

Inverse Roles Make Conjunctive Queries Hard

Carsten Lutz

Institute for Theoretical Computer Science, TU Dresden, Germany
lutz@tcs.inf.tu-dresden.de

Abstract. Conjunctive query answering is an important DL reasoning task. Although this task is by now quite well-understood, tight complexity bounds for conjunctive query answering in expressive DLs have never been obtained: all known algorithms run in deterministic double exponential time, but the existing lower bound is only an EXPTIME one. In this paper, we prove that conjunctive query answering in \mathcal{ALCT} is 2-EXPTIME-hard (and thus complete), and that it becomes NEXPTIME-complete under some reasonable assumptions.

1 Introduction

When description logic (DL) knowledge bases are used in applications with a large amount of instance data, ABox querying is the most important reasoning problem. The most basic query mechanism for ABoxes is *instance retrieval*, i.e., returning all the individuals from an ABox that are known to be instances of a given query concept. Instance retrieval can be viewed as a well-behaved generalization of subsumption and satisfiability, which are the standard reasoning problems on TBoxes. In particular, algorithms for the latter can typically be adapted to instance retrieval in a straightforward way, and the computational complexity coincides in almost all cases (see [13] for an exception). In 1998, Calvanese et al. introduced *conjunctive query answering* as a more powerful query mechanism for DL ABoxes. Since then, conjunctive queries have received considerable interest in the DL community, see for example the papers [2, 3, 5–8, 12]. In a nutshell, conjunctive query answering generalizes instance retrieval by admitting also queries whose relational structure is not tree-shaped. This generalization is both natural and useful because the relational structure of ABoxes is usually not tree-shaped as well.

In contrast to the case of instance retrieval, developing algorithms for conjunctive query answering is not merely a matter of extending algorithms for satisfiability, but requires developing new techniques. In particular, all hitherto known algorithms for DLs that include \mathcal{ALC} as a fragment run in deterministic double exponential runtime, in contrast to algorithms for deciding subsumption and satisfiability which require only exponential time even for DLs much more expressive than \mathcal{ALC} . Since the introduction of conjunctive query answering as a reasoning problem for DLs, it has remained an open question whether or not this increase in runtime can be avoided. In other words, it has not been clear whether generalizing instance retrieval to the more powerful conjunctive query answering is penalized by higher computational complexity. In this paper, we answer this question by showing that conjunctive query answering is computationally more expensive than instance retrieval when inverse roles are present. More precisely,

we prove the following two results about \mathcal{ALCI} , the extension of \mathcal{ALC} with inverse roles:

(1) Rooted conjunctive query answering in \mathcal{ALCI} is co-NEXPTIME-complete, where *rooted* means that conjunctive queries are required to be connected and contain at least one answer variable. The phrase “rooted” derives from the fact that every match of such a query is rooted in at least one ABox individual. The lower bound even holds for ABoxes of the form $\{C(a)\}$ and w.r.t. empty TBoxes.

(2) Conjunctive query answering in \mathcal{ALCI} is 2-EXPTIME-complete. The lower bound even holds for ABoxes of the form $\{C(a)\}$ and when queries do not contain any answer variables (or when they contain answer variables, but are not connected).

In the conference version of this paper, we will complement these results by showing that the high computational complexity of conjunctive query answering is indeed due to inverse roles. We will show that conjunctive query answering in \mathcal{ALC} and \mathcal{SHQ} , the fragment of \mathcal{SHIQ} without inverse roles, is only EXPTIME-complete. In this abstract, we concentrate on the lower bounds due to space limitations.

2 Preliminaries

We assume standard notation for the syntax and semantics of \mathcal{ALCI} knowledge bases [1]. In particular, a *TBox* is a set of concept inclusions $C \sqsubseteq D$ and a *knowledge base (KB)* is a pair $(\mathcal{T}, \mathcal{A})$ consisting of a TBox \mathcal{T} and an ABox \mathcal{A} . Let \mathbb{N}_V be a countably infinite set of *variables*. An *atom* is an expression $C(v)$ or $r(v, v')$, where C is an \mathcal{ALCI} concept, r is a (possibly inverse) role, and $v, v' \in \mathbb{N}_V$. A *conjunctive query* q is a finite set of atoms. We use $\text{Var}(q)$ to denote the set of variables occurring in the query q . Let \mathcal{A} be an ABox, \mathcal{I} a model of \mathcal{A} , q a conjunctive query, and $\pi : \text{Var}(q) \rightarrow \Delta^{\mathcal{I}}$ a total function. We write $\mathcal{I} \models^{\pi} C(v)$ if $(\pi(v)) \in C^{\mathcal{I}}$ and $\mathcal{I} \models^{\pi} r(v, v')$ if $(\pi(v), \pi(v')) \in r^{\mathcal{I}}$. If $\mathcal{I} \models^{\pi} at$ for all $at \in q$, we write $\mathcal{I} \models^{\pi} q$ and call π a *match* for \mathcal{I} and q . We say that \mathcal{I} *satisfies* q and write $\mathcal{I} \models q$ if there is a match π for \mathcal{I} and q . If $\mathcal{I} \models q$ for all models \mathcal{I} of a KB \mathcal{K} , we write $\mathcal{K} \models q$ and say that \mathcal{K} *entails* q . The *query entailment problem* is, given a knowledge base \mathcal{K} and a query q , to decide whether $\mathcal{K} \models q$. This is the decision problem corresponding to query answering (which is a search problem), see e.g. [6] for details.

3 Rooted Query Entailment in \mathcal{ALCI} is co-NEXPTIME-complete

Let $\mathcal{ALC}^{\text{rs}}$ be the variation of \mathcal{ALC} in which all roles are interpreted as reflexive and symmetric relations. Our proof of the lower bound stated as (1) above proceeds by first polynomially reducing rooted query entailment in $\mathcal{ALC}^{\text{rs}}$ w.r.t. the empty TBox to rooted query entailment in \mathcal{ALCI} w.r.t. the empty TBox. Then, we prove co-NEXPTIME-hardness of rooted query entailment in $\mathcal{ALC}^{\text{rs}}$.

Regarding the first step, we only sketch the basic idea, which is simply to replace each symmetric role r with the composition of r^- and r . Although r is not interpreted in a symmetric relation in \mathcal{ALCI} , the composition of r^- and r is clearly symmetric. To achieve reflexivity, we ensure that $\exists r^- . \top$ is satisfied by all relevant individuals and

for all relevant roles r . Thus, every individual can reach itself by first travelling r^- and then r , which corresponds to a reflexive loop. Since we are working without TBoxes and thus cannot use statements such as $\top \sqsubseteq \exists r^- . \top$, a careful manipulation of the ABox and query is needed. Details are given in appendix A.

Before we prove co-NEXPTIME-hardness of rooted query entailment in $\mathcal{ALC}^{\text{rs}}$, we discuss a preliminary. An interpretation \mathcal{I} of $\mathcal{ALC}^{\text{rs}}$ is *tree-shaped* if there is a bijection f from $\Delta^{\mathcal{I}}$ into the set of nodes of a finite undirected tree (V, E) such that $(d, e) \in s^{\mathcal{I}}$, for some role name s , implies that $d = e$ or $\{f(d), f(e)\} \in E$. The proof of the following result is standard, using unravelling of non-tree-shaped models.

Lemma 1. *If \mathcal{A} is an $\mathcal{ALC}^{\text{rs}}$ -ABox and q a conjunctive query, then $\mathcal{A} \not\models q$ implies that there is a tree-shaped model \mathcal{I} of \mathcal{A} such that $\mathcal{I} \not\models q$.*

Because $\mathcal{A} \models q$ clearly implies that $\mathcal{I} \models q$ for all tree-shaped models \mathcal{I} of \mathcal{A} , this lemma means that we can concentrate on tree-shaped interpretations when deciding conjunctive query entailment. We will exploit this fact to give an easier explanation of the reduction that is to follow.

We now give a reduction from a NEXPTIME-complete variant of the tiling problem to the complement of rooted query entailment in $\mathcal{ALC}^{\text{rs}}$.

Definition 1 (Domino System). A domino system \mathfrak{D} is a triple (T, H, V) , where $T = \{0, 1, \dots, k-1\}$, $k \geq 0$, is a finite set of tile types and $H, V \subseteq T \times T$ represent the horizontal and vertical matching conditions. Let \mathfrak{D} be a domino system and $c = c_0, \dots, c_{n-1}$ an initial condition, i.e. an n -tuple of tile types. A mapping $\tau : \{0, \dots, 2^{n+1}-1\} \times \{0, \dots, 2^{n+1}-1\} \rightarrow T$ is a solution for \mathfrak{D} and c iff for all $x, y < 2^{n+1}$, the following holds (where \oplus_i denotes addition modulo i):

- if $\tau(x, y) = t$ and $\tau(x \oplus_{2^{n+1}} 1, y) = t'$, then $(t, t') \in H$
- if $\tau(x, y) = t$ and $\tau(x, y \oplus_{2^{n+1}} 1) = t'$, then $(t, t') \in V$
- $\tau(i, 0) = c_i$ for $i < n$.

For a proof of NEXPTIME-hardness of this version of the domino problem, see e.g. Corollary 4.15 in [9].

We show how to translate a given domino system \mathfrak{D} and initial condition $c = c_0 \dots c_{n-1}$ into an ABox $\mathcal{A}_{\mathfrak{D},c}$ and query $q_{\mathfrak{D},c}$ such that each (tree-shaped) model \mathcal{I} of $\mathcal{A}_{\mathfrak{D},c}$ that satisfies $\mathcal{I} \not\models q_{\mathfrak{D},c}$ encodes a solution to \mathfrak{D} and c , and conversely each solution to \mathfrak{D} and c gives rise to a (tree-shaped) model of $\mathcal{A}_{\mathfrak{D},c}$ with $\mathcal{I} \not\models q_{\mathfrak{D},c}$. The ABox $\mathcal{A}_{\mathfrak{D},c}$ contains only the assertion $C_{\mathfrak{D},c}(a)$, with $C_{\mathfrak{D},c}$ a conjunction $C_{\mathfrak{D},c}^1 \sqcap \dots \sqcap C_{\mathfrak{D},c}^7$, whose conjuncts we define in the following. For convenience, let $m = 2n + 2$. The purpose of the first conjunct $C_{\mathfrak{D},c}^1$ is to enforce a binary tree of depth m whose leaves are labelled with the numbers $0, \dots, 2^m - 1$ of a binary counter implemented by the concept names A_0, \dots, A_{m-1} . We use concept names L_0, \dots, L_m to distinguish the different levels of the tree. This is necessary because we work with reflexive and symmetric roles. In the following $\forall s^i . C$ denotes the i -fold nesting $\forall s . \dots \forall s . C$. In particular, $\forall s^0 . C$ is C .

$$C_{\mathfrak{D},c}^1 := L_0 \sqcap \prod_{i < m} \forall s^i . (L_i \rightarrow (\exists s . (L_{i+1} \sqcap A_i) \sqcap \exists s . (L_{i+1} \sqcap \neg A_i))) \sqcap$$

$$\prod_{i < m} \forall s^i . \prod_{j < i} ((L_i \sqcap A_j) \rightarrow \forall s . (L_{i+1} \rightarrow A_j) \sqcap$$

$$(L_i \sqcap \neg A_j) \rightarrow \forall s . (L_{i+1} \rightarrow \neg A_j))$$

From now on, leaves in this tree are called L_m -nodes. Intuitively, each L_m -node corresponds to a position in the $2^{n+1} \times 2^{n+1}$ -grid that we have to tile: the counter A_x realized by the concept names A_0, \dots, A_n binarily encodes the horizontal position, and the counter A_y realized by A_{n+1}, \dots, A_m encodes the vertical position. We now extend the tree with some additional nodes. Every L_m -node gets three successor nodes labelled with F , and each of these F -nodes has a successor node labelled G . To distinguish the three different G -nodes below each L_m -node, we additionally label them with the concept names G_1, G_2, G_3 .

$$C_{\mathfrak{D},c}^2 := \forall s^m. (L_m \rightarrow (\prod_{1 \leq i \leq 3} \exists s. (F \sqcap \exists s. (G \sqcap G_i))))$$

We want that each G_1 -node represents the grid position identified by its ancestor L_m -node, the sibling G_2 node represents the horizontal neighbor position in the grid, and the sibling G_3 -node represents the vertical neighbor.

$$C_{\mathfrak{D},c}^3 := \forall s^m. (L_m \rightarrow (\prod_{i \leq n} ((A_i \rightarrow \forall s^2. (G_1 \sqcup G_3 \rightarrow A_i)) \sqcap (\neg A_i \rightarrow \forall s^2. (G_1 \sqcup G_3 \rightarrow \neg A_i))) \sqcap \prod_{n < i < m} ((A_i \rightarrow \forall s^2. (G_1 \sqcup G_2 \rightarrow A_i)) \sqcap (\neg A_i \rightarrow \forall s^2. (G_1 \sqcup G_2 \rightarrow \neg A_i))) \sqcap E_2 \sqcap E_3))$$

where E_2 is an \mathcal{ALC} -concept ensuring that the A_x value at each G_2 -node is obtained from the A_x -value of its G -node ancestor by incrementing modulo 2^{n+1} ; similarly, E_3 expresses that the A_y value at each G_3 -node is obtained from the A_y -value of its G -node ancestor by incrementing modulo 2^{n+1} . It is not hard to work out the details of these concepts, see e.g. [11] for more details. The *grid representation* that we have enforced is shown in Figure 1. To represent tiles, we introduce a concept name D_i for each $i \in T$ and put

$$C_{\mathfrak{D},c}^4 := \forall s^{m+2}. (G \rightarrow (\bigsqcup_{i \in T} D_i \sqcap \prod_{i,j \in T, i \neq j} \neg(D_i \sqcap D_j)))$$

The initial condition is easily guaranteed by

$$C_{\mathfrak{D},c}^5 := \prod_{i < n} \forall s^{m+2}. ((\prod_{j \leq n, \text{bit}_j(i)=0} \neg A_j \sqcap \prod_{j \leq n, \text{bit}_j(i)=1} A_j \sqcap \prod_{n < j < m} \neg A_j) \rightarrow T_{c_i}),$$

where $\text{bit}_j(i)$ denotes the value of the j -th bit in the binary representation of i . To enforce the matching conditions, we proceed in two steps. First we ensure that they are satisfied locally, i.e., among the three G -nodes below each L_m -node:

$$C_{\mathfrak{D},c}^6 := \forall s^{m+2}. (L_m \rightarrow (\prod_{i \in T} (\exists s^2. (G_1 \sqcap D_i) \rightarrow \forall s^2. (G_2 \rightarrow \bigsqcup_{(i,j) \in H} D_j)) \sqcap \prod_{i \in T} (\exists s^2. (G_1 \sqcap D_i) \rightarrow \forall s^2. (G_3 \rightarrow \bigsqcup_{(i,j) \in V} D_j))))$$

Second, we enforce the following condition, which together with local satisfaction of the matching conditions ensures their global satisfaction:

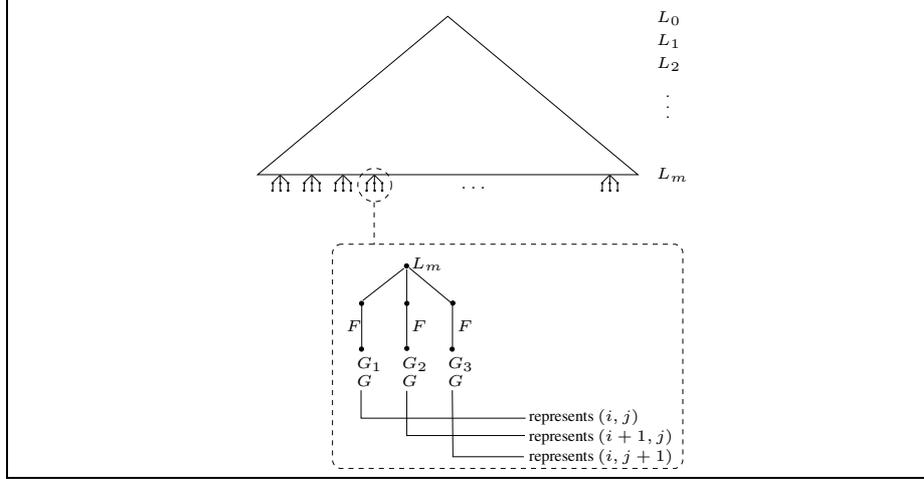


Fig. 1. The structure encoding the $2^{n+1} \times 2^{n+1}$ -grid.

(*) if the A_x and A_y -values of two G -nodes coincide, then their tile types coincide.

In (*), a G -node can be any of a G_1 -, G_2 -, or G_3 -node. To enforce (*), we use the query. Before we give details, let us finish the definition of the concept $C_{\mathcal{D},c}$. The last conjunct $C_{\mathcal{D},c}^T$ enforces two technical conditions that will be explained later: if d is an F -node and e its G -node successor, then

(T1) d and e are labelled dually regarding $A_i, \neg A_i$ for all $i < m$, i.e., d satisfies A_i iff e satisfies $\neg A_i$;

(T2) d and e are labelled dually regarding D_0, \dots, D_{k-1} , i.e., for all $j < k$, if d satisfies D_j , then e satisfies $D_0, \dots, D_{j-1}, \neg D_j, D_{j+1}, \dots, D_{k-1}$.

We use the following concept:

$$C_{\mathcal{D},c}^T := \forall s^{m+1}. (F \rightarrow (\prod_{i < m} (A_i \rightarrow \forall s. (G \rightarrow \neg A_i)) \sqcap (\neg A_i \rightarrow \forall s. (G \rightarrow A_i)) \sqcap \prod_{i \in T} \exists s. (G \sqcap D_i) \rightarrow (\neg D_i \sqcap \prod_{j < k, j \neq i} D_j)))$$

We now construct the query $q_{\mathcal{D},c}$ that does *not* match the grid representation iff (*) is satisfied. In other words, $q_{\mathcal{D},c}$ matches the grid representation if there are two G -nodes that agree on the value of the counters A_x and A_y , but are labelled with different tile types. Because of Lemma 1, we can concentrate on the grid representation as shown in Figure 1 while constructing $q_{\mathcal{D},c}$, and need not worry about models in which domain elements that are different in Figure 1 are identified.

The construction of $q_{\mathcal{D},c}$ is in several steps, starting with the query $q_{\mathcal{D},c}^i$ on the left-hand side of Figure 2, where $i \in \{0, \dots, m-1\}$. In the queries $q_{\mathcal{D},c}^i$, all the edges

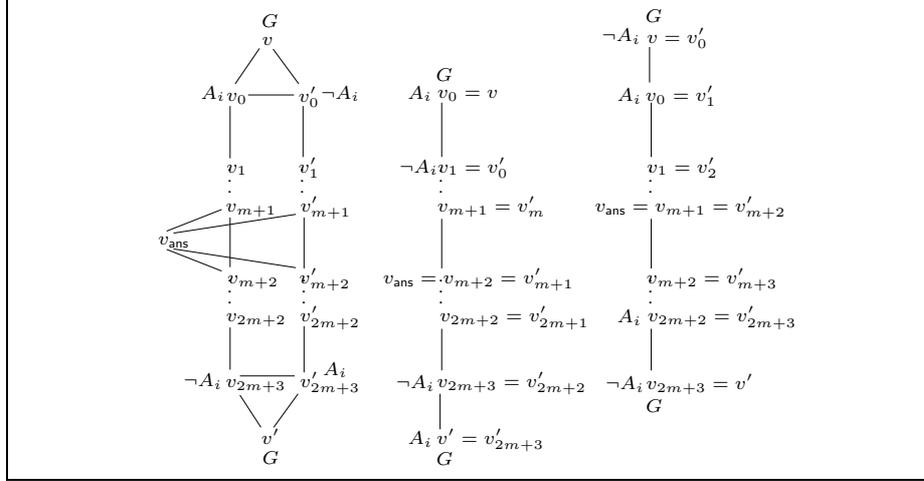


Fig. 2. The query $q_{D,a}^i$ (left) and two of its collapsings (middle and right).

represent the role s and v_{ans} is the only answer variable. The edges are undirected because we are working with symmetric roles. Formally,

$$\begin{aligned}
q_{D,c}^i := & \{ s(v_{i,0}, v_{i,1}), \dots, s(v_{i,2m+2}, v_{i,2m+3}), \\
& s(v'_{i,0}, v'_{i,1}), \dots, s(v'_{i,2m+2}, v'_{i,2m+3}), \\
& s(v_{i,0}, v'_{i,0}), s(v_{i,2m+3}, v'_{i,2m+3}), \\
& s(v, v_{i,0}), s(v, v'_{i,0}), \\
& s(v', v_{i,2m+3}), s(v', v'_{i,2m+3}), \\
& s(v_{ans}, v_{i,m+1}), s(v_{ans}, v_{i,m+2}), s(v_{ans}, v'_{i,m+1}), s(v_{ans}, v'_{i,m+2}), \\
& G(v), G(v'), A_i(v_{i,0}), \neg A_i(v'_{i,0}), \neg A_i(v_{i,2m+3}), A_i(v'_{i,2m+3}) \}
\end{aligned}$$

Observe that we dropped the index “ i ” to variables in Figure 2. Also observe that all the queries $q_{D,c}^i$, $i < m$, share the variables v , v' , and v_{ans} .

The purpose of the query $q_{D,a}^i$ is to relate any two G -nodes that agree on the value of the concept name A_i . To explain how this works, we need a few preliminaries. First, a *cycle* in a query is a sequence of distinct nodes v_0, \dots, v_{n-1} such that $n \geq 2$, and $s(v_i, v_{i+1}) \in q$ or $s(v_{i+1}, v_i) \in q$ for all $i < n$, where $v_n := v_0$. A query q' is a *collapsing* of a query q if q' is obtained from q by identifying variables. Each match of $q_{D,c}^i$ in our *tree-structured* grid representation gives rise to a collapsing of $q_{D,c}^i$ that does not comprise any cycles. To explain how $q_{D,c}^i$ works, it is helpful to analyze its cycle-free collapsings. We start with the two cycles v, v_0, v'_0 and v', v_{2m+3}, v'_{2m+3} . For eliminating each of these, we have two options:

- to remove the upper cycle, we can identify v with v_0 or v'_0 ;
- to remove the lower cycle, we can identify v' with v_{2m+3} or v'_{2m+3} .

Observe that if we identify v_0 and v'_0 (or v_{2m+3} and v'_{2m+3}) to collapse the cycle, there will be no matches of the query in any model.

Together, this gives four options for removing the two mentioned length-three cycles. However, two of these options are ruled out because the resulting collapsings have no match in the grid representation. The first such case is when we identify v with v_0 and v' with v_{2m+3} . Then v_0 and v_{2m+3} have to satisfy G . To continue our argument, we make a case distinction on the two options that we have for eliminating the cycle $\{v_{\text{ans}}, v_{m+1}, v_{m+2}\}$.

Case (1). If we identify v_{ans} and v_{m+1} , the path from the G -variable v_0 to v_{ans} is only of length $m+1$. In our grid representation, all paths from a G -node to an ABox individual (i.e., the root) are of length $m+2$, so there can be no match of this collapsing.

Case (2). If we identify v_{ans} and v_{m+2} , the path from v_{ans} to the G -variable v_{2m+3} is only of length $m+1$ and again there is no match.

We can argue analogously for the case where we identify v with v'_0 and v' with v'_{2m+3} . Therefore, the two remaining collapsings for eliminating the cycles $\{v, v_0, v'_0\}$ and $\{v', v_{2m+3}, v'_{2m+3}\}$ are the following:

- (a) identify v with v_0 and v' with v'_{2m+3} ;
- (b) identify v with v'_0 and v' with v_{2m+3} .

In the first case, we further have to identify v_{ans} with v_{m+2} and v'_{m+1} , for otherwise we can argue as above that there is no match. In the second case, we have to identify v_{ans} with v_{m+1} and v'_{m+2} . After this has been done, there is only one way to eliminate the cycle $v = v_0, \dots, v_{2m+3}, v' = v'_{2m+3}, \dots, v'_0$ such that the result is a chain of length $2m+4$ with the G -variables at both ends and the answer variable exactly in the middle (any other way to collapse means that there are no matches). The reflexive loops at the endpoints of the resulting chain and at v_{ans} can simply be dropped since we work with reflexive roles. The resulting cycle-free queries are shown in the middle and right part of Figure 2.

Note that the middle query has A_i at both ends of the chain, and the right one has $\neg A_i$ at the ends. According to our above argumentation, the original query $q_{\mathcal{D},c}^i$ has a match in the grid representation iff one of these two collapsings has a match. Thus, every match π of $q_{\mathcal{D},c}^i$ in the grid representation is such that $\pi(v)$ and $\pi(v')$ are (not necessarily distinct) instances of G that agree on the value of A_i . Informally, we say that $q_{\mathcal{D},c}^i$ connects G -nodes that have the same A_i -value.

At this point, a technical remark is in order. Observe that the two relevant collapsings of $q_{\mathcal{D},c}^i$ are such that the nodes next to the outer nodes are labelled dually w.r.t. A_i compared to the outer nodes. This is an artifact of query construction and cannot be avoided. It is the reason for introducing the F -nodes into our grid representation, and for ensuring that they satisfy Property (T1) from above.

Now set $q_{\text{cnt}} := \bigcup_{i < m} q_{\mathcal{D},c}^i$. It is easy to see that q_{cnt} connects G -nodes that have the same A_i -value, for all $i < m$. The query q_{cnt} is almost the desired query $q_{\mathcal{D},c}$. Recall that we want to enforce Condition (*) from above, and thus need to talk about tile types in the query. The query q_{tile} is given in the left-hand side of Figure 3 for the

because of paths lengths, we have to identify v_{ans} with $v_{i,m+1}$ and $v_{j,m+2}$. Next, we consider the cycle $v = w_{i,0}, \dots, w_{i,2m+3}, v' = w_{j,2m+3}, \dots, w_{j,0}$. As in the case of q_w^i , there is only one way to eliminate this cycle such that the result is a chain of length $2m+4$ with the G -variables at both ends and the answer variable exactly in the middle, and any other way to collapse means that there are no matches. It remains to eliminate the cycles $v = w_{i,0}, \dots, w_{i,2m+3}, v', w_{\ell,2m+3}, \dots, w_{\ell,0}$ with $\ell \neq j$. What is important here is that we have to identify $w_{i,1}$ with $w_{\ell,0}$ and $w_{i,2m+3}$ with $w_{\ell,2m+3}$. This is the case since the alternative (identifying $w_{i,0}$ with $w_{\ell,0}$ or $v' = 2_{j,2m+2}$ with $w_{\ell,2m+3}$) leads to a variable labelled with G , D_ℓ , and D_i (resp. D_j), and thus there is no match. Once these two identifications have been done, there is more than one way to identify the remaining nodes on the mentioned cycle, but the resulting query is always the same.

In summary, it is not hard to see that q_{tile} connects those G -nodes that are labelled by different tile types. Observe that we need property (T2) for this query to match at all.

Now, the desired query $q_{\mathcal{D},c}$ is simply the union of q_{cnt} and q_{tile} . From what was already said about q_{cnt} and q_{tile} , it is easily derived that $q_{\mathcal{D},c}$ does not match the grid representation iff Property (*) is satisfied. It is possible to show that there is a solution for \mathcal{D} and c iff $(\emptyset, \mathcal{A}_{\mathcal{D},c}) \not\models q_{\mathcal{D},c}$. We have thus proved that rooted query entailment in \mathcal{ALCI} is co-NEXPTIME-hard. A matching upper bound can be obtained by adapting the techniques in [6]. More details are given in the full version of this paper.

Theorem 1. *Rooted query entailment in \mathcal{ALCI} is co-NEXPTIME-complete. This holds even w.r.t. knowledge bases in which the TBox is empty and the ABox is a singleton.*

4 Boolean Query Entailment in \mathcal{ALCI} is 2-EXPTIME-complete

If we drop the requirement that queries are connected or that they have at least one answer variable, query entailment in \mathcal{ALCI} becomes 2-EXPTIME-complete. An upper bound can be taken e.g. from [6]. The lower bound can be proved by a reduction of the word problem of exponentially space bounded alternating Turing machines (ATMs) [4]. Because of space limitations, we can only give a very rough sketch of this result here. More details can be found in the extended version of this paper [10].

The main idea is to represent each configuration of an ATM by the leafs of a tree of depth n , very similar to the grid representation in Section 3. The trees representing configurations are then interconnected to a tree representing the computation. This is illustrated in Figure 4, where each of the T_i is a tree of depth n that is built using the role name s . The leafs of each T_i represent a configuration. The tree T_1 represents an existential configuration, and thus has only one successor configuration, which is represented by T_2 and connected via the same role name s also used inside the T_i trees. In contrast, the tree T_2 represents a universal configuration with two successor configurations T_2 and T_3 . The crucial point in the reduction is to relate the content of tape cells in one configuration to the content of the corresponding cells in the successor configurations. In principle, this is achieved using queries that are very similar to the query $q_{\mathcal{D},c}$ employed in the previous section. A few additional technical tricks are needed to achieve directedness (i.e., talking only about successor configurations, but not about predecessor configurations) since we work with symmetric roles.

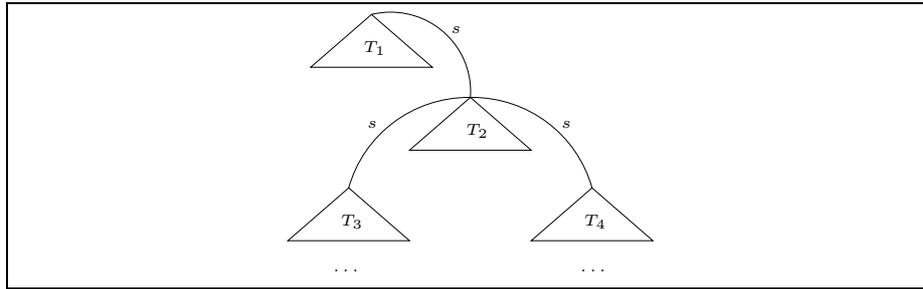


Fig. 4. Representing ATM computations.

Theorem 2. *Query entailment in \mathcal{ALCI} is 2-EXPTIME-complete. This holds even for queries without answer variables and w.r.t. knowledge bases in which the ABox is a singleton.*

5 Conclusion

We have shown that in the presence of inverse roles, conjunctive query answering is computationally more costly than instance checking. A Corresponding NEXPTIME upper bound for Theorem 1 and containment of conjunctive query entailment in EXPTIME for \mathcal{ALC} will be shown elsewhere. As (almost) remarked by a reviewer, the proof of Theorem 2 can easily be adapted to rooted query entailment if transitive roles and role hierarchies are present. Details on this will also be given elsewhere.

Acknowledgement We thanks the anonymous reviewers for valuable remarks on the submitted version of this paper.

References

1. F. Baader, D. L. McGuinness, D. Nardi, and P. Patel-Schneider. *The Description Logic Handbook: Theory, implementation and applications*. Cambridge University Press, 2003.
2. D. Calvanese, G. De Giacomo, and M. Lenzerini. On the decidability of query containment under constraints. In *Proceedings of the 17th ACM SIGACT SIGMOD SIGART Symposium on Principles of Database Systems (PODS'98)*, pages 149–158, 1998.
3. D. Calvanese, G. D. Giacomo, D. Lembo, M. Lenzerini, and R. Rosati. Data complexity of query answering in description logics. In P. Doherty, J. Mylopoulos, and C. Welty, editors, *Proceedings of the Tenth International Conference on Principles of Knowledge Representation and Reasoning (KR'06)*. AAAI Press, 2006.
4. A. K. Chandra, D. C. Kozen, and L. J. Stockmeyer. Alternation. *Journal of the ACM*, 28(1):114–133, 1981.
5. B. Glimm, I. Horrocks, and U. Sattler. Conjunctive query answering for description logics with transitive roles. In B. Parsia, U. Sattler, and D. Toman, editors, *Proceedings of the 2006 International Workshop on Description Logics (DL'06)*, volume 189 of *CEUR-WS*, 2006.

6. B. Glimm, C. Lutz, I. Horrocks, and U. Sattler. Answering conjunctive queries in the \mathcal{SHIQ} description logic. In M. Veloso, editor, *Proceedings of the Twentieth International Joint Conference on Artificial Intelligence (IJCAI'07)*, pages 299–404. AAAI Press, 2007.
7. I. Horrocks, U. Sattler, and S. Tobies. Reasoning with individuals for the description logic SHIQ. In D. MacAllester, editor, *Proceedings of the 17th International Conference on Automated Deduction (CADE-17)*, number 1831 in Lecture Notes in Computer Science, Germany, 2000. Springer Verlag.
8. U. Hustadt, B. Motik, and U. Sattler. Data complexity of reasoning in very expressive description logics. In L. P. Kaelbling and A. Saffiotti, editors, *Proceedings of the Nineteenth International Joint Conference on Artificial Intelligence (IJCAI'05)*, pages 466–471. Professional Book Center, 2005.
9. C. Lutz. *The Complexity of Reasoning with Concrete Domains*. PhD thesis, LuFG Theoretical Computer Science, RWTH Aachen, Germany, 2002.
10. C. Lutz. Inverse Roles Make Conjunctive Queries Hard. Available from <http://lat.inf.tu-dresden.de/~clu/papers/>
11. C. Lutz, C. Areces, I. Horrocks, and U. Sattler. Keys, nominals, and concrete domains. *Journal of Artificial Intelligence Research (JAIR)*, 23:667–726, 2005.
12. M. Ortiz, D. Calvanese, and T. Eiter. Data complexity of answering unions of conjunctive queries in \mathcal{SHIQ} . In B. Parsia, U. Sattler, and D. Toman, editors, *Proceedings of the 2006 International Workshop on Description Logics (DL'06)*, volume 189 of *CEUR-WS*, 2006.
13. A. Schaerf. On the complexity of the instance checking problem in concept languages with existential quantification. *Journal of Intelligent Information Systems*, 2:265–278, 1993.

A From $\mathcal{ALC}^{\text{rs}}$ to \mathcal{ALCI} without TBoxes

We show that rooted query entailment in $\mathcal{ALC}^{\text{rs}}$ w.r.t. the empty TBox can be polynomially reduced to rooted query entailment in \mathcal{ALCI} w.r.t. the empty TBox.

As already explained, the main idea behind the reduction is to replace each symmetric role r with the composition of r^- and r . Let \mathcal{A} be an $\mathcal{ALC}^{\text{rs}}$ ABox and q a conjunctive query. We assume w.l.o.g. that all concepts in \mathcal{A} are in negation normal form (NNF), i.e., that negation is applied only to concept names. Let $\text{Ind}(\mathcal{A})$ denote the set of all individual names occurring in \mathcal{A} , $\text{rol}(\mathcal{A})$ be the set of role names used in \mathcal{A} , and let $\text{rol}(q)$ be defined analogously. Fix a fresh concept name R . Intuitively, the purpose of R is to distinguish “real” domain elements from the auxiliary ones that serve as intermediate points in the composition of r^- and r . Also, define X as an abbreviation for $\bigwedge_{r \in \text{rol}(\mathcal{A}) \cup \text{rol}(q)} \exists r^- . \top$. We will enforce that X is satisfied by all relevant real individuals, thus achieving reflexivity.

We now present the details of the reduction. For each concept C in NNF, let $\delta(C)$ denote the result of replacing

- every subconcept $\exists r . C$ with $\exists r^- . \exists r . (C \sqcap R \sqcap X)$, and
- every subconcept $\forall r . C$ with $\forall r^- . \forall r . C$;

Now define an \mathcal{ALCI} ABox \mathcal{A}' and a query q' by manipulating \mathcal{A} and q as follows:

1. replace every concept assertion $C(a) \in \mathcal{A}$ with $\delta(C)(a)$;
2. for all $a \in \text{Ind}(\mathcal{A})$, add a concept assertion $R \sqcap X(a)$ to \mathcal{A} ;

3. replace every role assertion $r(a, b) \in \mathcal{A}$ with $r(c, a)$ and $r(c, b)$, where c is a fresh individual name;
4. for every variable v in q , add $R(v)$ to q ;
5. replace every role atom $r(v, v') \in q$ with $r(v^*, v)$ and $r(v^*, v')$, where v^* is a fresh variable.

The following lemma shows that our reduction is correct.

Lemma 2. $\mathcal{A} \not\models q$ iff $\mathcal{A}' \not\models q'$.

Proof. “ \Rightarrow ”. If $\mathcal{A} \not\models q$, then there is a model \mathcal{I} of \mathcal{A} such that $\mathcal{I} \not\models q$. Define a model \mathcal{I}' as follows:

- $\Delta^{\mathcal{I}'} = \Delta^{\mathcal{I}} \cup \{x_{d,r,e} \mid r \in \text{rol}(\mathcal{A}) \cup \text{rol}(q) \text{ and } (d, e) \in r^{\mathcal{I}}\}$;
- $r^{\mathcal{I}'} = \{(x_{d,r,e}, d), (x_{d,r,e}, e) \mid (d, e) \in r^{\mathcal{I}}\}$
- $A^{\mathcal{I}'} = A^{\mathcal{I}}$ for all concept names A except R ;
- $R^{\mathcal{I}'} = \Delta^{\mathcal{I}}$;
- $a^{\mathcal{I}'} = a^{\mathcal{I}}$ for all $a \in \text{Ind}(\mathcal{A})$;
- if c was introduced into \mathcal{A}' to split the assertion $r(a, b) \in \mathcal{A}$, set $c^{\mathcal{I}'} = x_{a^{\mathcal{I}}, r, b^{\mathcal{I}}}$.

It is readily checked that \mathcal{I}' is a model of \mathcal{A}' . In particular, $X^{\mathcal{I}'} = \Delta^{\mathcal{I}'}$ since roles are interpreted reflexively in \mathcal{I} . Furthermore, since $\mathcal{I} \not\models q$, we have $\mathcal{I}' \not\models q'$: suppose to the contrary that $\mathcal{I}' \models^{\pi} q'$ for some match π . Since q' contains the atom $R(v)$ for every variable $v \in \text{Var}(q)$, we have $\pi(v) \in \Delta^{\mathcal{I}'}$ for all $v \in \text{Var}(q)$. Let π' be the restriction of π to the variables in $\text{Var}(q)$. It is readily checked that $\mathcal{I} \models^{\pi'} q$, which is a contradiction.

“ \Leftarrow ”. If $\mathcal{A}' \not\models q'$, then there is a model \mathcal{I}' of \mathcal{A}' such that $\mathcal{I}' \not\models q'$. Define a model \mathcal{I} as follows:

- $\Delta^{\mathcal{I}} = (R \sqcap X)^{\mathcal{I}'}$;
- $r^{\mathcal{I}} = \{(d, e) \mid \exists f. (f, d) \in r^{\mathcal{I}'} \wedge (f, e) \in r^{\mathcal{I}'}\}$;
- $A^{\mathcal{I}} = A^{\mathcal{I}'} \cap \Delta^{\mathcal{I}'}$;
- $a^{\mathcal{I}} = a^{\mathcal{I}'}$ for all $a \in \text{Ind}(\mathcal{A})$.

Observe that $r^{\mathcal{I}}$ is reflexive (due to the choice of $\Delta^{\mathcal{I}}$ as a subset of $X^{\mathcal{I}'}$) and symmetric. Also observe that the interpretation of the individual names is well-defined: since \mathcal{A}' contains $R \sqcap X(a)$ for all $a \in \text{Ind}(\mathcal{A})$, $a^{\mathcal{I}'} \in \Delta^{\mathcal{I}'}$. Since $\mathcal{I}' \not\models q'$ and it is easily seen that $\mathcal{I} \models q$ would imply $\mathcal{I}' \models q'$, we have $\mathcal{I} \not\models q$. It remains to show that \mathcal{I} is a model of \mathcal{A} . This is a consequence of the following claim, which is easily proved by induction on the structure of C .

Claim. For all $d \in \Delta^{\mathcal{I}}$ and all $C \in \text{sub}(\mathcal{A})$, $d \in \delta(C)^{\mathcal{I}'}$ implies $d \in C^{\mathcal{I}}$.

We only do the two interesting cases.

- Let $C = \forall r. D$. Then $\delta(C) = \forall r^-. \forall r. \delta(D)$. Let $(d, e) \in r^{\mathcal{I}}$. We have to show that $e \in D^{\mathcal{I}}$. Since $(d, e) \in r^{\mathcal{I}}$, by definition of \mathcal{I} we have $(d, e) \in (r^-)^{\mathcal{I}'} \circ r^{\mathcal{I}'}$. Since $d \in \delta(C)^{\mathcal{I}'}$, we have $e \in D^{\mathcal{I}'}$ and it remains to apply the induction hypothesis.
- Let $C = \exists r. D$. Then $\delta(C) = \exists r^-. \exists r. (\delta(D) \sqcap R \sqcap X)$. Since $d \in \delta(C)^{\mathcal{I}'}$, there is an $e \in \Delta^{\mathcal{I}'}$ such that (i) $(d, e) \in (r^-)^{\mathcal{I}'} \circ r^{\mathcal{I}'}$ and (ii) $e \in (\delta(D) \sqcap R \sqcap X)^{\mathcal{I}'}$. By (ii), $d \in \Delta^{\mathcal{I}}$. By (i) and definition of \mathcal{I} , $(d, e) \in r^{\mathcal{I}}$. By (ii) and induction hypothesis, $d \in D^{\mathcal{I}}$ and we are done.

□