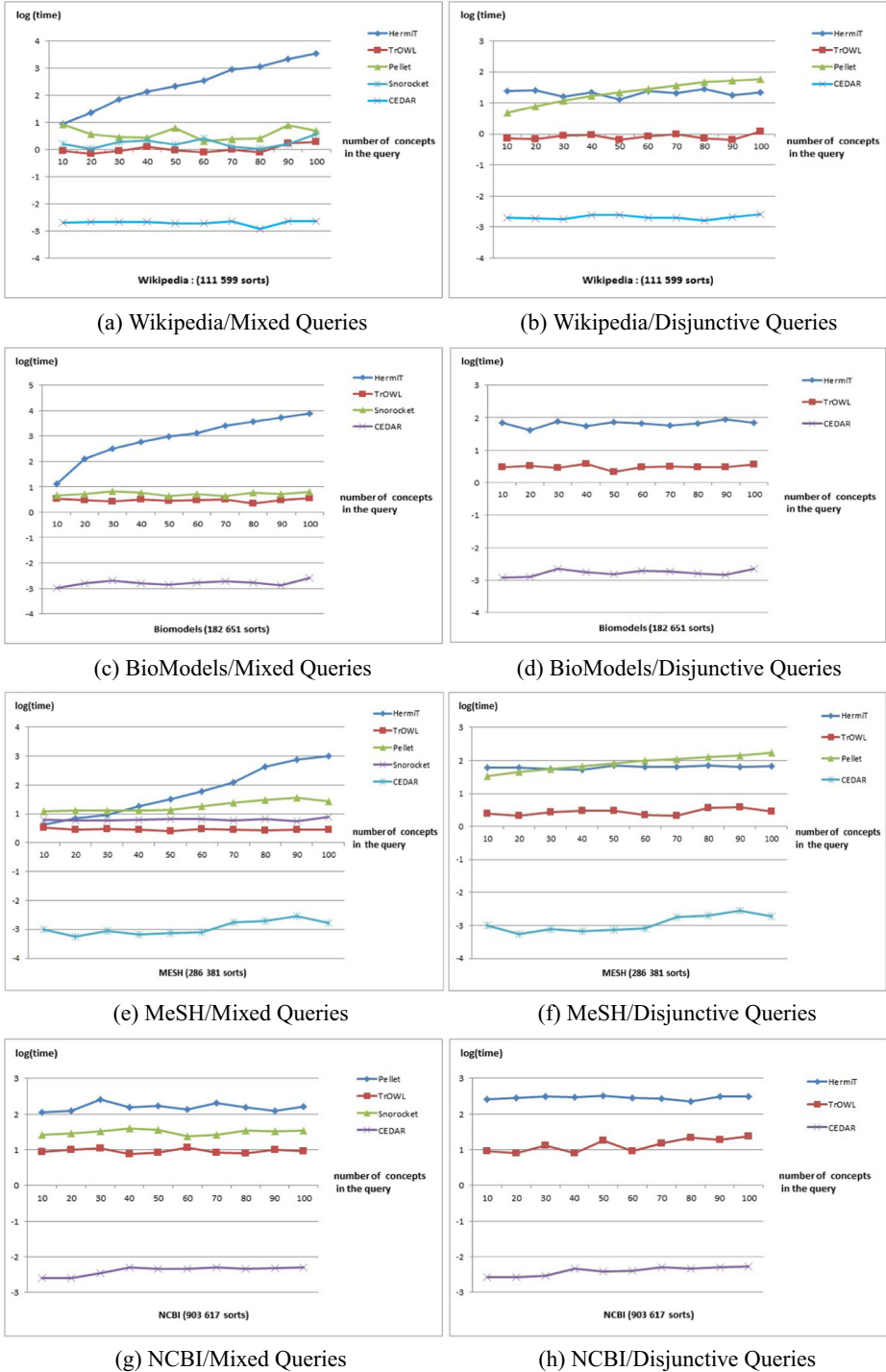


Fig. 3 Comparative performance graphs for query response times per taxonomy

Table 1 Relative normalized average performance times for mixed queries

Taxonomy	FaCT++	Hermit	TrOWL	Pellet	Racerpro	SnoRocket	<i>C&DAR</i>
Wikipedia		100	0.13	0.51		0.21	0.000233
BioModels		100	0.13	Error		0.24	0.000074
MeSH		100	1.17	8.29		2.60	0.000530
NCBI			5.78	100		19.75	0.002627

Bare taxonomies The reason why we limited our study to bare propositional reasoning is that this is (or ought to be) the most basic capability of any ontological reasoner. Any further capability in a more complete ontology based on a taxonomy (*e.g.*, reasoning with roles—existential and/or universal, cardinality constraints, *etc.*) must be conjugated with the basic propositional reasoning on taxonomic sorts. In fact, in such systems, a “complex” concept is typically a sort conjoined with some additional role-related expression. Thus, one can see bare Boolean taxonomic sort reasoning as sheer abstract interpretation of complex-concept reasoning (Cousot 1996).

Now, in terms of implementing such reasoners, it is evident that one must start with the simpler form of reasoning since it is part of all further reasoning. Still more important is that this most basic form of reasoning must be made as *efficient* as possible. In addition, factoring it out of more complex forms of reasoning (*e.g.*, for attributed taxonomies) makes the latter more efficient as it narrows it only to relevant concepts. This is simply justified as taking advantage of commutativity and associativity of conjunction (especially prior to distributing it over potential disjunctions). It is then a formal and efficient technique optimizing the process of *any* ontological reasoning, as complex as its conceptual expressions may be. Indeed, formally, taking any Boolean combination of expressions such of the form *s* & *s-properties*, ignoring the *s-properties* parts as a first pass will narrow the original expression to its essential remaining maximal consistent sorts.

Another reason that motivated us to experiment with bare taxonomies is that this is the case of several “real-life” ontologies (*e.g.*, NCBI, as already mentioned; but also, <http://www.h-its.org/en/research/nlp/wikitaxonomy/>, to name a couple).

Static processing Finally, it is worth pointing out that once a taxonomy has been classified, it may be saved on disk to be reloaded without any penalty and reused over and over.²¹ This is akin to compiling a program and not needing to recompile it for each use.

4 Pragmatic issues

In this section, we discuss two important pragmatic issues in the process of encoding and decoding a partial order. The first one concerns the detection of potential cycles when encoding taxonomies. The second one is about how to decode codes, which are elements of a Boolean lattice of binary words, into sorts which are elements of a partial order.

²¹ We implemented such a facility—see [Appendix](#).

Table 2 Relative normalized average performance times for disjunctive queries

Taxonomy	FaCT++	HermiT	TrOWL	Pellet	Racerpro	SnoRocket	\mathcal{CEDAR}
Wikipedia		74.65	2.86	100		N/A	0.00719
BioModels		100	4.68	Error		N/A	0.00258
MeSH		67.50	3.01	100		N/A	0.00141
NCBI		100	5.10			N/A	0.00141

4.1 Detecting cycles

As we found out, there is no guarantee that “real-life” taxonomies be acyclic—if only as the result of errors or inconsistent data. Detecting such potential cycles is an important issue since our encoding method’s correction relies on the taxonomy being a *dag*; *i.e.*, a directed acyclic graph. This section adds some information concerning the detection and identification of potential cycles in a set of “is-a” declarations specifying a taxonomy.

4.1.1 Problem

If the general $\mathcal{O}(n^3)$ Warshall algorithm is used to compute transitive-closure codes for all sorts in the taxonomy, then an existing maximal cycle will necessarily imply that all its elements are given equal codes. Indeed, by transitivity, the code of an element denotes the set of all its descendants. But all the elements of a cycle have the same set of descendants, and so their codes must be equal. Such a cycle is in fact an equivalence class for the least equivalence relation containing the declared “is-a” pairs.

One could eliminate all such cycles by collapsing them into a single sort (the class representative), obtaining the quotient set, which is then a dag. However, this is not desirable since such cycles are in all cases errors resulting from inconsistent declarations. In this case, they should be flagged as errors and their contents identified.

Reporting such cycles efficiently can be done by performing a topological reordering of the taxonomy according to the codes that would guarantee that sorts of equal codes are contiguous. Thus, a maximal cycle must be a maximal contiguous sequence of equal-coded sorts in this topological reordering of the taxonomy.

However, for the reasons discussed earlier in this paper, it is not feasible to use Warshall’s algorithm to compute the transitive closure on very large taxonomies. In addition, reordering such a very large taxonomy using QuickSort will be on average $\mathcal{O}(n \log n)$ with a prohibitive, although very rare, $\mathcal{O}(n^2)$ worst case (Hoare 1961).²²

Note that using our bottom-up encoding to compute the transitive closure of a taxonomy is correct only if it is a dag. If there are cycles in it, it will necessarily terminate with some of its elements left without code. This is because in Algorithm 1 a layer computed from a previous one [Line 8] removes any sort that has at least one child not encoded [Line 11]. This is correct for a dag since such sorts will always be reached later through a different longer path from \perp . But the existence of a cycle will make this assumption incorrect. For example, declaring both $s_1 < s_2$ and $s_2 < s_1$ will cause both s_1 and s_2 (and all their ancestors) to be removed from any layer to be encoded.

²²en.wikipedia.org/wiki/Quicksort

Therefore, the best we can expect with the bottom-up encoding method is that it always terminate in at most n iterations for a taxonomy of n elements. If the taxonomy is indeed a dag, all sorts will be correctly encoded. But if there are cycles, it will detect this to be the case (by checking that there remain non-encoded sorts upon termination). However, it does not have any possibility to identify how many maximal cycles there are and which sorts compose them. It is because all it knows is that bottom-up encoding left some elements non-encoded—which happens if and only if there are cycles. It does not have specific information allowing identification of which exact (maximal) cycle(s) they are.

4.1.2 Solution

In order to identify such potential maximal cycles, it is sufficient to collect all non-encoded sorts after a bottom-up encoding in a new set to be classified using Warshall's method and topologically reordered. While using Warshall's method is too costly on the full poset due to its large size, it becomes pragmatically feasible on the set of non-encoded sorts (if there are any) because such a set is always much smaller than the full declared taxonomy. In this way, all cycles can be identified as maximal contiguous elements and reported as errors.

Let us now define such a topological ordering. Recall that a taxonomy of n sorts is represented as an array of size n of `Sort` objects that are characterized by three fields; the sort's:

1. `index`—its offset in the array;
2. `code`—its bit-vector encoding;
3. `name`—its name.

Using this precedence test on sorts, we can thus obtain a unique topological (total) ordering of a taxonomy whereby a sort s is said to *precede* another sort s' iff, in this order:²³

1. $s.\text{code} < s'.\text{code}$; or,
2. $s.\text{code} = s'.\text{code}$ and $s.\text{index} < s'.\text{index}$; or,
3. s and s' are unrelated, and:
 - $|s.\text{code}| < |s'.\text{code}|$; or,
 - $|s.\text{code}| = |s'.\text{code}|$ and,

$$\text{firstDiff}(s.\text{code}, s'.\text{code}) < \text{firstDiff}(s'.\text{code}, s.\text{code}).$$

The expression $|c|$ for a code c denotes this code's cardinality (*i.e.*, its number of bits set to true). The expression $\text{firstDiff}(c, c')$ for two codes of equal cardinality c and c' denotes the lowest 1-bit position in c that is 0 in c' . So the last condition ranks sorts according to the number of descendants, and when such are equal, according to the descendant of lowest differing index.

This ordering on sorts will keep lesser sorts and sorts of lesser cardinality at lower ranks. Same-cardinality codes (*i.e.*, sort with same number of subsorts) are ranked according to lowest index of the subsort contained in one but not the other. For example, code $c = 1000111$ ($\{0, 1, 2, 6\}$) is topologically less than code $c' = 0101011$ ($\{0, 1, 3, 5\}$) because $\text{firstDiff}(c, c') = 2$ and $\text{firstDiff}(c', c) = 3$, and $2 < 3$.

It is not difficult to see that such a topological reordering will always end up with equally encoded elements being contiguous, while sorts with a greater number of lower bounds will

²³The ordering on bit-vector codes is simply defined as $c_1 \leq c_2$ iff $c_1 = c_1 \& c_2$.

be at higher ranks. In this manner, it is easy to identify all cycles in one single sweep of the taxonomy array as maximal sequences of contiguous equal codes of length at least 2.

On the Wikipedia taxonomy, the above method could identify 13 cycles: 12 of length 2 (i.e., x is-a y and y is-a x) and 1 of length 3 (i.e., x is-a y and y is-a z and z is-a x). There were also 49 warnings of cycles of length 1, which are harmless redundancies (i.e., x is-a x). All these cycles were removed in the final Wikipedia taxonomy we used in our tests and measurements.

4.2 Decoding

We now turn to decoding—i.e., relating bit-vector codes to the sorts they denote.

By construction of transitive closure, Algorithm 1, the bit vector of a sort at index i ($0 \leq i \leq n - 1$) has a 1 in position j ($0 \leq j \leq n - 1$) if and only if the sort at index j is its subsort. Therefore, a sort's bit vector has 1s in all and only the positions of its descendant sorts. For example, the taxonomy shown as Fig. 4 containing 12 sorts (other than \top and \perp) will result in the encoding shown as Table 3. Since there are 12 sorts in this taxonomy, all codes in this taxonomy have 12 bits. The top and bottom elements (\top and \perp) are implicit both in Fig. 4 and in Table 3. So the code for bottom is all 0s, and the code for \top is all 1s.

Let us first consider codes obtained without using negation. In other words, let us first restrict ourselves to decoding the result of only *positive* queries—i.e., ones involving sorts of an encoded taxonomy using a Boolean expression of its sorts' bit-vector codes using bitwise *and*, *or*—but not *not*. This always results in a bit vector. In order to determine what sorts this resulting bit vector corresponds to, there are two cases: either the resulting bit vector is that of an existing sort, or it is not.

In the first case, in order to speed up determining the sort of the bit vector, all codes are stored in a hash table mapping a code to its sort. In this way, evaluating for example “F & G” in the taxonomy of Fig. 4, which results in the bit-vector code 000000001011, the sort can be retrieved in this hash table to be associated with the code—sort D in our example.

In the second case, the code resulting from a query evaluation does not correspond to an existing sort. For example evaluating “I & J” in the taxonomy of Fig. 4 yields the

Fig. 4 Example of a small “is-a” taxonomy

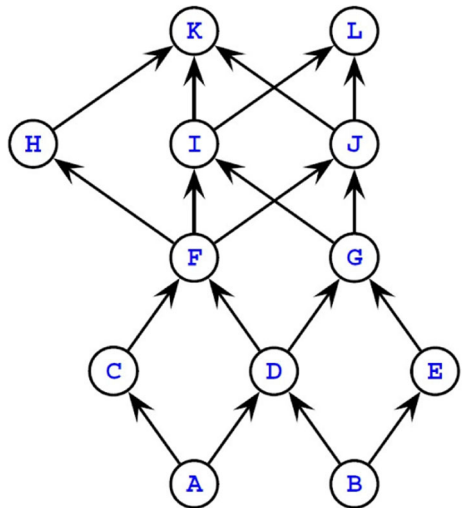


Table 3 Transitive-closure codes for sorts in Fig. 4

Index	Code	Sort
11	101110111111	L
10	011111111111	K
9	001010111111	I
8	000110111111	J
7	000010011011	G
6	000001101111	H
5	000000101111	F
4	000000011000	E
3	000000001000	B
2	000000000101	C
1	000000001011	D
0	000000000001	A

code 000010111111, which does not correspond to any specific sort in Table 3. However, semantically, this code is necessarily a minimal upper bound of the set denoted by the resulting sort if it existed. Hence, if we wish to express the resulting sort in terms of existing sorts, it is semantically the union of all the sorts whose codes are *maximal lower bounds* of the resulting code. In order to compute what sorts are in this set of maximal upper bounds, it suffices to retrieve all the sorts at index i such that there is a 1 at position i in the resulting code and keep only the maximal ones. In our example, the code 000010111111 has a 1 in positions 0, 1, 2, 3, 4, 5 and 7. This means that its subsorts are A, D, C, B, E, F, and G. However, among these, only F and G are maximal. Therefore, the result of the query “I & J” is the *disjunctive* sort {F ; G}.

While the above decoding scheme is correct for positive queries, it is not so however if the query made use of negation. To see this, let us consider the taxonomy shown as Fig. 5 and its encoding: shown as Table 4.

Let us now consider the (negative) query: “!canid.” The code resulting from evaluating this query is the complement of the sort “canid”—namely, 011001111. The decoding

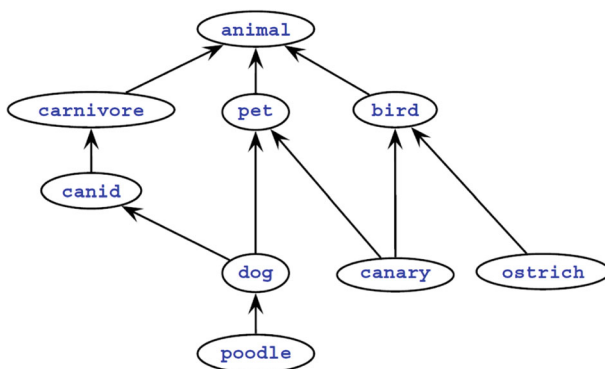


Fig. 5 Example of a small animal “is-a” taxonomy

Table 4 Transitive-closure codes for sorts in Fig. 5

Index	Code	Sort
8	100000000	poodle
7	010000000	canary
6	001000000	ostrich
5	100100000	dog
4	100110000	canid
3	110101000	pet
2	100110100	carnivore
1	111111111	animal
0	011000001	bird

method that we used above will yield the maximal elements in index set $\{0, 1, 2, 3, 6, 7\}$ —*viz.*, $\{\text{bird, animal, pet, ostrich, canary, poodle}\}$; namely, $\{\text{animal}\}$, which is obviously wrong. This decoding is incorrect because a negated code can no longer be interpreted as having a ‘1’ in a position corresponding to a subsort. Indeed, for such a code, a ‘1’ in position i means that sort of index i is either a supersort or unrelated. (This is because it comes from negating ‘0’ at position i in the complement where that meant “sort of index i is *not* a subsort of this sort.”)

This means that if a code to be decoded results from operations involving at least one negation, it must be ensured that all the 1s in it are “genuine” indicators of lower bounds (which may not be the case if *such a ‘1’ came from complementing a ‘0’ by negation*).

In order to do that, let us consider the case of negating a sort expression s . If we compute the code of $!s$ by switching all 1s to 0s and vice versa in the code of s , this will not work (as explained above). However, if we change a 1 in position i to a 0 in the code of $!s$ whenever the sort of index i is a supersort of a sort corresponding to a 0 in the code of $!s$, then we will be left with a code having a 1 only in positions of actual subsorts of $!s$. With such a code, the decoding method for positive-query coded can then compute the correct set of maximal lower bounds of $!s$.

Given a code to be decoded, we call its set of “*undesirables*” the set of (indices of) ancestors of sorts corresponding to 0s in this code. Clearly, by the very semantics of encoding, these indices cannot be 1s in the code to be decoded. For if they were, so should be their descendants—but such is not the case since they are precisely defined as ancestors of sorts corresponding to a 0 in the code. So, if we switch off any 1 in the code to be decoded that corresponds to an “undesirable,” we are left only with a code that can now be decoded with the method described above for codes involving no negation—since we are now guaranteed that all the 1s denote actual subsorts. Finally, note that for a code resulting from operations involving no negation, the set of “undesirables” is necessarily empty. Indeed, in this case, there can be no 0 in a position of a descendant of an actual ancestor.

In our example above, for the code 011001111 resulting from the evaluation of “! canid ,” the “undesirables” are the set of ancestors of sorts of indices $\{4, 5, 8\}$ (that is, the set of sorts $\{\text{canid, dog, poodle}\}$). This set is the set of indices $\{1, 2, 3, 4\}$ (or sorts $\{\text{animal, carnivore, pet, canid}\}$). Switching off undesirable bits gives the code: 011000001. This set of indices ($\{0, 6, 7\}$) denotes the set of sorts $\{\text{bird, ostrich, canary}\}$. Keeping only maximal elements, this yield the (correct) answer: “bird.”

For efficiency reasons, once a bit-vector code has been decoded, it is stored in a cache (a hash table) associating the code to the set of sorts whose codes are its maximal lower bounds. In this way, should the same code appear again as a result, it is first looked up in this cache to avoid the need to compute again its set of maximal lower bound sorts.

Finally, note that decoding a bit vector is only necessary for extracting the end result of a query in terms of defined sorts. All intermediate computation need not refer at all to the sorts and deal only with bit vectors, whether or not they correspond to defined sorts. There is no loss of information doing so as the encoding plunges the taxonomy in the minimal Boolean lattice containing it (Aït-Kaci et al. 1989).

5 Further work

In Amir and Aït-Kaci (2014a, b), we have extended this work to unification-based Knowledge Representation known as Order-Sorted Feature (OSF) constraint logic (Aït-Kaci 2007). While OSF logic uses functional features, we can use them to represent roles using aggregates. The advantage is that role-based reasoning is thus made simpler since it relies on Logic-Programming unification technology made possible by functional attributes (Aït-Kaci and Di Cosmo 1993; Aït-Kaci 1991). This is akin to compiling DL-based relational roles into aggregate-valued functional features. OSF sorts have also a “memoizing” effect whereby no property needs to be proven again once it has been established for any supersort (Aït-Kaci 2007).

As for the nature of codes (binary vectors), there are optimizations that may still be performed to our basic method. The technique known as *code modulation* (explained in detail in Aït-Kaci et al. (1989)) can take advantage of a taxonomy’s specific shape to minimize drastically code space. This makes so-modulated bit-vector encoding very scalable as explained in (Aït-Kaci et al. 1989). Indeed, modulated encoding reduced code size to a logarithmic function of non-modulated encoding. Code modulation being independent of the encoding technique, it can be applied to any method. Each module can in fact use different encoding methods each adapted to its specific topology. It can thus be used recursively (*i.e.*, one can modulate a module) with the encoding technique most appropriate to the topology of the module being encoded.

Another optimization to minimize classification time could be to perform lazy encoding.²⁴ In other words, one could only encode the sorts relevant to a query and cache intermediate results. The price to pay would be at query time, although only the first time a subset of the concepts it involves are used.

6 Conclusion

In this paper, we have presented an implementation of a Boolean Logic taxonomic reasoner based on bit-vector encoding. We have implemented such a reasoner in Java, and have compared its performance for pure taxonomic reasoning to that of six among the most renown Semantic Web reasoners. Focusing only on pure Boolean taxonomic reasoning—which is at the core of *every* SW reasoner—the results of our measurements show that our system achieves the best performance, by a large margin. This establishes beyond doubt that the

²⁴In the same manner as we have noticed that SnoRocket does.

best existing Semantic Web reasoners fail to live up to performance that can be easily achieved using bit-vector encoding. This is true even, and especially, when applied to very large taxonomies.

In order to carry out these experiments, we also developed a tool with an easy-to-use GUI that lets a user run these tests for any listed reasoner and taxonomy. This tool is available for download from the *CEDAR* Project's website for anyone to verify our results.²⁵ Also, a video clip of these demos showing these experiments *in vivo*, and a web service for running these demos on line are available online.²⁶ In this way, our results may hopefully not have to be taken on faith, but could be verified *de visu* by anyone who might wish to check them on their own.

Acknowledgments The authors wish to thank Prof. Mohand-Saïd Hacid and the anonymous referees for constructive feedback.

Conflict of interests The authors declare that they have no conflict of interest.

Funding This work was carried out as part of the *CEDAR* Project (Constraint Event-Driven Automated Reasoning) under the *Agence Nationale de la Recherche* (ANR) Chair of Excellence grant N° **ANR-12-CHEX-0003-01**.

Appendix

In Section A, we give an overview of an implementation specification for representing very large bit vectors, reducing memory-space consumption while retaining efficient operations. In our experiments, this alternative code representation was used only for saving an encoded taxonomy on disk and reloading it as a pre-encoded order. But for taxonomies of even larger size than those used in our experiments, it could be used for lattice operations as well.

A Compact codes

While the foregoing sections present a method for encoding elements of a partially ordered set based on transitive closure, the data structure it relies on is that of a binary word—*i.e.*, a bit vector. With such a structure, all boolean operations—*and*, *or*, *not*—are thus very efficient. This representation also eases computation of the transitive closure since setting a bit on or off is trivially accommodated.

However, while this representation is convenient and time-efficient for relatively small posets of the order of a few hundred elements, it quickly becomes space-inefficient for large posets of hundreds of thousands, or millions of elements.

In what follows, we define an alternative representation of indexed bit sets that offers the advantage of being more compact than bit vectors while retaining time-efficient boolean and bit-setting operations. It is also the format we use to save/load encoded taxonomies on/from disk. In Section A.1, the basic data structure is defined. In Section A.2, bit

²⁵cedar.liris.cnrs.fr/data/CEDAR-V1.0.zip

²⁶cedar.liris.cnrs.fr/demos.html

setting and unsetting operations are defined. In Section A.3, the three boolean operations—conjunction, disjunction, and negation—are defined. In Section A.4, some implementation considerations are discussed.

A.1 Bit code representation

The idea is intuitively simple. It consists of representing a bit vector as a finite array of k ($k \in \mathbb{N}$) pairs of integer indices $\langle l_i, u_i \rangle$, for $i = 0, \dots, k - 1$, such that, for all indices $i = 0, \dots, k - 2$:

$$0 \leq l_i < u_i < l_{i+1} < u_{k-1}. \tag{1}$$

We shall refer to such a sequence k pairs, $k \in \mathbb{N}$, $\{ \langle l_i, u_i \rangle \mid i = 0, \dots, k - 1 \}$ as a *compact code*. For $k = 0$, this is written as the empty sequence $\{ \}$.

Given a compact code representation of a bit vector V , each pair $\langle l, u \rangle$ represents a *maximal* contiguous sequence of 1's (hereafter referred to as a “packet”) in V . Thus, the i -th packet of a bit vector is represented as the pair of indices $\langle l_i, u_i \rangle$ such that l_i is the index of the lowest bit in the packet, and u_i is the index of the first 0-bit following the packet.

For example, the bit vector 0011111001111000000110000 corresponds to the compact code (*i.e.*, sequence of packet pairs):²⁷ $\{ \langle 4, 6 \rangle, \langle 12, 16 \rangle, \langle 18, 23 \rangle \}$.

The empty bit vector (containing all 0's) is represented as the empty sequence $\{ \}$. The *length* of a bit vector represented by a compact code sequence of k pairs (or packets) is u_k . The *size* of a compact code C of k pairs (or packets) is k (*i.e.*, its number of packets).

A.2 Bit operations

Let $C = \{ \langle l_i, u_i \rangle \mid i = 0, \dots, k \}$ be a compact code of k packets ($k > 0$). Given a number $n \in \mathbb{N}$, and a compact code C of k packets as defined above, we say that:²⁸

- n is *within a packet* of C iff $\exists i \in [0, k - 1]$ such that $l_i \leq n < u_i$ —in which case we shall write $C.\mathbf{packet}(n) = i$;
- n is *between packets* of C iff either one of the three statements holds:
 1. $n < l_0$; or,
 2. $u_{k-1} \leq n$; or,
 3. $\exists i \in [0, k - 2]$ such that $u_i \leq n < l_{i+1}$.

If a number n is between packets of a compact code C of size k , we define two functions $C.\mathbf{prev}(n)$ and $C.\mathbf{next}(n)$ for each of the three possible respective cases above as follows (where the symbol ‘?’ means “undefined”):

1. $C.\mathbf{prev}(n) \stackrel{def}{=} ?$ and $C.\mathbf{next}(n) \stackrel{def}{=} l_0$;
2. $C.\mathbf{prev}(n) \stackrel{def}{=} u_k$ and $C.\mathbf{next}(n) \stackrel{def}{=} ?$;
3. $C.\mathbf{prev}(n) \stackrel{def}{=} u_i$ and $C.\mathbf{next}(n) \stackrel{def}{=} l_{i+1}$.

For such a number n , we say that:

- n is *left-adjacent* in C if $n = C.\mathbf{next}(n) - 1$;

²⁷Recall that a bit vector is written with its lowest bit to the right.

²⁸In what follows, we shall use the “dot” notation of object-oriented methods to denote all functions or operations on codes.

- n is right-adjacent in C if $n = C.\mathbf{prev}(n)$;
- n is adjacent in C if it is both left-adjacent and right-adjacent in C .

Note that if n is between packets and adjacent, this necessarily means that the two packets on each side are only separated by a single 0-bit (the bit in position n in the denoted bit vector).

N.B.: In all the compact code expressions to follow, we use the implicit convention that a packet with undefinable bounds is simply omitted. Thus, we will always use the notation $\{ \langle l_0, u_0 \rangle, \dots, \langle l_{k-1}, u_{k-1} \rangle \}$ to denote a compact code, where $k \geq 0$ up to the above conventions regardless of the actual number of packets. For example, for $k = 0$ this will correspond to the empty code $\{ \}$, and for $k = 1$, this will correspond to the single-packet code $\{ \langle l_0, u_0 \rangle \}$.

We define the following bit-setting operations on C . These methods operate “in place” by modifying a code C that invokes them.

- $C.\mathbf{set}(n)$, for $n \in \mathbb{N}$, which sets the n -th bit of the bit vector denoted by C to 1.
- $C.\mathbf{set}(n, m)$, for $n, m \in \mathbb{N}, n < m$, which sets to 1 all the bits from position n (*inclusive*) to position m (*exclusive*) of the bit vector denoted by C .
- $C.\mathbf{unset}(n)$, for $n \in \mathbb{N}$, which sets the n -th bit of the bit vector denoted by C to 0.
- $C.\mathbf{unset}(n, m)$, for $n, m \in \mathbb{N}, n < m$, which sets to 0 all the bits from position n (*inclusive*) to position m (*exclusive*) of the bit vector denoted by C .

For $m \leq n$, both $C.\mathbf{set}(n, m)$ and $C.\mathbf{unset}(n, m)$ are no_ops—*i.e.*, they leave C unchanged. Since $\mathbf{set}(n)$ is equivalent to $\mathbf{set}(n, n + 1)$, we will just give the methods for $\mathbf{set}(n, m)$ and similarly for $\mathbf{unset}(n, m)$.

A.2.1 Bit setting

There are four cases to consider for which performing $C.\mathbf{set}(n, m)$ modifies C as follows.

1. If n is within a packet in C (say, $C.\mathbf{packet}(n) = i$) and m is within a packet in C (say, $C.\mathbf{packet}(m) = j$), then, if $i = j$, $C.\mathbf{set}(n, m)$ leaves C unchanged. Else (if $i < j$),²⁹ then C becomes:

$$\{ \dots, \langle l_i, u_j \rangle, \dots \}.$$

2. If n is within a packet in C (say, $C.\mathbf{packet}(n) = i$) and m is between packets in C , then C becomes:

$$\left\{ \begin{array}{l} \{ \dots, \langle l_i, u_j \rangle, \dots \} \\ \quad \text{if } m \text{ is left-adjacent in } C \text{ and } C.\mathbf{next}(m) = l_j; \\ \\ \{ \dots, \langle l_i, m \rangle, \dots \} \\ \quad \text{otherwise.} \end{array} \right. \quad (2)$$

3. If n is between packets in C and m is within a packet in C (say, $C.\mathbf{packet}(m) = j$), then C becomes:

$$\left\{ \begin{array}{l} \{ \dots, \langle l_i, u_j \rangle, \dots \} \\ \quad \text{if } n \text{ is right-adjacent in } C \text{ and } C.\mathbf{prev}(n) = u_i; \\ \\ \{ \dots, \langle n, u_j \rangle, \dots \} \\ \quad \text{otherwise.} \end{array} \right. \quad (3)$$

²⁹Note that $i \neq j$ implies necessarily that $i < j$ (by Condition (1) and since $n < m$).

4. If both n and m are between packets in C , then C becomes:

$$\left\{ \begin{array}{l} \{ \dots, \langle l_i, u_j \rangle, \dots \} \\ \quad \text{if } n \text{ is right-adjacent in } C \text{ and } C.\mathbf{prev}(n) = u_i, \text{ and} \\ \quad \text{if } m \text{ is left-adjacent in } C \text{ and } C.\mathbf{next}(m) = l_j; \\ \\ \{ \dots, \langle l_i, m \rangle, \dots \} \\ \quad \text{if } n \text{ is right-adjacent in } C \text{ and } C.\mathbf{prev}(n) = u_i, \text{ and} \\ \quad \text{if } m \text{ is not left-adjacent in } C; \\ \\ \{ \dots, \langle n, u_j \rangle, \dots \} \\ \quad \text{if } n \text{ is not right-adjacent in } C, \text{ and} \\ \quad \text{if } m \text{ is left-adjacent in } C \text{ and } C.\mathbf{next}(m) = l_j; \\ \\ \{ \dots, \langle n, m \rangle, \dots \} \\ \quad \text{otherwise.} \end{array} \right. \quad (4)$$

A.2.2 Bit unsetting

Here again, there are four cases to consider for which performing $C.\mathbf{unset}(n, m)$ modifies C as follows.

1. If both n and m are between packets in C : if $C.\mathbf{prev}(n) = C.\mathbf{prev}(m)$ (or, equivalently, if $C.\mathbf{next}(n) = C.\mathbf{next}(m)$), then $C.\mathbf{unset}(n, m)$ leaves C unchanged; else, C becomes:

$$\{ \dots, \langle l_i, C.\mathbf{prev}(n) \rangle, \langle C.\mathbf{next}(m), u_j \rangle, \dots \}.$$

2. If n is within a packet in C (say, $C.\mathbf{packet}(n) = i$) and m is between packets in C , then C becomes:

$$\left\{ \begin{array}{l} \{ \dots, \langle l_{i-1}, u_{i-1} \rangle, \langle C.\mathbf{next}(m), u_j \rangle, \dots \} \\ \quad \text{if } n = l_i, \text{ where } C.\mathbf{next}(m) = l_j; \\ \\ \{ \dots, \langle l_i, n \rangle, \langle C.\mathbf{next}(m), u_j \rangle, \dots \} \\ \quad \text{else (i.e., if } n > l_i), \text{ where } C.\mathbf{next}(m) = l_j. \end{array} \right. \quad (5)$$

3. If n is between packets in C and m is within a packet in C (say, $C.\mathbf{packet}(m) = j$), then C becomes:

$$\left\{ \begin{array}{l} \{ \dots, \langle l_i, C.\mathbf{prev}(n) \rangle, \langle l_{j+1}, u_{j+1} \rangle, \dots \} \\ \quad \text{if } m = u_j - 1, \text{ where } C.\mathbf{prev}(m) = u_i; \\ \\ \{ \dots, \langle l_i, C.\mathbf{prev}(n) \rangle, \langle m, u_j \rangle, \dots \} \\ \quad \text{else (i.e., if } m < u_j - 1), \text{ where } C.\mathbf{prev}(m) = u_i. \end{array} \right. \quad (6)$$

4. If both n is within a packet (say, $C.\mathbf{packet}(n) = i$), and m is within a packet (say, $C.\mathbf{packet}(m) = i$) in C , then C becomes:

$$\left\{ \begin{array}{l} \{ \dots, \langle l_{i-1}, u_{i-1} \rangle, \langle l_{j+1}, u_{j-1} \rangle, \dots \} \\ \text{if } n = l_i, \text{ and} \\ \text{if } m = u_{j-1}; \\ \\ \{ \dots, \langle l_{i-1}, u_{i-1} \rangle, \langle m, u_j \rangle, \dots \} \\ \text{if } n = l_i, \text{ and} \\ \text{if } m < u_{j-1}; \\ \\ \{ \dots, \langle l_i, n \rangle, \langle l_{j+1}, u_{j-1} \rangle, \dots \} \\ \text{if } n > l_i, \text{ and} \\ \text{if } m = u_{j-1}; \\ \\ \{ \dots, \langle l_i, n \rangle, \langle m, u_j \rangle, \dots \} \\ \text{if } n > l_i, \text{ and} \\ \text{if } m < u_{j-1}. \end{array} \right. \quad (7)$$

A.3 Boolean operations

Let:

$$C = \{ \langle l_i, u_i \rangle \mid i = 0, \dots, k - 1 \}$$

and:

$$C' = \{ \langle l'_i, u'_i \rangle \mid i = 0, \dots, k' - 1 \}$$

be two compact code pair sequences, with $k \geq 0$ and $k' \geq 0$.

A.3.1 Conjunction

Invoking $C.\mathbf{and}(C')$ will modify C according to C' by unsetting all the bits in C that are between packets in C' , leaving C' unchanged.

If $C = \{\}$, then C is left unchanged; else, if $C' = \{\}$, then C becomes $\{\}$.

Else (*i.e.*, if $k > 0$ and $k' > 0$), C is modified by invoking:

- $C.\mathbf{unset}(0, l'_0)$; and,
- $C.\mathbf{unset}(u'_i, l'_{i+1})$, for $i = 0$ up to $i = k' - 1$; and,
- $C.\mathbf{unset}(u'_{k'-1}, u_{k-1})$.

Note that in practice, when proceeding in the above order, as soon the first argument of any $\mathbf{unset}(\dots)$ is greater than or equal to u_{k-1} , there is no need to perform the unsetting nor proceed any further.

A.3.2 Disjunction

Invoking $C.\mathbf{or}(C')$ will modify C according to C' by setting all the bits in C that are within packets in C' , leaving C' unchanged.

If $C' = \{\}$, then C is left unchanged; else, if $C = \{\}$, then C becomes (a copy of) C' .

Else (*i.e.*, if $k > 0$ and $k' > 0$), C is modified by invoking:

- $C.\mathbf{set}(l'_i, u'_i)$, for $i = 0$ up to $i = k' - 1$.

A.3.3 Negation

Since a bit vector is open-ended, we may define its negation only up to a length at least greater than its highest 1-bit position. This operation is denoted as $C.\mathbf{not}(n)$. Thus, $\{\}. \mathbf{not}(n)$, is undefined for any $n \geq 0$.

Otherwise, for a non-empty code $C = \{ \langle l_0, u_0 \rangle, \dots, \langle l_{k-1}, u_{k-1} \rangle \}$ and $n \geq u_{k-1}$, $C.\mathbf{not}(n)$ modifies C to become:

$$\{ \langle 0, l_0 \rangle, \dots, \langle u_i, l_{i+1} \rangle, \dots, \langle u_{k-1}, n \rangle \}.$$

Again, following our convention, if $l_0 = 0$, $\langle 0, l_0 \rangle$ being undefinable, the first element of $C.\mathbf{not}(n)$ is $\langle u_0, l_1 \rangle$. Similarly, if $n = u_{k-1}$, then $\langle u_{k-1}, n \rangle$ is undefinable and the last element of $C.\mathbf{not}(n)$ is $\langle u_{k-2}, l_{k-1} \rangle$.

A.4 Implementation considerations

We need to come up with a data structure for representing a compact code that would enable retaining maximal efficiency in the bit setting and unsetting operations, and hence in the boolean operations that rely on them.

Most frequently used operations on such a data structure C for an integer n are:

- $C.\mathbf{packet}(n)$ —for n inside a packet in C , returning that packet number;
- $C.\mathbf{prev}(n)$ —for n between packets in C , returning the upper index of the packet preceding n ;
- $C.\mathbf{next}(n)$ —for n between packets in C , returning the lower index of the packet following n ;
- adding/removing a packet.

Because the elements of a code sequence are ranges rather than integers, one cannot expect hashed $\mathcal{O}(1)$ time access to find out whether a given integer lies within or between packets. So structures such as defined by the Java classes `java.util.HashSet` or `java.util.LinkedHashSet` cannot be used.

In order to make these operations at most $\mathcal{O}(\log(k))$ time for a compact code of size k , one way is to represent a compact code as a balanced binary tree of pairs of bit position spans $\langle l_i, u_i \rangle$, taking advantage of the ordering imposed by condition (1).

Thus, the `java.util.TreeSet` class looks like a convenient choice, since it offers the required data structure properties in addition to defining methods such as `first`, `last`, `higher`, `lower`, `add`, `remove`, *etc.*, as well as an order-respecting iterator.

On the other hand, the `java.util.TreeSet` is missing a `replace` method—which is critically needed for setting and unsetting bits. It is also missing an `insert` method that splices a new sequence of packet pairs into an existing compact code, which may also be often used. One must resort to several `add/remove` method invocations to replace or insert elements, which incur new searches (and possible intermediate rebalancing of the tree) each time. This is a waste since replacing and inserting can be done in $\mathcal{O}(1)$ time when having already found the required elements, and only one (final) $\mathcal{O}(\log(k))$ tree rebalancing.

Hence, rather than relying on the ready-to-use `java.util.TreeSet` class, it may be more beneficial to implement a new specific class for a compact code as a doubly linked list and a balanced binary tree adapted for the specific nature of its pair elements. This

would make transparent the double links of each pair element and ease replacement and insertion.³⁰

A.5 Discussion

A similar data structure was proposed by researchers in data and knowledge bases in (Agrawal et al. 1989). However, the authors did not use that representation for lattice operations as we do here. Instead, they focused on using it for obtaining more compact range-sequence codes for the transitive closure of the “is-a” relation of a taxonomy. That representation is equivalent to the one we specify here and to the one in (Aït-Kaci et al. 1989). Contrary to (Aït-Kaci et al. 1989), they define an element’s code as the union of index ranges from the post-order arrangement of the a spanning tree of the “is-parent-of” relation of a taxonomy. Each concept in the taxonomy (*i.e.*, each element in the poset) of post-order index j is then encoded as the interval $\langle i, j \rangle$ where i is the smallest post-order index of all its descendants. Although they did not do it, it is easy to show that their representation is equivalent to bit vectors. But they did not specify lattice operations on their data structures as we do for ours in this document. What they focused on was minimizing the total number of packets in range codes. In order to do so, they suggest generating codes based on the “optimal” spanning tree for generating the most compact set of codes. The data structure and algorithm for what they call “compressed transitive closure” do not maintain dynamic interval consistency caused by potential adjacency as we do here.³¹ Although they do cite (Aït-Kaci et al. 1989), Agrawal et al. do so only in the conclusion as they had just noticed its publication. They suggest that their approach and that exposed in Aït-Kaci et al. (1989) might be combined for processing large taxonomies. As far as we know, no follow-up on this suggestion was carried out.

It is clear how the work of (Agrawal et al. 1989), although orthogonal to ours, could be adapted to our needs as well in order to improve its space consumption. However, it is to be noted that their code-compaction method requires a topologically ordered poset. For very large taxonomies (over 1 million elements) the price of sorting the taxonomy might be worth spending only for once-for-all preprocessing prior to query time (Aït-Kaci and Amir 2013; Amir and Aït-Kaci 2013). Also, the question of incrementality is not addressed.

Finally, although this work has been motivated for obtaining a compact representation of binary codes encoding a partial order, it comes as evident that the data structure, and operations on it, specified in this appendix can represent any set of integers (or integer-indexed set) seen as a sequence of intervals. Set intersection is realized as the conjunction described in Section A.3.1; set union as the disjunction described in Section A.3.2; and, set complementation as the negation described in Section A.3.3. Therefore, it can readily be used for this purpose as well.

³⁰Actually, the `java.util.TreeSet` does maintain a doubly-linked list for its elements in order to ensure its two ordered iterators (ascending and descending). But this structure is not made public and one cannot splice in new elements from a given found element. But it is a simple matter to modify the source code of `java.util.TreeSet.java` and adapt it to what is needed.

³¹In fact, they see that only as a possible *a posteriori* optimization, but one that would cause their optimal spanning-tree finding algorithm to be incorrect if applied incrementally while it is executed.

References

- Agrawal, R., Borgida, A., & Jagadish, H.V. (1989). Efficient management of transitive relationships in large data and knowledge bases. In J. Clifford, B.G. Lindsay, & D. Maier (Eds.), *Proceedings of the ACM SIGMOD International Conference on Management of Data*, (Vol. 18(2) pp. 253–262). Portland: ACM, SIGMOD Record. [Available online http://dbs.informatik.uni-halle.de/Lehre/DBS.IIa_SS02/p253-agrawal.pdf].
- Ait-Kaci, H. (1991). *Warren's Abstract Machine: A Tutorial Reconstruction*. Cambridge: The MIT Press. [Available online <http://wambook.sourceforge.net/>].
- Ait-Kaci, H. (2007). Data models as constraint systems—A key to the Semantic Web. *Constraint Processing Letters*, 1(1), 33–88. [Available online www.cs.brown.edu/people/pvh/CPL/Papers/v1/hak.pdf].
- Ait-Kaci, H., & Amir, S. (2013). Classifying and querying very large taxonomies—a comparative study to the best of our knowledge. *CEDAR Technical Report Number 2*, CEDAR Project, LIRIS, Département d'Informatique, Université Claude Bernard Lyon 1, Villeurbanne, France. [Available online <http://cedar.liris.cnrs.fr/papers/ctr2.pdf>].
- Ait-Kaci, H., Boyer, R., Lincoln, P., & Nasr, R. (1989). Efficient implementation of lattice operations. *ACM Transactions on Programming Languages and Systems*, 11(1), 115–146. [Available online <http://hassan-ait-kaci.net/pdf/encoding-toplas-89.pdf>].
- Ait-Kaci, H., & Di Cosmo, R. (1993). *Compiling order-sorted feature term unification* Vol. 7. France: PRL Technical Note, Digital Paris Research Lab, Rueil-Malmaison. [Available online <http://hassan-ait-kaci.net/pdf/PRL-TN-7.pdf>].
- Amir, S., & Ait-Kaci, H. (2013). CEDAR: a fast taxonomic reasoner based on lattice operations. In *Proceedings of the ISWC 2013 Posters & Demonstrations Track, Sydney, Australia* (pp. 9–12).
- Amir, S., & Ait-Kaci, H. (2014a). CEDAR: efficient reasoning for the semantic web. In *Tenth International Conference on Signal-Image Technology and Internet-Based Systems, SITIS 2014* (pp. 157–163), Marrakech, Morocco.
- Amir, S., & Ait-Kaci, H. (2014b). Design and implementation of an efficient semantic web reasoner. CEDAR Technical Report Number 12, CEDAR Project, LIRIS, Département d'Informatique, Université Claude Bernard Lyon Vol. 1, Villeurbanne, France. [Available online <http://cedar.liris.cnrs.fr/papers/ctr12.pdf>].
- Baader, F., Brandt, S., & Lutz, C. (2005). Pushing the \mathcal{EL} envelope. In L. P. Kaelbling, & A. Saffiotti (Eds.), *Proceedings of the 19th International Joint Conference on Artificial Intelligence* (pp. 364–369). Edinburgh: IJCAI'05, Morgan Kaufmann Publishers. [Available online www.ijcai.org/papers/0372.pdf].
- Baader, F., Lutz, C., & Suntisrivaraporn, B. (2006). CEL—a polynomial-time reasoner for life science ontologies. In U. Furbach, & N. Shankar (Eds.), *Proceedings of the 3rd international joint conference on automated reasoning*, (Vol. 4130 pp. 287–291). Seattle: IJCAR'06, Springer-Verlag LNAI. [Available online <http://lat.inf.tu-dresden.de/research/papers/2006/BaaLutSun-IJCAR-06.pdf>].
- Ballard, G., Demmel, J., Holtz, J., & Schwartz, O. (2014). Communication costs of Strassen's matrix multiplication. *Communications of the ACM*, 57(2), 107–114. [Available online <https://aspire.eecs.berkeley.edu/wp/wp-content/uploads/2014/02/Communication-Costs-of-Strassen%E2%80%99s-Matrix-Multiplication.pdf>].
- Coppersmith, D., & Winograd, S. (1990). Matrix multiplication via arithmetic progressions. *Journal of Symbolic Computation*, 9(3), 251–280. [Available online <http://www.sciencedirect.com/science/article/pii/S0747717108800132>].
- Cousot, P. (1996). Abstract interpretation. *ACM Computing Surveys—Symposium on Models of Programming Languages and Computation*, 28(2), 324–328. Tutorial summary: [Available online <http://www.di.ens.fr/~cousot/AI/IntroAbsInt.html>].
- Fikes, R., Hayes, P., & Horrocks, I. (2004). OWL-QL—a language for deductive query answering on the Semantic Web. *Journal of Web Semantics*, 2(1), 19–29. [Available online <http://www.sciencedirect.com/science/article/pii/S1570826804000137>].
- Fischer, M.J., & Meyer, A.R. (1971). Boolean matrix multiplication and transitive closure. In *Proceedings of the 12th annual symposium on switching and automata theory, SWAT'71*, (pp. 129–131), Washington: IEEE Computer Society. [Available online rjlipton.files.wordpress.com/2009/10/matrix1971.pdf].
- Haarslev, V., Hidde, K., Möller, R., & Wessel, M. (2011). The RacerPro knowledge representation and reasoning system. *Semantic Web Journal*, 1, 1–11. [Available online http://www.semantic-web-journal.net/sites/default/files/swj109_3.pdf].
- Haarslev, V., & Möller, R. (2001). RACER system description. In R. Gore, A. Leitsch, & T. Nipkow (Eds.), *Proceedings of the 1st international joint conference on automated reasoning* (pp. 701–706). Siena: IJCAR'01, Springer-Verlag. [Available online <https://www.ifis.uni-luebeck.de/~moeller/papers/2001/HaMo01a.pdf>].

- Hoare, C.A.R. (1961). Algorithm 63: Partition, algorithm 64: Quicksort. *Communications of the ACM*, 4(7), 321–321. [Available online <http://comjnl.oxfordjournals.org/content/5/1/10.full.pdf>].
- Horrocks, I., & Sattler, U. (2007). A tableau decision procedure for \mathcal{SHIQ} . *Journal of Automated Reasoning*, 39(3), 249–276. [Available online <http://www.cs.ox.ac.uk/ian.horrocks/Publications/download/2007/HoS07a.pdf>].
- Kazakov, Y. (2009). Consequence-driven reasoning for horn \mathcal{SHIQ} ontologies. In C. Boutilier (Ed.), *Proceedings of the 21st international conference on artificial intelligence* (pp. 2040–2045). Pasadena: IJCAI'09, Association for the Advancement of Artificial Intelligence. [Available online <http://ijcai.org/papers09/Papers/IJCAI09-336.pdf>].
- Kazakov, Y., Krötzsch, M., & Simančík, F. (2011). Unchain my \mathcal{EL} reasoner. In R. Rosati, S. Rudolph, & M. Zakharyashev (Eds.), *Proceedings of the 24th international workshop on description logics*. Barcelona: DL'11, CEUR Workshop Proceedings. [Available online http://ceur-ws.org/Vol-745/paper_54.pdf].
- Lawley, M.J., & Bousquet, C. (2010). Fast classification in Protégé: Snorocket as an OWL 2 EL reasoner. In T. Meyer, M. A. Orgun, & K. Taylor (Eds.), *Proceedings of the 2nd Australasian ontology workshop: Advances in ontologies* (pp. 45–50). Adelaide: AOW'10, ACS. [Available online <http://kr.meraka.org.za/~aow2010/Lawley-etal.pdf>].
- Manna, Z., & Waldinger, R. (1991). Fundamentals of deductive program synthesis. In A. Apostolico, & Z. Galil (Eds.), *Combinatorial algorithms on words, NATO ISI Series*: Springer. [Available online <http://lara.epfl.ch/~kuncak/t/MannaWaldingerTSE.pdf>].
- Motik, B., Shearer, R., & Horrocks, I. (2009). Hypertableau reasoning for description logics. *Journal of Artificial Intelligence Research*, 36(1), 165–228. [Available online www.jair.org/media/2811/live-2811-4689-jair.pdf].
- Shearer, R., Motik, B., & Horrocks, I. (2008). HerMiT: A highly-efficient OWL reasoner. In U. Sattler, & C. Dolbear (Eds.), *Proceedings of the 5th international workshop on OWL experiences and directions*. Karlsruhe: OWLED'08, CEUR Workshop Proceedings. [Available online <http://www.cs.ox.ac.uk/ian.horrocks/Publications/download/2008/ShMH08b.pdf>].
- Sirin, E., Parsia, B., Grau, B.C., Kalyanpur, A., & Katz, Y. (2007). Pellet: A practical OWL-DL reasoner. *Journal of Web Semantics*, 5(2), 51–53. This is a summary; full paper: [Available online <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.96.5433&rep=rep1&type=pdf>].
- Stothers, A. (2010). *On the complexity of matrix multiplication*. Edinburgh: PhD thesis, University of Edinburgh. [Available online <http://www.maths.ed.ac.uk/assets/files/pgrexternalfiles/theses/probability/stothers.pdf>].
- Strassen, V. (1969). Gaussian elimination is not optimal. *Numerische Mathematik*, 13, 354–356.
- Thomas, E., Pan, J.Z., & Ren, Y. (2010). TrOWL: Tractable OWL 2 reasoning infrastructure. In L. Aroyo, G. Antoniou, E. Hyvnen, A. ten Teije, H. Stuckenschmidt, L. Cabral, & T. Tudorache (Eds.), *Proceedings of the 7th extended semantic web conference* (pp. 431–435). Heraklion: ESWC'10, Springer. [Available online <http://homepages.abdn.ac.uk/jeff.z.pan/pages/pub/TPR2010.pdf>].
- Tsarkov, D., & Horrocks, I. (2006). FaCT++ description logic reasoner: System description. In U. Furbach, & N. Shankar (Eds.), *Proceedings of the 3rd international joint conference on automated reasoning* (pp. 292–297). Seattle: IJCAR'06, Springer. [Available online <http://www.cs.ox.ac.uk/Ian.Horrocks/Publications/download/2006/TsHo06a.pdf>].
- Warren Jr, H.S. (1975). A modification of Warshall's algorithm for the transitive closure of binary relations. *Communications of the ACM*, 18(4), 218–220. [Available online <http://dl.acm.org/citation.cfm?id=360746>].
- Warshall, S. (1962). A theorem on Boolean matrices. *Journal of the ACM*, 9(1), 11–12.
- Williams, V.V. (2011). *Breaking the Coppersmith-Winograd barrier*: University of California at Berkeley and Stanford University. [Available online <http://www.cs.rit.edu/~rlc/Courses/Algorithms/Papers/matrixMult.pdf>].