# Privacy-Preserving Ontology Publishing: The Case of Quantified ABoxes w.r.t. a Static Cycle-Restricted $\mathcal{EL}$ TBox[⋆]

Franz Baader[1], Patrick Koopmann[1], Francesco Kriegel[1],
Adrian Nuradiansyah[1], and Rafael Peñaloza[2]

[1] Theoretical Computer Science, Technische Universität Dresden, Dresden, Germany
`firstname.lastname@tu-dresden.de`
[2] University of Milano-Bicocca, Milano, Italy
`rafael.penaloza@unimib.it`

**Abstract.** We review our recent work on how to compute optimal repairs, optimal compliant anonymizations, and optimal safe anonymizations of ABoxes containing possibly anonymized individuals. The results can be used both to remove erroneous consequences from a knowledge base and to hide secret information before publication of the knowledge base, while keeping as much as possible of the original information.

**Keywords:** Repair · Compliance · Safety · Privacy-preserving ontology publishing · Optimality · Complexity · Practical algorithm

## 1 Introduction

In contrast to most of the work in description logic (DL), which is about how to derive consequences of a DL knowledge base (KB) efficiently, this paper is about how to get rid of consequences. The reason for this wish can, on the one hand, be that a certain consequence is incorrect, and thus one wants to repair the KB to get rid of this error. On the other hand, one may want to remove a consequence since it is considered to be private information that is not supposed to be made public. In both cases, the new KB should not introduce new consequences (i.e., it should be entailed by the original one), and it should be optimal in the sense that a minimal amount of consequences is removed (i.e., it should be as close as possible to the original one w.r.t. the entailment relation).

Though both scenarios share the wish to remove consequences, there are some differences. On the technical side, in the context of repairs one usually considers a single consequence or a finite set of consequences of the form $C(a)$, i.e., one wants to get rid of instance relationships for specific individuals [3].[1] The resulting KB is then called a repair of the original one. In the context of privacy,

---

[1] We restrict the attention here to derived instance relationships, though repairs for subsumption relationships have also been considered in the literature [7].

one usually considers a policy $\mathcal{P}$, consisting of one or a finite number of concepts, and wants to get rid of all consequences of the form $C(a)$ for $C \in \mathcal{P}$ and $a$ a named individual [8,9,17,18]. The resulting KB is then said to be a $\mathcal{P}$-compliant anonymization of the original one. Instead of changing the data before publishing it, one could also provide only restricted access through queries, whose answers are monitored by a "censor", which may decide not to give an answer or even lie if needed to satisfy the policy. There has been quite some work in this direction in the database and the DL community [12, 14–16, 23], but this approach is not the topic of the present paper.

On the intentional side, achieving compliance is not always sufficient to guarantee privacy [9, 17, 18]. In fact, an attacker may already have some knowledge, which does not imply the secret, but which together with a published compliant anonymization may be used to derive the secret information. Thus, in the context of privacy, one is interested in computing anonymization that are safe in the sense that, even if extended with an arbitrary compliant KB, they do not imply $C(a)$ for $C \in \mathcal{P}$ and $a$ a named individual.

In the general setting of a DL KB consisting of a TBox and an ABox, optimal repairs (optimal compliant/safe anonymizations) need not exist [3, 7]. There are two ways to overcome this problem. On the one hand, one can weaken the notion of optimality and restrict the attention to repairs (anonymizations) that can be obtained from the original KB by applying certain repair (anonymization) steps. This approach is, e.g., followed in [17, 18] in the setting of privacy and in [7,19,22,29] for the repair scenario. Classical repair approaches that completely remove axioms rather than just weakening them also fall under this category [11, 21, 24, 25, 27, 28].

On the other hand, one can stick with the quest for optimality, and restrict the considered KBs such that optimality can be achieved. Our first work in this direction [6, 10] considered compliance and safety in the very restricted setting of an $\mathcal{EL}$ instance store [20], i.e., where there is no TBox and the ABox does not contain role assertions. In the first paper, the attacker's knowledge is considered to be a set of $\mathcal{EL}$ concept assertions (an $\mathcal{EL}$ instance store) whereas in the second also other DLs are used to represent the attacker's knowledge. In [8] we extended the results of [6] to ABoxes with role assertions (and still no TBox), but restricted the attention to compliance for $\mathcal{EL}$ policies. In [9] we investigated safety in the same setting, but had to restrict the policies to ones consisting of a single $\mathcal{EL}$ concept (singleton policies). Finally, in [3] we extended the results of [8] in two directions, but formulate the new results in the (more general) context of repairs rather than compliance. On the one hand, we add a TBox, which must however be cycle-restricted. On the other hand, we develop a more practical algorithm for computing optimal repairs.

This paper summarizes the results obtained in our previous publications [3,8, 9], but presents them uniformly in the setting of privacy. In addition, it extends the results of [9] by developing a more practical algorithm for computing optimal safe anonymizations in the setting without a TBox. Finally, we show that using TBoxes one can reduce safety for general policies to safety for singleton policies.

Since characterizing safety for general policies is an open problem, this shows that extending our results for safety to the case with TBox is a non-trivial problem. Proofs of our new results can be found in the technical report [5].

## 2   Preliminaries

We use concepts of the lightweight DL $\mathcal{EL}$ both to define TBoxes and to formulate which consequences are unwanted. The data are represented in the form of quantified ABoxes, which are atomic ABoxes (i.e., ones not containing assertions for complex concepts) in which some individual names are assumed to be anonymized. While such anonymous individuals do not belong to the standard DL repertoire, they are actually available in OWL. Also, such ABoxes have already been used in previous work on privacy-preserving ontology publishing [17, 18]. Finally, note that concept and role assertions involving anonymous individuals can be used to express concept assertions for complex concepts.

**$\mathcal{EL}$ concepts and TBoxes.** We assume basic knowledge about DLs [2]. Specifically, we consider the DL $\mathcal{EL}$, defined over a fixed *signature* $\Sigma$, which is the disjoint union of the countably infinite sets $\Sigma_{\mathsf{O}}$, $\Sigma_{\mathsf{C}}$, and $\Sigma_{\mathsf{R}}$ of *object names*, *concept names*, and *role names*. $\mathcal{EL}$ concepts are built using the concept constructors $\top$, $\sqcap$ and $\exists$. We treat conjunctions as sets, that is, they do not contain duplicates and the order is irrelevant. $\mathcal{EL}$ TBoxes, in the following just called TBoxes, are defined as usual as sets of concept inclusions (CIs) $C \sqsubseteq D$. We use the notation $C \sqsubseteq^{\mathcal{T}} D$ (alternatively $\mathcal{T} \models C \sqsubseteq D$) to denote that $C \sqsubseteq D$ holds in all models of $\mathcal{T}$. A TBox is called *cycle-restricted* if there is no non-empty sequence of role names $r_1, \ldots, r_k$ and no $\mathcal{EL}$ concept $C$ such that $C \sqsubseteq^{\mathcal{T}} \exists r_1. \cdots \exists r_k.C$. Cycle-restrictedness of a given TBox can be decided in polynomial time [1].

An *atom* is of the form $A$ or $\exists r.C$, where $A \in \Sigma_{\mathsf{C}}$, $r \in \Sigma_{\mathsf{R}}$, and $C$ is a concept. Every $\mathcal{EL}$ concept $C$ is a conjunction of atoms (with $\top$ as empty conjunction), called the *top-level conjunction* of $C$. We denote the set of atoms occurring in it as $\mathsf{Conj}(C)$. Given a TBox $\mathcal{T}$ and a set $\mathcal{C}$ of concepts, we use $\mathsf{Sub}(\mathcal{T}, \mathcal{C})$ to denote the set of concepts occurring in $\mathcal{T}$ and $\mathcal{C}$ (as elements or subconcepts), $\mathsf{Atoms}(\mathcal{T}, \mathcal{C})$ to denote the set of atoms occurring in $\mathcal{T}$ and $\mathcal{C}$, and similarly for $\mathsf{Sub}(\mathcal{C})$ and $\mathsf{Atoms}(\mathcal{C})$ for the concepts and atoms occurring in $\mathcal{C}$. Given two sets of $\mathcal{EL}$ concepts $\mathcal{K}$ and $\mathcal{L}$, we say that $\mathcal{K}$ *is covered by* $\mathcal{L}$ (written $\mathcal{K} \leq \mathcal{L}$) if, for every $C \in \mathcal{K}$, there is $D \in \mathcal{L}$ s.t. $C \sqsubseteq^{\emptyset} D$.

**Quantified ABoxes.** We use a generalisation of ABoxes called *quantified ABoxes (qABoxes)* to adequately represent anonymous individuals as in OWL and *nulls* common in database systems, which play a central role in anonymization [18]. To illustrate, consider the ABox $\{r(a, b), A(a), B(b)\}$, and assume we want to hide the fact that $b$ is an instance of $B$. Quantified ABoxes allow us to achieve this in a better way than by just deleting the fact $B(b)$, namely by additionally adding an anonymous copy of $b$, resulting in the quantified ABox

$\exists \{x\}. \{r(a,b), A(a), r(a,x), B(x)\}$, for which $a$ is still an instance of $\exists r.B$. In fact, this qABox is equivalent to the ABox $\{r(a,b), (A \sqcap \exists r.B)(a)\}$, which uses a concept assertion involving the complex concept $A \sqcap \exists r.B$.

Essentially, qABoxes are syntactic variants of conjunctive queries. Formally, a qABox is of the form $\exists X.\mathcal{A}$, where $X$ is a finite subset of $\Sigma_\mathsf{O}$, the elements of which are called *variables*, and $\mathcal{A}$ is the *matrix*, a finite set of concept assertions $A(u)$ where $u \in \Sigma_\mathsf{O}$ and $A \in \Sigma_\mathsf{C}$, and of role assertions $r(u,v)$ where $u, v \in \Sigma_\mathsf{O}$ and $r \in \Sigma_\mathsf{R}$. Without loss of generality, we assume different qABoxes to use disjoint sets of variables. A non-variable object name in $\exists X.\mathcal{A}$ is called an *individual name*, and the set of all these names is denoted as $\Sigma_\mathsf{I}(\exists X.\mathcal{A})$. We further set $\Sigma_\mathsf{O}(\exists X.\mathcal{A}) \coloneqq \Sigma_\mathsf{I}(\exists X.\mathcal{A}) \cup X$. Traditional DL ABoxes are qABoxes where $X = \emptyset$; we then write $\mathcal{A}$ instead of $\exists \emptyset.\mathcal{A}$. The matrix $\mathcal{A}$ of a qABox $\exists X.\mathcal{A}$ is such a traditional ABox. An interpretation $\mathcal{I}$ is a *model* of a qABox $\exists X.\mathcal{A}$ if there is an interpretation $\mathcal{J}$ such that $\Delta^\mathcal{I} = \Delta^\mathcal{J}$, the interpretation functions $\cdot^\mathcal{I}$ and $\cdot^\mathcal{J}$ coincide on $\Sigma \setminus X$, and $u^\mathcal{J} \in A^\mathcal{J}$ for each $A(u) \in \mathcal{A}$ as well as $(u^\mathcal{J}, v^\mathcal{J}) \in r^\mathcal{J}$ for each $r(u,v) \in \mathcal{A}$.

Let $\mathcal{T}$ be a TBox and $\exists X.\mathcal{A}$, $\exists Y.\mathcal{B}$ two qABoxes. We write $\exists X.\mathcal{A} \models^\mathcal{T} \exists Y.\mathcal{B}$ to express that every model of $\mathcal{T}$ and $\exists X.\mathcal{A}$ is also a model of $\exists Y.\mathcal{B}$, in which case we say $\exists Y.\mathcal{B}$ *is entailed by* $\exists X.\mathcal{A}$ *w.r.t.* $\mathcal{T}$. Entailment of traditional ABoxes from a qABox can be decided in polynomial time, while entailment between qABoxes is NP-complete.

## 3   Computing Optimal Compliant Anonymizations

A *policy* is a finite set of $\mathcal{EL}$ concepts. Intuitively, a policy says that one should not be able to derive that any of the individuals of a qABox belongs to a concept in the policy. To make a given qABox compliant to a policy, we compute an anonymization of it, which is a compliant qABox entailed by it. Intuitively, such an anonymization is optimal if it does not remove more information than necessary.

**Definition 1.** *Let $\mathcal{T}$ be a TBox, $\mathcal{P}$ be a policy, and $\exists X.\mathcal{A}$, $\exists Y.\mathcal{B}$ be qABoxes.*

1. *$\exists X.\mathcal{A}$ is* compliant with $\mathcal{P}$ w.r.t. $\mathcal{T}$ *if, for each $a \in \Sigma_\mathsf{I}(\exists X.\mathcal{A})$ and $C \in \mathcal{P}$, $\exists X.\mathcal{A} \not\models^\mathcal{T} C(a)$,*
2. *$\exists Y.\mathcal{B}$ is a $\mathcal{P}$-compliant anonymization of $\exists X.\mathcal{A}$ w.r.t. $\mathcal{T}$ if $\exists X.\mathcal{A} \models^\mathcal{T} \exists Y.\mathcal{B}$ and $\exists Y.\mathcal{B}$ is compliant with $\mathcal{P}$ w.r.t. $\mathcal{T}$;*
3. *$\exists Y.\mathcal{B}$ is an* optimal $\mathcal{P}$-compliant anonymization of $\exists X.\mathcal{A}$ w.r.t. $\mathcal{T}$ *if additionally $\exists Z.\mathcal{C} \models^\mathcal{T} \exists Y.\mathcal{B}$ implies $\exists Y.\mathcal{B} \models^\mathcal{T} \exists Z.\mathcal{C}$ for every $\mathcal{P}$-compliant anonymization $\exists Z.\mathcal{C}$ of $\exists X.\mathcal{A}$ w.r.t. $\mathcal{T}$.*

Since, in $\mathcal{EL}$, entailment of concept assertions (viewed as singleton ABoxes) is in P, we can decide compliance in polynomial time. More interesting is the question of how to compute a (preferably optimal) anonymization for a given qABox. This problem is investigated in [8] for the case without TBox, and in [3] for the case with TBoxes. These works also consider a weaker version of entailment, called

*IQ-entailment*, for the case where we are only interested in instance queries, and [3] considers a generalisation of anonymizations called *ABox repairs*, where instead of a policy, a set of assertions is given that should not be entailed. For brevity, we focus here on the version of anonymizations defined above.

To guarantee existence of optimal anonymizations, we restrict ourselves to cycle-restricted TBoxes. An example where the TBox is not cycle-restricted and where no optimal repairs exist is as follows. Consider the traditional ABox $\{A(a)\}$, the TBox $\{A \sqsubseteq \exists r.\,A, \ \exists r.\,A \sqsubseteq A\}$, and the policy $\{A\}$. Intuitively, an optimal anonymization would have to entail any qABox of the form $\exists\{x_0, \ldots, x_n\}.\{r(a, x_0), r(x_i, x_{i+1}) \mid 0 \leq i \leq n-1\}$ for $n \geq 0$, which is not possible for a qABox entailed by $\{A(a)\}$ w.r.t. $\mathcal{T}$. A formal proof that there is no optimal compliant anonymization in this case can be found in [7]. As shown in [3], this problem can be avoided by considering IQ-entailment, which we do not discuss here.

Next, we present a class of anonymizations called canonical anonymizations, which cover all optimal anonymizations. They are given by a rather elegant direct definition, but may be hard to compute in practice. We then present an optimized approach that computes smaller representations of them.

### 3.1   Canonical Compliant Anonymizations

If the TBox $\mathcal{T}$ is cycle-restricted, it is possible to compute (in exponential time) its *saturation*, i.e., a qABox $\mathsf{sat}^{\mathcal{T}}(\exists X.\mathcal{A})$ such that for every qABox $\exists Y.\mathcal{B}$, $\exists X.\mathcal{A} \models^{\mathcal{T}} \exists Y.\mathcal{B}$ iff $\mathsf{sat}^{\mathcal{T}}(\exists X.\mathcal{A}) \models^{\emptyset} \exists Y.\mathcal{B}$. The saturation integrates into the qABox all relevant information that can be inferred using the TBox, so that entailments can be decided without use of the TBox (see [3] for how to compute $\mathsf{sat}^{\mathcal{T}}(\exists X.\mathcal{A})$).

In our approach, we first compute the saturation, and then perform the actual anonymization based on *repair types* and *compliance seed functions*. For convenience, we fix in the following the TBox $\mathcal{T}$, policy $\mathcal{P}$ and qABox $\exists X.\mathcal{A}$ given as input, and abbreviate $\Sigma_{\mathsf{I}}(\exists X.\mathcal{A})$ as $\Sigma_{\mathsf{I}}$. A repair type specifies for a given object which entailments are to be removed by the anonymization.

**Definition 2.** *Let* $\exists Y.\mathcal{B} \coloneqq \mathsf{sat}^{\mathcal{T}}(\exists X.\mathcal{A})$ *and* $u \in \Sigma_{\mathsf{O}}(\exists Y.\mathcal{B})$. *A* repair type for $u$ *is a subset* $\mathcal{K}$ *of* $\mathsf{Atoms}(\mathcal{P}, \mathcal{T})$ *that satisfies the following:*

1. *$\mathcal{B} \models^{\emptyset} C(u)$ for each atom $C \in \mathcal{K}$,*
2. *if $C, D$ are distinct atoms in $\mathcal{K}$, then $C \not\sqsubseteq^{\emptyset} D$,*
3. *$\mathcal{K}$ is* premise-saturated *w.r.t. $\mathcal{T}$, i.e., for all $C \in \mathsf{Sub}(\mathcal{P}, \mathcal{T})$ s.t. $\mathcal{B} \models^{\emptyset} C(u)$ and $C \sqsubseteq^{\mathcal{T}} D$ for some $D \in \mathcal{K}$, there is $E \in \mathcal{K}$ such that $C \sqsubseteq^{\emptyset} E$.*

Condition 1 makes sure the concepts in the repair type are indeed entailed for the given individual. Condition 2 avoids redundancies, and Condition 3 ensures that removing the corresponding assertions is effective also in presence of the TBox. The compliance seed function now assigns to every named individual a repair type based on the given policy.

**Definition 3.** *A* compliance seed function *is a function $s$ that maps each individual name $b \in \Sigma_I$ to a repair type $s(b)$ for $b$ such that, if $C \in \mathcal{P}$ and $\mathsf{sat}^{\mathcal{T}}(\exists X.\mathcal{A}) \models^{\emptyset} C(b)$, then there is $D \in s(b)$ such that $C \sqsubseteq^{\emptyset} D$.*

Each compliance seed function induces a compliant anonymization defined next. Intuitively, for concept names $A \in s(a)$, we simply remove the concept assertion $A(a)$ from $\mathcal{A}$. For atoms of the form $\exists r.C \in s(a)$, we need to modify the role successors of $a$ such that $\exists r.C(a)$ is no longer entailed. To avoid losing more information than necessary, we do not just remove assertions from the objects in $\mathcal{A}$, but also create copies of objects by introducing new variables, which are based on the set of repair types for each object name.

**Definition 4.** *Given a compliance seed function $s$, we define the* canonical compliant anonymization *$\mathsf{ca}^{\mathcal{T}}(\exists X.\mathcal{A}, s)$ induced by $s$ as the qABox $\exists Y.\mathcal{B}$ where:*

1. *The set $Y$ consists of the variables $y_{u,\mathcal{K}}$ s.t. $u$ is an object name in $\mathsf{sat}^{\mathcal{T}}(\exists X.\mathcal{A})$ and $\mathcal{K}$ is a repair type for $u$, except for the case where $u$ is an individual name and $\mathcal{K} = s(u)$. In the latter case, we keep the individual name $u$, but use $y_{u,s(u)}$ as a synonym for $u$ in the definition of $\mathcal{B}$ below.*
2. *The matrix $\mathcal{B}$ consists of the following assertions:*
   (a) *$A(y_{u,\mathcal{K}})$ if $A(u)$ occurs in $\mathsf{sat}^{\mathcal{T}}(\exists X.\mathcal{A})$ and $A \notin \mathcal{K}$, and*
   (b) *$r(y_{u,\mathcal{K}}, y_{v,\mathcal{L}})$ if $r(u,v)$ occurs in $\mathsf{sat}^{\mathcal{T}}(\exists X.\mathcal{A})$ and for each $\exists r.C \in \mathcal{K}$ s.t. the matrix of $\mathsf{sat}^{\mathcal{T}}(\exists X.\mathcal{A})$ entails $C(v)$, there exists $D \in \mathcal{L}$ s.t. $C \sqsubseteq^{\emptyset} D$.*

Every qABox $\mathsf{ca}^{\mathcal{T}}(\exists X.\mathcal{A}, s)$ induced by a seed function $s$ is a compliant anonymization of $\exists X.\mathcal{A}$, but it need not be optimal. However, every optimal compliant anonymization is induced (up to equivalence) by some seed function. Thus, we can compute all optimal compliant anonymizations (modulo equivalence) by computing all canonical compliant anonymizations and then removing the non-optimal ones. The latter requires testing entailment between quantified ABoxes.

**Theorem 5 ([3]).** *There is a deterministic, exponential time algorithm with access to an* NP *oracle that computes the set of all optimal compliant anonymizations of $\exists X.\mathcal{A}$ for $\mathcal{P}$ w.r.t. $\mathcal{T}$.*

### 3.2   Optimality Using Minimal Seed Functions

The NP oracle in Theorem 5 is needed for the NP-complete entailment test, which is applied to exponentially large qABoxes. If it is sufficient to compute some, rather than all, optimal compliant anonymizations, we can dispense with the NP oracle and instead utilize a (polynomial time decidable) partial order on seed functions [8]. For two compliance seed functions $s$ and $t$, we say that *$s$ is covered by $t$* (written $s \leq t$) if $s(a)$ is covered by $t(a)$ for every $a \in \Sigma_I$, i.e., for every $C$ in $s(a)$ there is $D$ in $t(a)$ s.t. $C \sqsubseteq^{\emptyset} D$.

**Proposition 6 ([8]).** *If $\mathsf{ca}^{\mathcal{T}}(\exists X.\mathcal{A}, s) \models^{\mathcal{T}} \mathsf{ca}^{\mathcal{T}}(\exists X.\mathcal{A}, t)$ for two compliance seed functions $s$ and $t$, then $s \leq t$.*

This was shown in [8] for the case without a TBox, but the proof can easily be extended to the case considered here.

The proposition implies that each minimal seed function induces an optimal anonymization. Since there is always at least one minimal seed function and since $\leq$ can be decided in polynomial time, we can draw the following conclusion.

**Theorem 7 ([8]).** *A non-empty set of optimal compliant anonymizations of $\exists X.\mathcal{A}$ for $\mathcal{P}$ w.r.t. $\mathcal{T}$ can be computed in exponential time.*

### 3.3   Smaller Optimal Compliant Anonymizations

Since the number of variables introduced in a canonical compliant anonymization is always exponential in the size of the TBox and the policy,[2] computing even one of them in practice quickly becomes infeasible. The exponential blow-up is in general not avoidable, already for the very limited case without TBox and where the qABox corresponds to an $\mathcal{EL}$ instance store [6]. However, in many practical cases, we can compute a compliant anonymization that is significantly smaller than the canonical compliant anonymization, but logically equivalent to it [3, 4]. The idea is to avoid introducing unnecessary variables by starting with the individual names and unmodified single copies of all object names, and then incrementally determining which variables of the canonical anonymization need to be included, where in each step we only look at the immediate role-successors of each object name and the requirements expressed in the associated repair type.

To be more precise, let $s$ be a repair seed function and $\exists Y.\mathcal{B} := \mathsf{ca}^{\mathcal{T}}(\exists X.\mathcal{A}, s)$. According to Definition 4, we have $r(y_{t,\mathcal{K}}, y_{u,\mathcal{L}}) \in \mathcal{B}$ iff $\mathsf{sat}^{\mathcal{T}}(\exists X.\mathcal{A})$ contains the role assertion $r(t, u)$ and the repair type $\mathcal{L}$ covers

$$\mathsf{Succ}(\mathcal{K}, r, u) := \{\, C \mid \exists r.C \in \mathcal{K} \text{ and the matrix of } \mathsf{sat}^{\mathcal{T}}(\exists X.\mathcal{A}) \text{ entails } C(u) \,\}.$$

Our procedure produces a sequence $Y_0, Y_1, \ldots, Y_m$ of subsets $Y_i$ of $Y$ such that $\exists Y.\mathcal{B}$ is equivalent to $\exists Y_m.\mathcal{B}_m$, where $\mathcal{B}_m$ is the subset of $\mathcal{B}$ that uses only objects from $\Sigma_\mathsf{I} \cup Y_m$. We start with the set

$$Y_0 := \{\, y_{t,\emptyset} \mid t \text{ is an object name occurring in } \mathsf{sat}^{\mathcal{T}}(\exists X.\mathcal{A}) \,\}.$$

The subsequent sets are obtained by exhaustively applying the following rule:

**Compliant Anonymization Rule.**
   **If** (i) $y_{t,\mathcal{K}}$, $y_{u,\mathcal{L}} \in \Sigma_\mathsf{I} \cup Y_i$,  (ii) $r(t, u)$ occurs in $\mathsf{sat}^{\mathcal{T}}(\exists X.\mathcal{A})$,  (iii) $\mathcal{L}$ does not cover $\mathsf{Succ}(\mathcal{K}, r, u)$,  (iv) there is a covering-minimal repair type $\mathcal{M}$ for $u$ that covers $\mathcal{L} \cup \mathsf{Succ}(\mathcal{K}, r, u)$,  and (v) $y_{u,\mathcal{M}} \notin \Sigma_\mathsf{I} \cup Y_i$,
   **then** set $Y_{i+1} := Y_i \cup \{y_{u,\mathcal{M}}\}$.

Since each rule application adds a variable, the exhaustive application of the Compliant Anonymization Rule must terminate after finitely many steps with a set $Y_m \subseteq Y$ of variables. We call $\exists Y_m.\mathcal{B}_m$ the *optimized compliant anonymization* of $\exists X.\mathcal{A}$ w.r.t. $\mathcal{T}$ induced by the seed function $s$.

---

[2] However, canonical anonymizations can be computed in polynomial time w.r.t. data complexity, i.e., if only the size of the qABox counts (TBox and policy fixed).
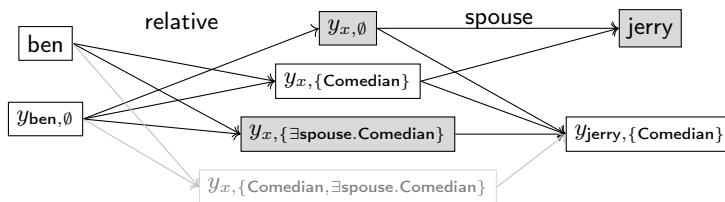
Fig. 1: Canonical anonymization (all nodes) and the subset that is the optimized anonymization (non-shadowed nodes). Gray nodes denote instances of Comedian.

**Theorem 8 ([3]).** *For each compliance seed function $s$, the optimized compliant anonymization induced by $s$ is equivalent to $\mathsf{ca}^{\mathcal{T}}(\exists X.\mathcal{A}, s)$.*

To compute $\mathcal{B}_m$ we do not need to compute the larger matrix $\mathcal{B}$ first. Instead, we directly apply the definition of the matrix (Definition 4) to the object names in $\Sigma_{\mathsf{I}} \cup Y_m$. Experiments with an implementation[3] of this procedure (for the more general case of ABox repairs) indicate that applying this optimized procedure reduces the size of the computed compliant anonymizations considerably [3, 4].

*Example 9.* To illustrate both kinds of anonymizations, consider an empty TBox, policy $\mathcal{P} := \{P\}$ for $P := \exists \mathsf{relative}.(\mathsf{Comedian} \sqcap \exists \mathsf{spouse}.\mathsf{Comedian})$, and qABox $\exists X.\mathcal{A} := \exists\{x\}.\{\mathsf{relative}(\mathsf{ben}, x), \mathsf{Comedian}(x), \mathsf{spouse}(x, \mathsf{jerry}), \mathsf{Comedian}(\mathsf{jerry})\}$. As seed function, we select $s$ s.t. $s(\mathsf{ben}) = \{P\}$ and $s(\mathsf{jerry}) = \emptyset$. Fig. 1 depicts both the canonical and the optimized compliant anonymization.

## 4    Safety of Quantified ABoxes, Mainly Without TBox

To guarantee privacy, policy compliance is not always sufficient since an attacker may have additional knowledge that, by itself, does not reveal the secret, but which, together with the to be published compliant information, would violate the privacy policy. This is captured by the notion of *safety*: a qABox $\exists X.\mathcal{A}$ is *safe* for a given policy $\mathcal{P}$ if for every $\mathcal{P}$-compliant $\exists Y.\mathcal{B}$, the union $\exists(X \cup Y).(\mathcal{A} \cup \mathcal{B})$ is compliant with $\mathcal{P}$ as well.

This definition is based on the assumption that the additional knowledge possessed by the attacker is also in the form of a qABox (the qABox $\exists Y.\mathcal{B}$ in the formal definition). Since we do not know which additional knowledge the attacker has, we need to consider all possible compliant qABoxes $\exists Y.\mathcal{B}$. Non-compliant qABoxes $\exists Y.\mathcal{B}$ need not be considered here: in fact, it is useless trying to hide the secret information from such an attacker that already knows it.

For instance, the canonical compliant anonymization shown in Figure 1 is not safe since one could add the compliant qABox

$$\exists\{y\}.\{\mathsf{relative}(\mathsf{ben}, y), \mathsf{Comedian}(y), \mathsf{spouse}(y, \mathsf{jerry})\}.$$

---

[3] https://github.com/de-tu-dresden-inf-lat/abox-repairs-wrt-static-tbox

In the resulting qABox, Ben is an instance of the policy concept $P$.

In [9], we give a characterization for safety of qABoxes for *singleton policies*,[4] which are of the form $\{P\}$ for an $\mathcal{EL}$ concept $P$. Specifically, safety for $\{P\}$ is violated if (1) $A(a) \in \mathcal{A}$ for some individual name $a$ and $A \in \mathsf{Atoms}(P)$, or (2) $r(a, u) \in \mathcal{A}$ and $\exists r.D \in \mathsf{Atoms}(P)$ such that a part of the concept $D$ can be found in $\exists X.\mathcal{A}$ at the specific object $u$ — in both cases we can construct attacking compliant qABoxes as certificates for non-safety. The second condition is captured by the notion of *partial homomorphisms* (cf. Definition 3.6 in [9]). Intuitively, a partial homomorphism from a concept $D$ to a qABox is "almost" a homomorphism,[5] but which only maps all those nodes of the syntax tree of $D$ that are between the root and a "cut." Figure 2 shows an example: the "cut" is depicted as the green line. These two conditions are not only necessary but also sufficient for safety.
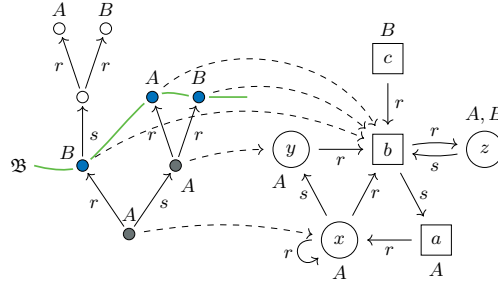


Fig. 2: A partial homomorphism from the concept on the left to the quantified ABox on the right at variable $x$, depicted by the dashed arrows.

**Proposition 10 ([9]).** $\exists X.\mathcal{A}$ *is safe for* $\{P\}$ *iff, for each individual name* $a$, *the following holds: (1) if* $A \in \mathsf{Atoms}(\{P\})$, *then* $A(a) \notin \mathcal{A}$ *and (2) if* $r(a, u) \in \mathcal{A}$ *and* $\exists r.D \in \mathsf{Atoms}(\{P\})$, *then there is no partial homomorphism from* $D$ *to* $\exists X.\mathcal{A}$ *at* $u$.

Since the existence of a partial homomorphism can be decided in polynomial time [9], we obtain the following complexity result.

**Theorem 11.** *Safety of qABox w.r.t. singleton* $\mathcal{EL}$ *policies is in* $\mathsf{P}$.

### 4.1   Canonical Safe Anonymizations

If a qABox turns out not to be safe, we again want to compute an anonymization that is safe and that preserves as much information from the original qABox as possible. We say that a qABox $\exists Y.\mathcal{B}$ is a $\{P\}$-*safe anonymization* of $\exists X.\mathcal{A}$ if

---

[4] Characterizing safety for general policies is an open problem.
[5] Homomorphisms come into play since they characterize the instance problem in $\mathcal{EL}$.

$\exists X.\mathcal{A} \models \exists Y.\mathcal{B}$ and $\exists Y.\mathcal{B}$ is safe for $\{P\}$. Such an anonymization is *optimal* if there is no $\{P\}$-safe anonymization $\exists Z.\mathcal{C}$ of $\exists X.\mathcal{A}$ that lies strictly between $\exists X.\mathcal{A}$ and $\exists Y.\mathcal{B}$ w.r.t. the entailment order. In [9], we presented an approach for computing a unique *optimal safe anonymization* in exponential time. The approach computes a qABox called *canonical safe anonymization* that entails each $\{P\}$-safe anonymization of $\exists X.\mathcal{A}$.

**Definition 12.** *The* canonical safe anonymization $\mathsf{sa}(\exists X.\mathcal{A}, \{P\})$ *of* $\exists X.\mathcal{A}$ *w.r.t.* $\{P\}$ *is defined as the qABox* $\exists Y.\mathcal{B}$ *such that*

1. *the set $Y$ consists of the variables $y_{t,\mathcal{K}}$ where $t$ is an object name occurring in $\exists X.\mathcal{A}$ and $\mathcal{K}$ is a subset of $\mathsf{Atoms}(\{P\})$ that does not contain $\sqsubseteq_\emptyset$-comparable atoms, and*
2. *the matrix $\mathcal{B}$ consists of the following assertions:*
   (a) *$A(y_{t,\mathcal{K}})$ if $A(t)$ occurs in $\mathcal{A}$ and $A \notin \mathcal{K}$,*
   (b) *$r(y_{t,\mathcal{K}}, y_{u,\mathcal{L}})$ provided $r(t, u) \in \mathcal{A}$ and, for each $\exists r.C \in \mathcal{K}$, there is $D \in \mathcal{L}$ with $C \sqsubseteq_\emptyset D$,*
   (c) *$r(y_{t,\mathcal{K}}, b)$ if $r(t, b)$ occurs in $\mathcal{A}$ and there is no $\exists r.C \in \mathcal{K}$.*
   *In these conditions, the first object name $y_{t,\mathcal{K}}$ may also stand for an individual name $a$, which is then treated like the variable $y_{a,\mathsf{Max}(\mathsf{Atoms}(\{P\}))}$, where $\mathsf{Max}(\mathcal{K})$ collects the subsumption-maximal elements of $\mathcal{K}$ modulo equivalence.*

As in the case of compliance, the canonical safe anonymizations introduce an exponential number of copies for each object in the input, which may make a computation infeasible in practice.

### 4.2 Making It Smaller Again

Similar to the case of compliant anonymizations, we can reduce the number of variables in the safe anonymization by creating copies only when needed. According to Definition 12, $r(y_{t,\mathcal{K}}, y_{u,\mathcal{L}}) \in \mathcal{B}$ iff $r(t, u) \in \mathcal{A}$ and $\mathcal{L}$ covers $\mathsf{Succ}(\mathcal{K}, r) := \{C \mid \exists r.C \in \mathcal{K}\}$. To compute the optimized safe anonymization, we again produce a sequence $Y_0, \ldots, Y_m$ of subsets of $Y$. Starting with the set $Y_0 := \{y_{t,\emptyset} \mid t \in \Sigma_\mathsf{O}(\exists X.\mathcal{A})\}$, and applying the following two rules exhaustively.

**Safe Anonymization Rule 1.**
   **If** (i) $y_{t,\mathcal{K}}, y_{u,\mathcal{L}} \in Y_i$, (ii) $r(t, u) \in \mathcal{A}$, (iii) $\mathcal{L}$ does not cover $\mathsf{Succ}(\mathcal{K}, r)$, (iv) $\mathcal{M}$ is a cover-minimal set of atoms covering $\mathcal{L} \cup \mathsf{Succ}(\mathcal{K}, r)$, but (v) $y_{u,\mathcal{M}} \notin Y_i$,
   **then** set $Y_{i+1} := Y_i \cup \{y_{u,\mathcal{M}}\}$
**Safe Anonymization Rule 2.**
   **If** (i) $a \in \Sigma_\mathsf{I}$ and $y_{u,\mathcal{L}} \in Y_i$, (ii) $r(a, u) \in \mathcal{A}$, (iii) $\mathcal{L}$ does not cover $\mathsf{Succ}(\mathsf{Max}(\mathsf{Atoms}(\{P\})), r)$, (iv) $\mathcal{M}$ is a cover-minimal set of atoms covering $\mathcal{L} \cup \mathsf{Succ}(\mathsf{Max}(\mathsf{Atoms}(\{P\})), r)$, but (v) $y_{u,\mathcal{M}} \notin Y_i$,
   **then** set $Y_{i+1} := Y_i \cup \{y_{u,\mathcal{M}}\}$

After generating the set of variables, we construct the matrix of the optimized safe anonymization based on Definition 12.
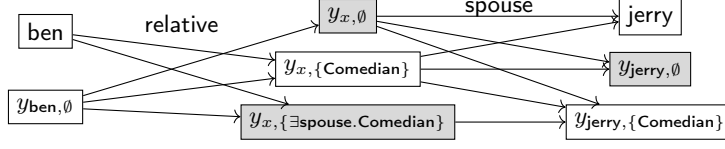
Fig. 3: Optimized safe anonymization of the example ABox. Gray nodes denote instances of Comedian.

**Definition 13.** *Let $Y_m \subseteq Y$ be the set of all variables obtained by exhaustive applications of Safe Anonymization Rule 1 and Rule 2. The optimized $\{P\}$-safe anonymization of $\exists X.\mathcal{A}$ is the qABox $\exists Y_m.\mathcal{B}_m$, where $\mathcal{B}_m$ contains all assertions in the matrix of $\mathsf{sa}(\exists X.\mathcal{A}, \{P\})$ involving only object names in $\Sigma_I \cup Y_m$.*

*Example 14.* For the policy and qABox in Example 9, the canonical safe anonymization would contain 24 variables, while the optimized safe anonymization is much smaller.[6] Applying the Safe Anonymization Rule 2 to the pair $b$ and $y_{x,\emptyset}$ exhaustively, we obtain the variables $y_{x,\{C\}}$ and $y_{x,\{\exists s.C\}}$, and then applying the Safe Anonymization Rule 1 to the pair $y_{x,\{C,\exists s.C\}}$ and $y_{j,\emptyset}$ generates $y_{j,\{C\}}$. On the resulting set of objects, no rule is applicable, and our procedure terminates. Thus, the optimized safe anonymization contains only 8 objects in total. Using the matrix construction in Definition 12, we obtain the optimal safe anonymization $\exists Y_m.\mathcal{B}_m$ whose matrix is depicted in Figure 3.

$\mathcal{B}_m$ is a subset of the matrix of $\mathsf{sa}(\exists X.\mathcal{A}, \{P\})$, which implies that the former is entailed by the latter. It immediately follows that $\exists Y_m.\mathcal{B}_m$ is a $\{P\}$-safe anonymization of $\exists X.\mathcal{A}$. We can also show the other direction.

**Proposition 15.** *The optimized $\{P\}$-safe anonymization of $\exists X.\mathcal{A}$ entails $\mathsf{sa}(\exists X.\mathcal{A}, \{P\})$.*

We thus obtain the following theorem, which shows that we can work with the smaller anonymization.

**Theorem 16.** *Given a qABox $\exists X.\mathcal{A}$ and a singleton policy $\{P\}$, the optimized $\{P\}$-safe anonymization $\exists Y_m.\mathcal{B}_m$ and $\mathsf{sa}(\exists X.\mathcal{A}, \{P\})$ are equivalent.*

### 4.3  Static $\mathcal{EL}$ TBoxes and General Policies

So far, our methods for testing for and achieving safety can only deal with singleton policies without a TBox. Safety w.r.t. a TBox is defined as follows: the qABox $\exists X.\mathcal{A}$ is *safe for $\mathcal{P}$ w.r.t. $\mathcal{T}$* if for each quantified ABox $\exists Y.\mathcal{B}$ that is compliant with $\mathcal{P}$ w.r.t. $\mathcal{T}$, the union $\exists X.\mathcal{A} \cup \exists Y.\mathcal{B}$ is also compliant with $\mathcal{P}$ w.r.t. $\mathcal{T}$. Interestingly, TBoxes can be used to express general policies by singleton policies.

---

[6] To save space and increase legibility, we abbreviate names by their first letters.

**Proposition 17.** *Consider a quantified ABox $\exists X.\mathcal{A}$, an $\mathcal{EL}$ TBox $\mathcal{T}$, and a policy $\mathcal{P}$. Further let $A$ be a fresh concept name not occurring in $\exists X.\mathcal{A}$, in $\mathcal{T}$, or in $\mathcal{P}$, and define the extended TBox $\mathcal{T}_{\mathcal{P}} := \mathcal{T} \cup \{\, P \sqsubseteq A \mid P \in \mathcal{P} \,\}$. Then $\exists X.\mathcal{A}$ is safe for $\mathcal{P}$ w.r.t. $\mathcal{T}$ iff $\exists X.\mathcal{A}$ is safe for $\{A\}$ w.r.t. $\mathcal{T}_{\mathcal{P}}$.*

By setting $\mathcal{T} := \emptyset$ in this proposition, we see that safety for an arbitrary policy $\mathcal{P}$ (but without TBox) can be reduced to safety for the singleton policy $\{A\}$ w.r.t. to a non-empty cycle-restricted TBox. As shown in [9], such a reduction cannot exist without a TBox. Until now, we do not have a characterization of safety akin to Proposition 10 for non-singleton policies without TBox. The proposition shows that dealing with (cycle-restricted) TBoxes, even for singleton policies, is at least as hard as dealing with general policies.

Nevertheless, by using ideas from [17, 18], we can find a coNP decision procedure for safety for a general policy w.r.t. an $\mathcal{EL}$ TBox. This complexity result extends the one given in [9] (Proposition 3.16) for the case without a TBox, and at the same time corrects a typo in the formulation of that proposition.

**Proposition 18.** *The safety problem for general policies w.r.t. static $\mathcal{EL}$ TBoxes is in coNP.*

## 5   Conclusions

The work reviewed in this paper shows that, under some restrictions, optimality can indeed be achieved when computing repairs as well as compliant and safe anonymizations. What remains open is the question of how to deal with general policies and/or cycle-restricted TBoxes in the context of safety. For general TBoxes, optimality is not always achievable, but one can of course ask whether the existence of an optimal repair or an optimal compliant/safe anonymization is decidable, and whether one can then compute such an optimal ABox if it exists. Using conjunctive queries rather then $\mathcal{EL}$ concepts is also an interesting topic for future research.

Classical repairs (which are based on removing axioms) have been used to define inconsistency-tolerant semantics. Basically, instead of replacing an inconsistent ABox by one of its repairs, one reasons w.r.t. all optimal classical repairs in a certain well-defined way [13, 26]. It would be interesting to see what happens if optimal classical repairs are replaced with optimal repairs (in the sense introduced in the present paper) in such inconsistency-tolerant semantics.

## References

1. Baader, F., Borgwardt, S., Morawska, B.: Extending unification in $\mathcal{EL}$ towards general TBoxes. In: Proc. of the 13th Int. Conf. on Principles of Knowledge Representation and Reasoning (KR 2012). pp. 568–572. AAAI Press/The MIT Press (2012)
2. Baader, F., Horrocks, I., Lutz, C., Sattler, U.: An Introduction to Description Logic. Cambridge University Press (2017)
3. Baader, F., Koopmann, P., Kriegel, F., Nuradiansyah, A.: Computing optimal repairs of quantified ABoxes w.r.t. static $\mathcal{EL}$ TBoxes. In: Platzer, A., Sutcliffe, G. (eds.) Proceedings of the 28th International Conference on Automated Deduction (CADE-28). Lecture Notes in Computer Science, vol. 12699, pp. 309–326 (2021)
4. Baader, F., Koopmann, P., Kriegel, F., Nuradiansyah, A.: Computing optimal repairs of quantified ABoxes w.r.t. static $\mathcal{EL}$ TBoxes (extended version). LTCS-Report 21-01, Chair of Automata Theory, Institute of Theoretical Computer Science, Technische Universität Dresden, Dresden, Germany (2021), `https://lat.inf.tu-dresden.de/research/reports/2021/BaKoKrNu-LTCS-21-01.pdf`
5. Baader, F., Koopmann, P., Kriegel, F., Nuradiansyah, A., Peñaloza, R.: Privacy-Preserving Ontology Publishing: The Case of Quantified ABoxes w.r.t. a Static Cycle-Restricted $\mathcal{EL}$ TBox (Extended Version). LTCS-Report 21-04, Chair of Automata Theory, Institute of Theoretical Computer Science, Technische Universität Dresden, Dresden, Germany (2021), `https://lat.inf.tu-dresden.de/research/reports/2021/BaKoKrNuPe-LTCS-21-04.pdf`
6. Baader, F., Kriegel, F., Nuradiansyah, A.: Privacy-preserving ontology publishing for $\mathcal{EL}$ instance stores. In: Calimeri, F., Leone, N., Manna, M. (eds.) Logics in Artificial Intelligence - 16th European Conference, JELIA 2019, Rende, Italy, May 7-11, 2019, Proceedings. Lecture Notes in Computer Science, vol. 11468, pp. 323–338. Springer (2019)
7. Baader, F., Kriegel, F., Nuradiansyah, A., Peñaloza, R.: Making repairs in description logics more gentle. In: Thielscher, M., Toni, F., Wolter, F. (eds.) Principles of Knowledge Representation and Reasoning: Proceedings of the Sixteenth International Conference, KR 2018, Tempe, Arizona, 30 October - 2 November 2018. pp. 319–328. AAAI Press (2018)
8. Baader, F., Kriegel, F., Nuradiansyah, A., Peñaloza, R.: Computing compliant anonymisations of quantified ABoxes w.r.t. $\mathcal{EL}$ policies. In: Pan, J.Z., Tamma, V.A.M., d'Amato, C., Janowicz, K., Fu, B., Polleres, A., Seneviratne, O., Kagal, L. (eds.) The Semantic Web - ISWC 2020 - 19th International Semantic Web Conference, Athens, Greece, November 2-6, 2020, Proceedings, Part I. Lecture Notes in Computer Science, vol. 12506, pp. 3–20. Springer (2020)
9. Baader, F., Kriegel, F., Nuradiansyah, A., Peñaloza, R.: Safety of quantified ABoxes w.r.t. singleton $\mathcal{EL}$ policies. In: Hung, C., Hong, J., Bechini, A., Song, E. (eds.) SAC '21: The 36th ACM/SIGAPP Symposium on Applied Computing. pp. 863–872. ACM (2021)
10. Baader, F., Nuradiansyah, A.: Mixing description logics in privacy-preserving ontology publishing. In: Benzmüller, C., Stuckenschmidt, H. (eds.) KI 2019: Advances in Artificial Intelligence - 42nd German Conference on AI, Kassel, Germany, September 23-26, 2019, Proceedings. Lecture Notes in Computer Science, vol. 11793, pp. 87–100. Springer (2019)
11. Baader, F., Suntisrivaraporn, B.: Debugging SNOMED CT using axiom pinpointing in the description logic $\mathcal{EL}^+$. In: Proceedings of the International Conference

on Representing and Sharing Knowledge Using SNOMED (KR-MED'08). Phoenix, Arizona (2008)

12. Benedikt, M., Grau, B.C., Kostylev, E.V.: Source information disclosure in ontology-based data integration. In: Singh, S.P., Markovitch, S. (eds.) Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA. pp. 1056–1062. AAAI Press (2017)

13. Bienvenu, M., Bourgaux, C., Goasdoué, F.: Computing and explaining query answers over inconsistent DL-Lite knowledge bases. J. Artif. Intell. Res. **64**, 563–644 (2019)

14. Biskup, J., Bonatti, P.A.: Controlled query evaluation for enforcing confidentiality in complete information systems. Int. J. Inf. Sec. **3**(1), 14–27 (2004)

15. Bonatti, P.A., Sauro, L.: A confidentiality model for ontologies. In: Alani, H., Kagal, L., Fokoue, A., Groth, P., Biemann, C., Parreira, J.X., Aroyo, L., Noy, N.F., Welty, C., Janowicz, K. (eds.) The Semantic Web - ISWC 2013 - 12th International Semantic Web Conference, Sydney, NSW, Australia, October 21-25, 2013, Proceedings, Part I. Lecture Notes in Computer Science, vol. 8218, pp. 17–32. Springer (2013)

16. Grau, B.C., Kharlamov, E., Kostylev, E.V., Zheleznyakov, D.: Controlled query evaluation for datalog and OWL 2 profile ontologies. In: Yang, Q., Wooldridge, M.J. (eds.) Proceedings of the Twenty-Fourth International Joint Conference on Artificial Intelligence (IJCAI 2015). pp. 2883–2889. AAAI Press (2015)

17. Grau, B.C., Kostylev, E.V.: Logical foundations of privacy-preserving publishing of linked data. In: Schuurmans, D., Wellman, M.P. (eds.) Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, February 12-17, 2016, Phoenix, Arizona, USA. pp. 943–949. AAAI Press (2016)

18. Grau, B.C., Kostylev, E.V.: Logical foundations of linked data anonymisation. J. Artif. Intell. Res. **64**, 253–314 (2019)

19. Horridge, M., Parsia, B., Sattler, U.: Laconic and precise justifications in OWL. In: Sheth, A.P., Staab, S., Dean, M., Paolucci, M., Maynard, D., Finin, T.W., Thirunarayan, K. (eds.) The Semantic Web - ISWC 2008, 7th International Semantic Web Conference, ISWC 2008, Karlsruhe, Germany, October 26-30, 2008. Proceedings. Lecture Notes in Computer Science, vol. 5318, pp. 323–338. Springer (2008)

20. Horrocks, I., Li, L., Turi, D., Bechhofer, S.: The instance store: DL reasoning with large numbers of individuals. In: Haarslev, V., Möller, R. (eds.) Proceedings of the 2004 International Workshop on Description Logics (DL2004), Whistler, British Columbia, Canada, June 6-8, 2004. CEUR Workshop Proceedings, vol. 104. CEUR-WS.org (2004)

21. Kalyanpur, A., Parsia, B., Horridge, M., Sirin, E.: Finding all justifications of OWL DL entailments. In: Proc. of ISWC'07. Lecture Notes in Computer Science, vol. 4825, pp. 267–280. Springer-Verlag (2007)

22. Lam, J.S.C., Sleeman, D.H., Pan, J.Z., Vasconcelos, W.W.: A fine-grained approach to resolving unsatisfiable ontologies. J. Data Semant. **10**, 62–95 (2008)

23. Lembo, D., Rosati, R., Savo, D.F.: Revisiting controlled query evaluation in description logics. In: Kraus, S. (ed.) Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence (IJCAI 2019). pp. 1786–1792. ijcai.org (2019)

24. Meyer, T., Lee, K., Booth, R., Pan, J.Z.: Finding maximally satisfiable terminologies for the description logic $\mathcal{ALC}$. In: Proc. of the 21st Nat. Conf. on Artificial Intelligence (AAAI 2006). AAAI Press/The MIT Press (2006)

25. Parsia, B., Sirin, E., Kalyanpur, A.: Debugging OWL ontologies. In: Ellis, A., Hagino, T. (eds.) Proc. of the 14th International Conference on World Wide Web (WWW'05). pp. 633–640. ACM (2005)
26. Rosati, R.: On the complexity of dealing with inconsistency in description logic ontologies. In: Walsh, T. (ed.) Proceedings of the 22nd International Joint Conference on Artificial Intelligence (IJCAI 2011). pp. 1057–1062. IJCAI/AAAI (2011)
27. Schlobach, S., Cornet, R.: Non-standard reasoning services for the debugging of description logic terminologies. In: Gottlob, G., Walsh, T. (eds.) Proc. of the 18th Int. Joint Conf. on Artificial Intelligence (IJCAI 2003). pp. 355–362. Morgan Kaufmann, Los Altos, Acapulco, Mexico (2003)
28. Schlobach, S., Huang, Z., Cornet, R., Harmelen, F.: Debugging incoherent terminologies. J. Automated Reasoning **39**(3), 317–349 (2007)
29. Troquard, N., Confalonieri, R., Galliani, P., Peñaloza, R., Porello, D., Kutz, O.: Repairing ontologies via axiom weakening. In: McIlraith, S.A., Weinberger, K.Q. (eds.) Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence (AAAI-18). pp. 1981–1988. AAAI Press (2018)