Tableau-based Forgetting in ALC Ontologies

Zhe Wang¹ and **Kewen Wang**¹ and **Rodney Topor**¹ and **Xiaowang Zhang**²

Abstract. In this paper, we propose two new approaches to forgetting for \mathcal{ALC} based on the well-known tableau algorithm. The first approach computes the result of forgetting by rolling up tableaux, and also provides a decision algorithm for the existence of forgetting in \mathcal{ALC} . When the result of forgetting does not exist, we provide an incremental method for computing approximations of forgetting. This second approach uses variable substitution to refine approximations of forgetting and eventually obtain the result of forgetting. This approach is capable of preserving structural information of the original ontologies enabling readability and comparison. As both approaches are based on the tableau algorithm, their implementations can make use of the mechanisms and optimization techniques of existing description logic reasoners.

1 Introduction

An ontology is a formal definition for a common vocabulary (or signature) that is used to describe and represent an application domain. Ontologies can be used by automated tools to provide advanced services such as more accurate web search, intelligent software agents and knowledge management. An example of large biomedical ontology is SNOMED CT [15]. Ontology editing and maintaining tools, such as Protégé [14], are supported by efficient reasoners based on tableau algorithms [10] for description logics (DLs). However, as shown in [3], the existing reasoners provide limited reasoning support for ontology modifications, which largely restricts the wide use of ontologies in the Semantic Web.

For ontology modifications, an essential task is to reduce the vocabulary of an ontology \mathcal{T} , that is, to *forget* a sub-vocabulary S of the vocabulary of \mathcal{T} and transform \mathcal{T} into a new ontology \mathcal{T}' containing no symbol in S and sharing with \mathcal{T} the same logical consequences that do not use S.

In AI and mathematical logic, *forgetting* or *uniform interpolation* has been well investigated in classical logics [7, 8], modal logics [11], and logic programming [4]. More recently, the technique of forgetting has been proposed for ontology module extraction and ontology reuse [12, 6, 5]. Forgetting was also shown to be useful in ontology composition, decomposition, revision and summarization.

Algorithms for forgetting in simple DLs have been developed [12, 5]. Forgetting for the more complex DL ALC was investigated in [13], where an algorithm based on ALC concept normal form was proposed. However, the problem of forgetting for ALC remained unsolved, for the following reasons. (1) The result of forgetting in an ALC TBox may not exist, and the decidability of the existence of forgetting was open. (2) When the result of forgetting does not exist, an incremental algorithm for computing approximations of forgetting

¹ Griffith University, Australia. {jack.wang, k.wang, r.topor}@griffith.edu.au
² Peking University, China. zxw@is.pku.edu.cn

was missing. (3) As the algorithm transforms each TBox into a normal form, the structural information of the original TBox is lost after forgetting. (4) The computation cannot make use of off-the-shelf DL reasoners which are based on tableau algorithms.

In this paper, we propose two different approaches for forgetting in ALC based on the tableau algorithm. We first introduce a calculus called *rolling up* of tableaux, and show its applications to computing forgetting in concept descriptions and TBoxes. We provide an algorithm to decide the existence of forgetting in TBoxes, and to compute the result of forgetting whenever it exists. The algorithm can also be used to compute approximations of forgetting in an incremental manner. However, this approach is unable to preserve structural information of the original TBox. Inspired by a method used for computing least common subsumers in [2], we provide a different approach for computing forgetting, by introducing a general tableau-based calculus via substitutions on concept terms containing concept variables. To this end, a set of substitution rules are introduced, and then algorithms are developed for forgetting both in concept descriptions and in TBoxes. When the result of TBox forgetting exists, the algorithm is capable of computing the result, in which the original structural information is preserved.

2 Preliminaries

In this section, we briefly recall some basics of ALC. The reader is referred to [1] for further details.

A concept description (or concept) in ALC is built up with concept names and role names. N_C denotes the set of concept names and N_R the set of role names. The syntax of ALC -concept descriptions is defined inductively as follows.

$$\begin{array}{c} B \longleftarrow A \mid \top \mid \bot \\ C, D \longleftarrow B \mid \neg C \mid C \sqcap D \mid C \sqcup D \mid \exists r.C \mid \forall r.C \end{array}$$

where $A \in N_C$ and $r \in N_R$. B is called an *atomic concept*. A *literal* is an atomic concept or its negation. Given a literal L, L^+ denotes its concept name.

A concept is in *negation normal form* (NNF) if negations occur only in front of concept names. A concept can be transformed into its NNF in linear time. In the rest of the paper, without loss of generality, we assume all concept descriptions are in NNF.

A *TBox* is a set of axioms of the form $C \sqsubseteq D$ (*C* is *subsumed* by *D*). $C \equiv D$ is the abbreviation of both $C \sqsubseteq D$ and $D \sqsubseteq C$.

The signature of a concept description C, written sig(C), is the set of all concept and role names in C. Similarly, we can define $sig(\mathcal{T})$ for a TBox \mathcal{T} .

The tableau based approach for DL reasoning is well established, and is the basis for most DL reasoners.

A *tableau* **T** is a set of trees $\{\mathfrak{T}_1, \ldots, \mathfrak{T}_n\}$. Each node x in \mathfrak{T}_i is labeled with a set of concepts $\mathcal{L}_i(x)$, and each edge $\langle x, y \rangle$ in \mathfrak{T}_i is

labeled with a set of roles $\mathcal{L}_i(\langle x, y \rangle)$. When $\langle x, y \rangle$ is labeled with a set containing role name r, we say x is a *r*-predecessor of y, and y a *r*-successor of x.

The initial state of the tableau algorithm is a labeled root node x, denoted $\mathcal{L}_0(x) = \{C_1, \ldots, C_n\}$, which is then expanded by tableau expansion rules (T-rules, ref. Table 1) [1].

Table 1. Tableau expansion rules for \mathcal{ALC} (T-rules)

□-rule:	if	$C_1 \sqcap C_2 \in \mathcal{L}(x)$, and $\{C_1, C_2\} \not\subseteq \mathcal{L}(x)$
	then	set $\mathcal{L}(x) = \mathcal{L}(x) \cup \{C_1, C_2\}$
⊔-rule:	if	$C_1 \sqcup C_2 \in \mathcal{L}(x), \text{ and } \{C_1, C_2\} \cap \mathcal{L}(x) = \emptyset$
	then	create a copy of the tree with label function \mathcal{L}' ,
		and set $\mathcal{L}'(x) = \mathcal{L}(x) \cup \{C_1\}$ and
		$\mathcal{L}(x) = \mathcal{L}(x) \cup \{C_2\}$
∃-rule:	if	$\exists r. C \in \mathcal{L}(x), \text{ and }$
		x has no r-successor y with $C \in \mathcal{L}(y)$
	then	create a new node y, and set $\mathcal{L}(\langle x, y \rangle) = \{r\}$
		and $\mathcal{L}(y) = \{C\}$
∀-rule:	if	$\forall r.C \in \mathcal{L}(x), \text{ and }$
		there is an r-successor y of x with $C \notin \mathcal{L}(y)$
	then	set $\mathcal{L}(y) = \mathcal{L}(y) \cup \{C\}$
\Box -rule:	if	$C_1 \sqsubseteq C_2 \in \mathcal{T}$, and $\neg C_1 \sqcup C_2 \notin \mathcal{L}(x)$
	then	set $\mathcal{L}(x) = \mathcal{L}(x) \cup \{\neg C_1 \sqcup C_2\}$

Note that \sqcup -rule splits a tree \mathfrak{T} into two, \mathfrak{T} and \mathfrak{T}' . For convenience of marking corresponding nodes and edges in these two trees, we denote corresponding nodes using the same names (*e.g.*, *x* in the table). We assume each new node created by other T-rules has a new name. A subtree whose root is named by *x* is called a *x*-rooted subtree.

A tableau is *complete* if no T-rule applies to it. A label is said to *clash* if it contains complementary literals A and $\neg A$ or \bot . A complete tableau **T** is *open* if at least one tree in **T** is clash free. Otherwise, **T** is *closed*. $T \models C \sqsubseteq D$ iff the complete tableau of expanding $\mathcal{L}_0 = \{C \sqcap \neg D\}$ w.r.t. T is closed. If $T = \emptyset$, we simply write $\models C \sqsubseteq D$.

Two notions of forgetting for ALC are defined in [13], one for concept descriptions and one for TBoxes.

Definition 1 (C-Forgetting) Let C be an ALC -concept description and S a set of concept and role names. A concept description C' on the signature sig(C) - S is a result of c-forgetting about S in C if the following conditions are satisfied:

 $\begin{array}{l} (\mathbf{CF1}) \models C \sqsubseteq C'. \\ (\mathbf{CF2}) \models C \sqsubseteq D \text{ implies} \models C' \sqsubseteq D \text{ for any } \mathcal{ALC} \text{ -concept } D \text{ with} \\ \operatorname{sig}(D) \subseteq \operatorname{sig}(C) - \mathcal{S}. \end{array}$

Given a signature S, the result of c-forgetting about S in any ALC -concept C always exists and is unique up to concept equivalence, denoted cforget(C, S).

In [13], only concept descriptions with finite expressions are considered, but the definition can be extended to infinite cases.

Definition 2 (TBox Forgetting) Let T be an ALC -TBox and S be a set of concept and role names. A TBox T' over the signature sig(T) - S is a result of TBox forgetting about S in T if the following conditions are satisfied:

(TF1) $\mathcal{T} \models \mathcal{T}'$; **(TF2)** $\mathcal{T} \models C \sqsubseteq D$ implies $\mathcal{T}' \models C \sqsubseteq D$ for any \mathcal{ALC} -concepts C, D s.t. $\operatorname{sig}(C \sqcup D) \subseteq \operatorname{sig}(\mathcal{T}) - \mathcal{S}$.

It is shown in [13] that some TBoxes may not have finite result of forgetting, even if S contains only concept names. If the result

of forgetting about S in T is expressible as a (finite) ALC TBox, the result of forgetting is unique up to TBox equivalence, denoted forget(T, S). In this case, we say S is *forgettable* in T.

The *existence problem* of forgetting is, for given S and T, to decide whether S is forgettable in T.

3 Forgetting by Rolling Up

In this section, we introduce a forgetting algorithm based on tableau expansion and *rolling up* techniques. We first introduce rolling up of tableaux for \mathcal{ALC} concept descriptions, and show its application in computing c-forgetting. After that, we introduce the algorithm for computing TBox forgetting based on the rolling up techniques for c-forgetting. With the help of tableau theory and rolling up, we also show the decidability of the existence of forgetting for \mathcal{ALC} .

3.1 Forgetting in Concept Descriptions

The tableau algorithm for ALC expands a (collection of the conjuncts of) concept C into a complete tableau **T**, called a *tableau* of C. An observation is that **T** contains the complete information about C, so that we can restore the concept C from **T**, up to concept equivalence. We call such a process the *rolling up* of a tableau.

Before presenting the formal definition of rolling up, we need to first introduce an additional expansion rule. It is because \forall -rule is only applicable when the node has a successor, otherwise the universal quantified concept C in $\forall r.C$ cannot be expanded. In order to fully expand a tableau including all universal quantified concepts, we introduce a new role r_{\forall} for each role r, and a new expansion rule \forall *-rule as follows.

 $\forall *\text{-rule If } \forall r.C \in \mathcal{L}(x), \text{ and } (1) \text{ if } x \text{ has no } r_{\forall}\text{-successor, then} \\ \text{create a new node } y, \text{ and set } \mathcal{L}(\langle x, y \rangle) = \{r_{\forall}\} \text{ and } \mathcal{L}(y) = \{C\}; \\ \text{or } (2) \text{ if } y \text{ is the } r_{\forall}\text{-successor of } x \text{ and } C \notin \mathcal{L}(y), \text{ then set } \mathcal{L}(y) = \\ \mathcal{L}(y) \cup \{C\}. \end{cases}$

Now a tableau is complete if no T-rule or \forall *-rule applies to it. A complete tableau **T** is open if there is a tree in **T** that is either clash free or with clashes occurring only in subtrees whose roots are r_{\forall} -successors. Otherwise, **T** is closed. The \forall *-rule does not affect the termination and correctness of the tableau algorithm.

Now we introduce the definition of rolling up for a tableau.

Definition 3 (Rolling Up) *Given a signature* S *and a tableau* $\mathbf{T}_x = \{\mathfrak{T}_1, \ldots, \mathfrak{T}_n\}$ where each \mathfrak{T}_i is x-rooted, we define the rolling up of \mathbf{T}_x over S to be a concept

$$\operatorname{roll}(\mathbf{T}_{x}, \mathcal{S}) = \bigsqcup_{1 \leq i \leq n} \left(\begin{array}{cc} \prod_{L^{+} \in \mathcal{S}, L \in \mathcal{L}_{i}(x)} & L \sqcap \\ \prod_{r \in \mathcal{S}, r \in \mathcal{L}_{i}(\langle x, y \rangle)} & \exists r.\operatorname{roll}(\mathbf{T}_{y}, \mathcal{S}) \sqcap \\ \prod_{r \in \mathcal{S}, r \neq \in \mathcal{L}_{i}(\langle x, z \rangle)} & \forall r.\operatorname{roll}(\mathbf{T}_{z}, \mathcal{S}) \end{array} \right)$$

where L is a literal with L^+ its concept name, and \mathbf{T}_y and \mathbf{T}_z are the sets of y-rooted and z-rooted subtrees in \mathbf{T}_x , respectively.

Example 1 Let **T** be the complete tableau of expanding $\mathcal{L}_0(x) = \{ (A \sqcup \exists r. \neg B) \sqcap \forall r. (B \sqcup C) \}$. Then $\mathsf{roll}(\mathbf{T}, \{A, B, C, r\})$ is $(A \sqcap \forall r. (B \sqcup C)) \sqcup (\exists r. (\neg B \sqcap C) \sqcap \forall r. (B \sqcup C))$, and $\mathsf{roll}(\mathbf{T}, \{A, C, r\})$ is $A \sqcup \exists r. C$.

The following result states the correctness of rolling up. That is, the rolling up of an arbitrary tableau of a concept is equivalent to the concept itself. **Proposition 1** Given an ALC concept C, let **T** be an arbitrary complete tableau obtained by expanding $\mathcal{L}_0(x) = \{C\}$. Then we have $\models C \equiv \operatorname{roll}(\mathbf{T}, \operatorname{sig}(C))$.

We can define an equivalence relation over all complete tableaux: two complete tableaux $\mathbf{T}_1, \mathbf{T}_2$ are *equivalent* over signature S if \models roll(\mathbf{T}_1, S) \equiv roll(\mathbf{T}_2, S). Given an \mathcal{ALC} -concept C, define $\mathbf{T}(C)$ to be a tableau obtained by expanding $\mathcal{L}_0(x) = \{C\}$ and removing all closed trees. We can show that $\mathbf{T}(C)$ is unique up to equivalence. Note that $\mathbf{T}(C) = \emptyset$ if $\models C \equiv \bot$. The fact that $\mathbf{T}(C)$ contains no closed trees is necessary for the correctness of the following theorem.

Theorem 1 Given an ALC concept C and a set S of concept and role names, we have $cforget(C, S) = roll(\mathbf{T}(C), sig(C) - S)$.

3.2 Forgetting in TBoxes

With non-empty TBoxes, the tableau algorithm requires the \sqsubseteq -rule, which may cause the tableaux to be infinite. An example is $\mathcal{T} = \{\top \sqsubseteq \exists r.B\}$. As $\exists r.B$ is added into the label of each node, the \exists -rule may apply indefinitely. A *blocking* condition is used to ensure termination of the tableau algorithm. However, for the tableau of concept C w.r.t. a TBox \mathcal{T} to capture the complete information of C and \mathcal{T} , in this subsection we assume blocking is not applied.

Since the rolling up of an infinite tableau is an infinite expression, we generalize the classical definition of concept description to allow *infinite concepts*, which are concept descriptions with infinite expressions. The subsumption between infinite concepts is a natural extension of that between finite concepts. We extend Definition 1 to infinite concepts by allowing both C and C' to be possibly infinite. The results for c-forgetting in the previous subsection, *e.g.*, Theorem 1, apply to infinite concepts.

In what follows, we show that TBox forgetting can be characterized by c-forgetting in infinite cases, by introducing a concept encoding for the \Box -rule.

Given a TBox \mathcal{T} , define concept $con(\mathcal{T}) = \prod_{C \sqsubseteq D \in \mathcal{T}} (\neg C \sqcup D)$. Note that each TBox \mathcal{T} can be transformed into an equivalent TBox of the form $\{\top \sqsubseteq con(\mathcal{T})\}$, and thus can be uniquely characterized by the concept $con(\mathcal{T})$. From the finiteness of TBoxes, $con(\mathcal{T})$ is always finite.

Given a finite concept C and a number $n \ge 0$, define

$$C^{(n)} = \prod_{k=0}^{n} \prod_{r_1,\ldots,r_k \in \mathcal{R}} \forall r_1 \cdots \forall r_k.C,$$

where \mathcal{R} is the set of role names in C (and other concepts considered). Note that $C^0 = C$, and $\models C^{(n+1)} \sqsubseteq C^{(n)}$ for any $n \ge 0$. $C^{(\infty)}$ is the limit of $C^{(n)}$ when n goes to ∞ , which is an infinite concept description.

Intuitively, applying \forall -rule to $con(\mathcal{T})^{(n)}$ in the label of a root node adds $con(\mathcal{T})$ to the labels of all the node within depth n of the tableau, as exactly \sqsubseteq -rule does. Thus, for any concept C, the expansion of $\mathcal{L}_0(x) = \{C \sqcap con(\mathcal{T})^{(n)}\}$ without \sqsubseteq -rule imitates the expansion of $\mathcal{L}'_0(x) = \{C\}$ with \sqsubseteq -rule applied to all the nodes within depth n. In what follows, we use $\mathbf{T}(C \sqcap con(\mathcal{T})^{(\infty)})$ to denote the tableau of concept C w.r.t. \mathcal{T} (with \sqsubseteq -rule), and in this way, each label is finite in the tableau. From Theorem 1, $cforget(C \sqcap con(\mathcal{T})^{(\infty)}, \mathcal{S})$ is the rolling up of $\mathbf{T}(C \sqcap con(\mathcal{T})^{(\infty)})$ over $sig(C) \cup sig(\mathcal{T}) - \mathcal{S}$.

The following lemma is an extension of Lemma 9 in [9] to the infinite case.

Lemma 1 Let \mathcal{T} be an \mathcal{ALC} -TBox, and C_1 and C_2 be two finite \mathcal{ALC} -concepts. Then $\mathcal{T} \models C_1 \sqsubseteq C_2$ iff $\models con(\mathcal{T})^{(\infty)} \sqcap C_1 \sqsubseteq C_2$.

The following result connects TBox forgetting with c-forgetting, and also shows the decidability of the existence of TBox forgetting.

Theorem 2 Let T be an ALC -TBox and S be a set of concept and role names. Then

(1) S is forgettable in T iff there exists a number $n \ge 0$ such that forget $(T, S) = \{ \top \sqsubseteq cforget(con(T)^{(n)}, S) \}.$

(2) The existence problem of TBox forgetting is decidable.

To see the correctness of Theorem 2, we first show an equivalent definition of TBox forgetting using concept relations in infinite case. From Lemma 1, (TF1) is $\mathcal{T} \models \top \sqsubseteq con(\mathcal{T}')$ and is equivalent to $\models con(\mathcal{T})^{(\infty)} \sqsubseteq con(\mathcal{T}')$. Since each concept inclusion can be transformed into the form of $\top \sqsubseteq D$, (TF2) is equivalent to say that $\models con(\mathcal{T})^{(\infty)} \sqsubseteq D$ implies $\models con(\mathcal{T}')^{(\infty)} \sqsubseteq D$ for any finite \mathcal{ALC} -concept D with sig $(D) \subseteq sig(\mathcal{T}) - \mathcal{S}$.

Note that $C = con(\mathcal{T}')$ must be finite and over sig $(\mathcal{T}) - S$. From (CF1) and (CF2) of c-forgetting, we have the following result.

Proposition 2 Given an ALC -TBox T and a set S of concept and role names, S is forgettable in T iff there exists a finite concept description C over sig(T) - S satisfying

$$(E1) \models \mathsf{cforget}(con(T)^{(\infty)}, \mathcal{S}) \sqsubseteq C, and$$

(E2) $\models C^{(\infty)} \sqsubseteq \mathsf{cforget}(con(\mathcal{T})^{(\infty)}, \mathcal{S}).$

In this case, $forget(\mathcal{T}, \mathcal{S}) = \{\top \sqsubseteq C\}.$

(E1) and (E2) correspond to (TF1) and (TF2), respectively. To see (E2), note that (TF2) requires $\models con(\mathcal{T}')^{(\infty)} \sqsubseteq cforget(con(\mathcal{T})^{(n)}, \mathcal{S})$ for all $n \ge 0$.

From Proposition 2, the existence of TBox forgetting is reduced to the existence of a finite concept C satisfying (E1) and (E2). The following result shows that the form of such a concept can be further restricted.

Lemma 2 Let \mathcal{T} and \mathcal{S} be as in Proposition 2. If there exists a finite concept satisfying (E1) and (E2), then there is a number $n \ge 0$ such that $C = \mathsf{cforget}(con(\mathcal{T})^{(n)}, \mathcal{S})$ satisfies (E1) and (E2).

Note that $C = \text{cforget}(con(\mathcal{T})^{(n)}, \mathcal{S})$ is finite for any n and always satisfies (E1). The existence of TBox forgetting is further reduced to the existence of a number $n \ge 0$ such that $C = \text{cforget}(con(\mathcal{T})^{(n)}, \mathcal{S})$ satisfies (E2).

We use sub(C) to denote the negation closure of the set of all subconcepts occurring in C. And $sub(\mathcal{T})$ is $sub(con(\mathcal{T}))$. It is clear that each label $\mathcal{L}(x)$ in tableau $\mathbf{T}(C^{(\infty)} \sqcap D)$ satisfies $\mathcal{L}(x) \subseteq sub(C \sqcup D)$. D). There are at most $2^{|sub(C \sqcup D)|}$ different labels in $\mathbf{T}(C^{(\infty)} \sqcap D)$. We can a tableau \mathbf{T} is closed at death n for $n \geq 0$ if each trac in

We say a tableau T is closed at depth n for $n \ge 0$ if each tree in T is closed with a clash within depth n.

Lemma 3 Let C, D be two finite ALC -concepts, S a set of concept and role names, and $N = 2^{|sub(C \sqcup D)|}$. Then $\models C^{(\infty)} \sqsubseteq$ $cforget(D^{(\infty)}, S)$ iff the tableau of $C^{(N)} \sqcap \neg cforget(D^{(N)}, S)$ is closed at depth N.

The above lemma shows that it is decidable for each n whether $C = \text{cforget}(con(\mathcal{T})^{(n)}, \mathcal{S})$ satisfies (E2). And it also suggests that we only need to check up to some large enough n.

Lemma 4 Let C, D be two finite ALC -concepts, S a set of concept and role names, and $N = 2^{|sub(C \sqcup D)|}$. Then for any k > 1, the tableau of $(cforget(C^{(N+k)}, S))^{(N)} \sqcap \neg D$ is closed at depth N iff the tableau of $(cforget(C^{(N)}, S))^{(N)} \sqcap \neg D$ is closed at depth N. The correctness Theorem 2 of is clear from Proposition 2 and Lemmas 2 to 4. These results also provide a fixed number $n = 2^{|sub(\mathcal{T})|}$ to Theorem 2 for deciding the existence and computing TBox forgetting.

Now we introduce an algorithm that decides the existence of forgetting and computes forget(\mathcal{T}, \mathcal{S}) if it exists (ref. Figure 1).

Algorithm 1

Input: An ALC -TBox T, a set S of concept and role names. **Output**: forget(T, S) if S is forgettable in T; otherwise "S is not forgettable in T".

Method: Initially, let n = 0, $\mathbf{T} = \{\mathcal{L}_0(x) = \emptyset\}$, and $\mathcal{T}_0 = \emptyset$. Let $N = 2^{|sub(\mathcal{T})|}$ and $\mathbf{T}' = \mathbf{T}$.

Step 1. Complete \mathbf{T}' with T-rules and $\forall *$ -rule, applying \sqsubseteq -rule only to the nodes within depth N. Compute $D = \operatorname{roll}(\mathbf{T}', \operatorname{sig}(\mathcal{T}) - \mathcal{S})$.

Step 2. Repeat the following steps until n > N:

1. Complete **T** with T-rules and \forall *-rule, applying \sqsubseteq -rule only to the nodes on depth *n*.

2. Compute $C = \operatorname{roll}(\mathbf{T}, \operatorname{sig}(\mathcal{T}) - \mathcal{S})$ and let $\mathcal{T}_{n+1} = \{\top \sqsubseteq C\}$.

3. If $\mathcal{T}_n \models \mathcal{T}_{n+1}$ and the tableau of $C^{(N)} \sqcap \neg D$ is closed at depth N, then return \mathcal{T}_n as forget $(\mathcal{T}, \mathcal{S})$.

4. Assign n = n + 1.

Step 3. Return "S is not forgettable in T".

Figure 1. Compute and approximate TBox forgetting via rolling up.

Step 2 of Algorithm 1 computes incrementally for $0 \le n \le N$, $C = cforget(con(\mathcal{T})^{(n)}, S)$. We do not assign n = N directly because in most of the cases, the existence of forgetting can be decided (with the result also computed) with very small n. In 3 of Step 2, the algorithm checks whether C satisfies (E2) only when \mathcal{T}_{n+1} is not strictly stronger (*w.r.t.* logical consequences) than \mathcal{T}_n . And if it is the case, \mathcal{T}_n is the result of forgetting and does not change (logically) for increasing n. If n exceeds N and no concept C satisfying (E2) is found, then the result of forgetting does not exist.

The correctness of Algorithm 1 is easily seen from the previous discussion.

Theorem 3 Given an ALC -TBox T and a set S of concept and role names, then Algorithm 1 returns forget(T, S) if it exists, and returns "S is not forgettable in T" otherwise.

Algorithm 1 is also an anytime approximation algorithm for TBox forgetting, as each T_n obtained is an approximation of the result of forgetting. We have the following results for T_n .

Proposition 3 Given an ALC -TBox T and a set S of concept and role names, for any $n \ge 0$, we have

- *1.* $\mathcal{T} \models \mathcal{T}_{n+1} \models \mathcal{T}_n$.
- 2. $\mathcal{T} \models C \sqsubseteq D$ implies $\mathcal{T}_n \models C \sqsubseteq D$ for any \mathcal{ALC} -concepts C, Ds.t. sig $(C \sqcup D) \subseteq$ sig $(\mathcal{T}) - S$ and $2^{|sub(C \sqcup D) \cup sub(\mathcal{T})|} < n$.

Each \mathcal{T}_{n+1} is a better approximation of the result of forgetting than \mathcal{T}_n , and \mathcal{T}_{n+1} can be computed by further expanding the tableau for computing \mathcal{T}_n . Thus, Algorithm 1 is also an incremental approximation algorithm for TBox forgetting.

4 Forgetting by Variable Substitution

As easily seen, the forgetting algorithm in the previous section does not preserve the structure of the initial TBox. However, in many applications, it is desirable to preserve the structure for readability and comparison. By introducing concept variables into the TBox, the following approach is capable of preserving structural information in the initial TBox during forgetting.

Let N_X be the set of *concept variables* (or variables), which can be quantified over and can be replaced with concepts in ALC. A *concept term* is is a concept description that allows concept variables to be used as atomic concepts. E, F denote concept terms. We assume all concept terms are also in NNF, *i.e.*, negations only occur in front of concept names and variables. We generalize the syntax of TBoxes to allow variables and axioms of the form $E \sqsubseteq F$.

A substitution σ is a set of pairs of the form $X \mapsto E$ with $X \in N_X$ and E a concept term. A substitution is ground if every E contains no variable. We say σ is *over* a signature S, if E contains only concept and role name in S. Given a concept term F, a TBox \mathcal{T} or a tableau **T** containing variables, $\sigma(F)$, $\sigma(\mathcal{T})$ and $\sigma(\mathbf{T})$ can be defined in a natural way.

4.1 Forgetting in Concept Descriptions

In this approach, we start with a general concept D over sig(C) - S such that $\models C \sqsubseteq D$, and replace D with stronger and stronger concepts (*w.r.t.* subsumption) to approximate forget(C, S).

A natural way of strengthening concepts is by introducing new conjuncts. In [2], a notion of *concept decoration* is defined. Given an ALC -concept description C in NNF, dec(C) is a concept term defined inductively as follows:

- if C is a literal, dec(C) = C;
- if C is $C_1 \sqcap C_2$, $dec(C) = dec(C_1) \sqcap dec(C_2)$; if C is $C_1 \sqcup C_2$, $dec(C) = dec(C_1) \sqcup dec(C_2)$;
- if C is $\exists r.E, dec(C) = \exists r.(X \sqcap E)$ with X a new variable; if C is $\forall r.E, dec(C) = \forall r.(X \sqcap E)$ with X a new variable.

Define $C_{dec} = X_0 \sqcap dec(C)$ with X_0 a new variable.

The decoration process simply adds one variable conjunct under each quantifier in the concept description. For example, the decoration of concept description $(A \sqcup \exists r. \neg B) \sqcap \forall r. (B \sqcup C)$ is $X_0 \sqcap (A \sqcup \exists r. (X_1 \sqcap \neg B)) \sqcap \forall r. (X_2 \sqcap (B \sqcup C)).$

Since variables are added as conjuncts, it is easy to see that $\models \sigma(C_{dec}) \sqsubseteq C$ for any concept C and ground substitution σ . Moreover, for any concept D with $\models D \sqsubseteq C$, there exists a ground substitution σ such that $\models \sigma(C_{dec}) \equiv D$. Such a σ can be simply constructed as $\{X_0 \mapsto D\} \cup \{X \mapsto \top \mid X \neq X_0, X \text{ in } C_{dec}\}$.

With the notion of decoration and substitution, we can present an equivalent characterization of c-forgetting as follows.

Proposition 4 Given ALC -concepts C, D and a set S of concept and role names, we have cforget(C, S) = D iff $(1) sig(D) \subseteq$ $sig(C) - S, (2) \models C \sqsubseteq D$, and (3) the following formula is false:

$$\exists \sigma. \{ \sigma \text{ is ground and is over } sig(C) - S, \\ s.t. \models C \sqsubseteq \sigma(D_{dec}) \text{ and } \nvDash D \sqsubseteq \sigma(D_{dec}) \}.$$
(*)

To check whether cforget(C, S) = D, we want to decide whether such a substitution σ satisfying (*) exists or not, and to construct σ if it exists. Using the tableau algorithm, we need to expand two tableaux, \mathbf{T}_{cl} and \mathbf{T}_{op} , which are defined to be $\mathcal{L}_0(x) = \{C, \neg D_{dec}\}$ and $\mathcal{L}'_0(x) = \{D, \neg D_{dec}\}$, respectively.

In order to construct σ , we introduce a set of *Substitution rules* (Srules, ref. Table 2). All the rules are applicable only if \mathcal{L} is the label of an open tree in \mathbf{T}_{cl} . As disjunctions need to be handled, the Srules proposed here are more complex than those in [2]. X, Y (with

 Table 2.
 Substitution rules (S-rules)

U-rule:	if	$\{\neg X, \neg Y\} \subseteq \mathcal{L}(x)$
	then	apply $\sigma = \{X \mapsto Y\}$ to \mathbf{T}_{cl} and \mathbf{T}_{op}
⊤-rule:	if	$\neg X \in \mathcal{L}(x)$
	then	apply $\sigma = \{X \mapsto \top\}$ to \mathbf{T}_{cl} and \mathbf{T}_{op}
L-rule:	if	$\{\neg X, L_i\} \subseteq \mathcal{L}_i(x)$ with $i \in I$ and each $L_i^+ \notin \mathcal{S}$
	then	apply $\sigma = \{X \mapsto \bigsqcup_{i \in I} L_i\}$ to \mathbf{T}_{cl} and \mathbf{T}_{op}
Q-rule:	if	$\{\neg X, L_i\} \subseteq \mathcal{L}_i(x)$ with $i \in I$ and each $L_i^+ \notin \mathcal{S}$,
		$\{\neg X, \exists r_j.C_j\} \subseteq \mathcal{L}_j(x) \text{ with } j \in J \text{ and } r_j \notin \mathcal{S},$
		$\{\neg X, \forall r_k.C_k\} \subseteq \mathcal{L}_k(x) \text{ with } k \in K \text{ and } r_k \notin S$
	then	apply to \mathbf{T}_{cl} and \mathbf{T}_{op} substitution $\sigma =$
		$\{X \mapsto \bigsqcup_{i \in I} L_i \sqcup \bigsqcup_{j \in J} \exists r_j . Y_j \sqcup \bigsqcup_{k \in K} \forall r_k . Y_k\}$
		where each Y_j, Y_k are new variables

subscripts) are concept variables, and *I*, *J*, *K* are mutual disjoint sets of numbers.

Each variable X only occurs negated in $\langle \mathbf{T}_{cl}, \mathbf{T}_{op} \rangle$, and thus we only consider $\neg X$ in each label. When both T-rules and S-rules are applicable, T-rules always have precedence over S-rules. U-rule unifies any two variables in a label, and have precedence over all other S-rules, to ensure that other S-rules apply at most once in each label $\mathcal{L}(x)$. Also, the \top -rule and L-rule have precedence over the Q-rule, as we want to introduce as few new variables as possible. We call a tableau S-complete if no T-rule or S-rule is applicable, and we talk about openness and closeness only for S-complete tableaux.

The following theorem states the soundness and completeness of the S-rules.

Proposition 5 Given ALC -concepts C, D with $\models C \sqsubseteq D$, and a set S of concept and role names s.t. $sig(D) \subseteq sig(C) - S$, then

(1) The application of T-rules and S-rules to \mathbf{T}_{cl} and \mathbf{T}_{op} always terminates; and

(2) Formula (*) holds iff there is a way of applying S-rules to obtain a ground substitution σ , such that $\sigma(\mathbf{T}_{cl})$ is closed and $\sigma(\mathbf{T}_{op})$ is open.

Now with T-rules and S-rules, we are able to show an approximating algorithm for computing c-forgetting (ref. Figure 2). We call a concept description an *S*-literal if it is of the form A or $\neg A$ with $A \in S$, or $\exists r.C$ or $\forall r.C$ with $r \in S$.

Algorithm 2

Input: An ALC -concept C and a set S of concept and role names. **Output**: cforget(C, S).

Method:

Step 1. Let D be the concept obtained from C by replacing all S-literals in C with \top .

Step 2. Repeat the following steps until D does not change: 1. Assign \mathbf{T}_{cl} to be $\mathcal{L}_0(x) = \{C, \neg D_{dec}\}$ and \mathbf{T}_{op} to be $\mathcal{L}'_0(x) = \{D, \neg D_{dec}\}$.

2. Complete \mathbf{T}_{cl} and \mathbf{T}_{op} with T-rules and S-rules.

3. If a ground substitution σ is found s.t. $\sigma(\mathbf{T}_{cl})$ is closed and $\sigma(\mathbf{T}_{op})$ is open, then assign $D = \sigma(D_{dec})$.

Step 3. Return D as cforget(C, S).

Figure 2. Compute c-forgetting via variable substitution

In Step 1 of Algorithm 2, D is taken as the first approximation of cforget(C, S), as it is not hard to verify that $\models C \sqsubseteq D$ and $sig(D) \subseteq sig(C) - S$. In Step 2, D is refined by $\sigma(D_{dec})$ repeatedly until D is the strongest w.r.t. subsumption, which is the result of c-forgetting. Since c-forget always exists for any ALC -concept description C, Algorithm 2 always terminates. The correctness of Algorithm 2 is stated as follows.

Theorem 4 Given an ALC -concept description C and a set S of concept and role names, then Algorithm 2 always terminates and returns cforget(C, S).

4.2 Forgetting in TBoxes

To generalize the approach in the previous subsection to TBox forgetting, we start with a weak TBox \mathcal{T}' over $sig(\mathcal{T}) - S$ such that $\mathcal{T} \models \mathcal{T}'$, and approximate the result of forgetting by replacing \mathcal{T}' with stronger and stronger TBoxes (*w.r.t.* logical consequence).

Intuitively, a TBox axiom is strengthened via introducing new disjuncts into the left-hand side of each axiom, and/or new conjuncts into the right-hand side. We define *TBox decoration* with the help of concept decoration.

Definition 4 (TBox Decoration) Given an ALC -TBox T, dec(T)is obtained from T by replacing each axiom $C \sqsubseteq D$ in T with $\neg dec(\neg C) \sqsubseteq X \sqcap dec(D)$, where each X is a new variable. $T_{dec} = dec(T) \cup \{\top \sqsubseteq X_0\}$ with X_0 being a new variable.

For example, the decoration of TBox { $A \sqcap \exists r.B \sqsubseteq \forall r.C, C \sqsubseteq D$ } is { $\top \sqsubseteq X_0, A \sqcap \exists r.(\neg X_1 \sqcup B) \sqsubseteq X_2 \sqcap \forall r.(X_3 \sqcap C), C \sqsubseteq X_4 \sqcap D$ }.

Lemma 5 Given an ALC -TBox T, we have (1) $\sigma(T_{dec}) \models T$ for any ground substitution σ ; and (2) for any ALC -TBox T' with $T' \models T$, there exists a ground substitution σ such that $\sigma(T_{dec}) \equiv T'$.

Similar to Proposition 4, we have the following characterization for TBox forgetting.

Proposition 6 Given ALC -TBoxes $\mathcal{T}, \mathcal{T}'$ and a set S of concept and role names, then $\text{forget}(\mathcal{T}, S) = \mathcal{T}'$ iff $(1) \operatorname{sig}(\mathcal{T}') \subseteq \operatorname{sig}(\mathcal{T}) - S$, $(2) \mathcal{T} \models \mathcal{T}'$, and (3) the following formula is false:

$$\exists \sigma. \{ \sigma \text{ is ground and is over } sig(\mathcal{T}) - \mathcal{S}, \\ s.t. \ \mathcal{T} \models \sigma(\mathcal{T}'_{dec}), \text{ and } \mathcal{T}' \not\models \sigma(\mathcal{T}'_{dec}) \}. \quad (**)$$

To check whether $\operatorname{forget}(\mathcal{T}, \mathcal{S}) = \mathcal{T}'$, in contrast to c-forgetting, TBox tableaux expansion requires the additional \sqsubseteq -rule, and the classical tableau blocking condition, called T-blocking. \mathbf{T}_{cl} and \mathbf{T}_{op} both are initialized to be $\mathcal{L}_0(x) = \{\bigsqcup_{E \sqsubseteq F \in \mathcal{T}'_{dec}} E \sqcap \neg F\}$. The difference is that they are expanded *w.r.t.* different TBoxes. In particular, \mathbf{T}_{cl} and \mathbf{T}_{op} are expanded *w.r.t.* \mathcal{T} and \mathcal{T}' , respectively. By applying T-rules (including the \sqsubseteq -rule) and S-rules to \mathbf{T}_{cl} and \mathbf{T}_{op} , the algorithm tries to construct a ground substitution σ such that $\sigma(\mathbf{T}_{cl})$ is closed and $\sigma(\mathbf{T}_{op})$ is open.

However, because of the inter-operation of \sqsubseteq -rule and S-rules, the algorithm may not terminate. In particular, when \mathcal{T}' is already the result of forgetting, \sqsubseteq -rule may keep adding concepts of the form $\exists r.C$ into the labels and introducing new nodes which trigger the application of Q-rule. As labels are changed after applying Q-rule, T-blocking may fail.

A blocking condition is needed here, which is similar to that used in [2]. Application of Q-rule with substitution $X \mapsto E$ is S-blocked in $\mathcal{L}(x)$ if in a previous state, Q-rule has been applied with $X' \mapsto E'$ in $\mathcal{L}'(x')$ such that: (1) E = E', (2) $\mathcal{L}(x) = \mathcal{L}'(x')$, and (3) for each r occurring in E and each r-successor y of x, there is a r-successor y' of x' with $\mathcal{L}(y) = \mathcal{L}(y')$. The equations in (1) – (3) are regardless of variable name variations.

The following result states the termination, soundness and completeness of S-rules regarding TBoxes. **Proposition 7** Given ALC -TBoxes T, T' and a set S of concept and role names s.t. $T \models T'$ and $sig(T') \subseteq sig(T) - S$, then

(1) The application of T-rules and S-rules to \mathbf{T}_{cl} and \mathbf{T}_{op} always terminates (with T-blocking and S-blocking); and

(2) Formula (**) holds iff there is a way of applying S-rules to obtain a ground substitution σ , s.t. $\sigma(\mathbf{T}_{cl})$ is closed and $\sigma(\mathbf{T}_{op})$ is open.

Now we present the algorithm for computing TBox forgetting based on T-rules, S-rules, and blocking conditions (ref. Figure 3). We assume all concepts in the TBox are in NNF.

Algorithm 3

Input: An ALC -TBox T and a set S of concept and role names. **Output**: forget(T, S).

Method:

. . .

Step 1. \mathcal{T}' is obtained from \mathcal{T} by replacing all \mathcal{S} -literals on the lefthand sides of the axioms with \bot , and those on the right-hand sides with \top .

Step 2. Repeat the following steps until \mathcal{T}' does not change:

Assign both T_{cl} and T_{op} to be L₀(x) = {∐_{E⊑F∈T'_{dec}} E □ ¬F}.
 Complete T_{cl} and T_{op} by T-rules and S-rules w.r.t. T and T', respectively.

3. If a ground substitution σ is found *s.t.* $\sigma(\mathbf{T}_{cl})$ is closed and $\sigma(\mathbf{T}_{op})$ is open, then assign $\mathcal{T}' = \sigma(\mathcal{T}'_{dec})$. Step 3. Return \mathcal{T}' as forget $(\mathcal{T}, \mathcal{S})$.

Figure 3. Compute TBox forgetting via variable substitution.

Note that Algorithm 3 may not terminate, even with the blocking conditions. In particular, when the result of forgetting does not exist, there is an infinite sequence of stronger and stronger TBoxes derived as better approximations of the result of forgetting. That is, although blocking condition guarantees each execution of Step 2 to terminate, Step 2 can be repetitively executed infinite times.

The following example shows the effect of executing Algorithm 3 in a case where TBox forgetting does not exist.

Example 2 Let $\mathcal{T} = \{ A \sqsubseteq B, B \sqsubseteq \exists r.B \sqcap C \}$ and $\mathcal{S} = \{B\}$. Then Algorithm 3 starts with $\mathcal{T}' = \{ A \sqsubseteq \top, \bot \sqsubseteq \exists r.\top \sqcap C \}$ and $\mathcal{T}'_{dec} = \{ \top \sqsubseteq X_0, A \sqsubseteq X_1, \bot \sqsubseteq X_2 \sqcap \exists r.(X_3 \sqcap \top) \sqcap C \}$.

Denote σ_i and T'_i to be the substitution and resulting TBox, respectively, in the *i*-th iteration of Step 2. We omit the pairs in each σ_i of the form $X \mapsto \top$ and trivial axioms in each T'_i . Then we have

 $\sigma_1 = \{X_1 \mapsto C\} \text{ and } \mathcal{T}'_1 = \{A \sqsubseteq C\}, \text{ and thus } \mathcal{T}'_{dec} = \{\top \sqsubseteq X_0, A \sqsubseteq X_1 \sqcap C\};\$

 $\sigma_2 = \{X_1 \mapsto \exists r. Y, Y \mapsto C\} \text{ and } \mathcal{T}'_2 = \{A \sqsubseteq C \sqcap \exists r. C\}, \text{ and } thus \mathcal{T}'_{dec} = \{\top \sqsubseteq X_0, A \sqsubseteq X_1 \sqcap C \sqcap \exists r. (X_2 \sqcap C) \};$

 $\sigma_3 = \{X_1 \mapsto \top, X_2 \mapsto \exists r.Y, Y \mapsto C\} \text{ and } \widetilde{\mathcal{T}}'_3 = \{A \sqsubseteq C \sqcap \exists r.(C \sqcap \exists r.C)\};$

$$\mathcal{T}'_n = \{ A \sqsubseteq \underbrace{C \sqcap \exists r. (C \sqcap \exists r. (C \dotsm \exists r. C)) \} \text{ for } n \ge 1.}_{n C's} \}$$

However, whenever the result of forgetting exists, the termination and correctness of Algorithm 3 are guaranteed.

Theorem 5 Given ALC -TBox T and a set S of concept and role names, if S is forgettable in T, then Algorithm 3 always terminates and returns forget(T, S).

5 Conclusion

We have presented two approaches for computing the result of forgetting in both ALC concept descriptions and TBoxes, based on the tableau algorithm for ALC. The first approach is based on the technique of rolling up tableaux. Compared to the algorithm introduced in [13], this new method allows successive approximations of forgetting to be computed incrementally, which is desirable when the result of forgetting does not exist. An important application of this method is to show that the existence problem of forgetting in ALC TBoxes is decidable, However, the first method cannot guarantee that the structural information of the original ontology (TBox) is preserved after forgetting. As a result, we have developed a second, different method for forgetting in \mathcal{ALC} . This method consists of running two tableau-based procedures in parallel. The second new method possesses several advantages: (1) it is an incremental computation algorithm; (2) it can be implemented using an off-the-shelf reasoner for ALC; and (3) it preserves the structural information of the original ontologies (TBoxes). For future research, it would be useful to find lower bounds on the complexity of forgetting. It would be also useful to generalize the forgetting algorithms for more expressive DLs than ALC. It would be interesting to implement our algorithms and incorporate them into ontology editors.

Acknowledgements: The authors would like to thank the referees for their helpful and constructive comments. This work was partially supported by the Australia Research Council (ARC) under DP0666107 and DP1093652.

REFERENCES

- F. Baader, D. Calvanese, D. L. McGuinness, D. Nardi, and P. F. Patel-Schneider, eds. *The Description Logic Handbook*. (2003).
- [2] F. M. Donini, S. Colucci, T. Di Noia, and E. Di Sciascio, 'A tableauxbased method for computing least common subsumers for expressive description logics', in *Proc. 21st IJCAI*, pp. 739–745, (2009).
- [3] M. Dzbor, E. Motta, C. Buil, J. M. Gomez, O. Görlitz, and H. Lewen, 'Developing ontologies in OWL: an observational study', in *Proc. Workshop on OWL: Experiences and Directions*, (2006).
- [4] T. Eiter and K. Wang, 'Semantic forgetting in answer set programming', *Artificial Intelligence*, 14, 1644–1672, (2008).
- [5] B. Konev, F. Wolter, and M. Zakharyaschev, 'Forgetting and uniform interpolation in large-scale description logic terminologies', in *Proc.* 20th IJCAI, pp. 830–835, (2009).
- [6] R. Kontchakov, F. Wolter, and M. Zakharyaschev, 'Can you tell the difference between DL-Lite ontologies?', in *Proc. 11th KR*, pp. 285–295, (2008).
- [7] J. Lang, P. Liberatore, and P. Marquis, 'Propositional independence: Formula-variable independence and forgetting.', *J. Artif. Intell. Res.*, 18, 391–443, (2003).
- [8] F. Lin and R. Reiter, 'Forget it', in *Proc. AAAI Fall Symposium on Rel-evance*, pp. 154–159. New Orleans (LA), (1994).
- [9] B. ten Cate, W. Conradie, M. Marx, and Y. Venema, 'Definitorially complete description logics', in *Proc. 10th KR*, pp. 79–89, (2006).
- [10] D. Tsarkov, I. Horrocks, and P. F. Patel-Schneider, 'Optimizing terminological reasoning for expressive description logics', *J. Autom. Reasoning*, **39**(3), 277–316, (2007).
- [11] Albert Visser, 'Uniform interpolation and layered bisimulation', in *Proc. Gödel'96*, pp. 139–164, (1996).
- [12] Z. Wang, K. Wang, R. Topor, and J. Z. Pan, 'Forgetting concepts in DL-Lite', in *Proc. 5th ESWC*, pp. 245–257, (2008).
- [13] K. Wang, Z. Wang, R. Topor, J. Z. Pan, and G. Antoniou, 'Concept and role forgetting in ALC-ontologies', in *Proc. 8th ISWC*, pp. 666–681, (2009).
- [14] 'Protégé', http://protege.stanford.edu, (2010).
- [15] 'SNOMED CT', http://www.fmrc.org.au/snomed/, (2007).