

# Approximation and Semantic Tree-width of Conjunctive Regular Path Queries\*

Diego Figueira ✉🏠

Univ. Bordeaux, CNRS, Bordeaux INP, LaBRI, UMR5800, F-33400 Talence, France

Rémi Morvan ✉🏠

Univ. Bordeaux, CNRS, Bordeaux INP, LaBRI, UMR5800, F-33400 Talence, France

---

## Abstract

We show that the problem of whether a query is equivalent to a query of tree-width  $k$  is decidable, for the class of Unions of Conjunctive Regular Path Queries with two-way navigation (UC2RPQs). A previous result by Barceló, Romero, and Vardi [3] has shown decidability for the case  $k = 1$ , and here we show that decidability in fact holds for any arbitrary  $k > 1$ . The algorithm is in 2EXPSpace, but for the restricted but practically relevant case where all regular expressions of the query are of the form  $a^*$  or  $(a_1 + \dots + a_n)$  we show that the complexity of the problem drops to  $\Pi_2^P$ .

We also investigate the related problem of approximating a UC2RPQ by queries of small tree-width. We exhibit an algorithm which, for any fixed number  $k$ , builds the maximal under-approximation of tree-width  $k$  of a UC2RPQ. The maximal under-approximation of tree-width  $k$  of a query  $q$  is a query  $q'$  of tree-width  $k$  which is contained in  $q$  in a maximal and unique way, that is, such that for every query  $q''$  of tree-width  $k$ , if  $q''$  is contained in  $q$  then  $q''$  is also contained in  $q'$ .

**2012 ACM Subject Classification** Information systems → Query languages for non-relational engines; Theory of computation → Parameterized complexity and exact algorithms; Theory of computation → Database query processing and optimization (theory)

**Keywords and phrases** graph databases, conjunctive regular path queries, semantic optimization, tree-width, containment, approximation

*This document contains internal hyperlinks: clicking on a keyword or mathematical notation jumps to its definition.<sup>1</sup>*

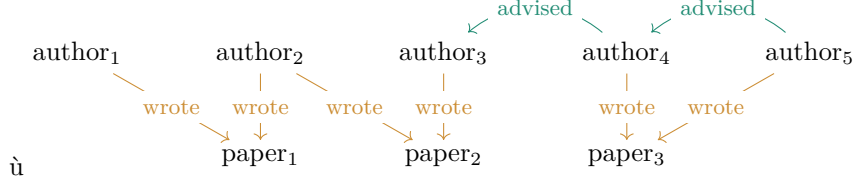
---

\* November 2022. Accepted to ICDT 2023. Licensed under CC BY 4.0.

<sup>1</sup> This result was achieved by using the [knowledge](#) package and its companion tool [knowledge-clustering](#).

## 1 Introduction

*Graph databases* are abstracted as edge-labeled directed graphs  $G = \langle V(G), E(G) \rangle$ , where nodes of  $V(G)$  represent entities and labeled edges  $E(G) \subseteq V(G) \times \mathbb{A} \times V(G)$  represent relations between these entities, with  $\mathbb{A}$  being a fixed finite alphabet. For instance, Figure 1 depicts a *graph database*, whose nodes are authors and papers, on the alphabet  $\mathbb{A} = \{\text{wrote}, \text{advised}\}$ . Edges  $x \xrightarrow{\text{wrote}} y$  indicate that the person  $x$  wrote the paper  $y$ , while edges  $x \xrightarrow{\text{advised}} y$  indicate that person  $x$  was the Ph.D. advisor of person  $y$ .



■ **Figure 1** A *graph database* with eight nodes and eight edges on a two-letter alphabet.

Being a subclass of relational databases, *graph databases* can be queried by the predominant query language of *conjunctive queries*, *a.k.a.* *CQs*, which consists of the closure under projection of conjunctions of atoms of the form  $x \xrightarrow{a} y$  for some letter  $a \in \mathbb{A}$ . For instance, the *conjunctive query*

$$\gamma_1(x, y) = x \xrightarrow{\text{wrote}} z \wedge y \xrightarrow{\text{wrote}} z$$

returns, when *evaluated* on the *graph database*  $G$  defined in Figure 1, all pairs of nodes  $(u, v)$  such that  $u$  is a co-author of  $v$ . Each variable not appearing in the left-hand side of the definition of a *conjunctive query* (in this example,  $z$ ) is existentially quantified. Note that every *CQ* can be seen as a *graph database*, where each atom is an edge; hence, we sometimes use *graph database* terminology for *CQs*.

The expressive power of *CQs* is somewhat limited, since *CQs* cannot express, for example, transitive closure. Since the ability to navigate paths is of importance in many *graph database* scenarios, most modern graph query languages support, as a central querying mechanism, *conjunctive regular path queries*, or *CRPQs* for short. *CRPQs* are defined analogously to *conjunctive queries*, except that their atoms are now of the form  $x \xrightarrow{L} y$  where  $L$  is an arbitrary regular language over the alphabet  $\mathbb{A}$ . For instance the *evaluation* of the *CRPQ*  $\gamma_2(x, y) = x \xrightarrow{\text{wrote}} z \wedge z' \xrightarrow{\text{wrote}} z \wedge y \xrightarrow{(\text{advised})^*} z'$  on  $G$  yields every pair of persons  $(u, v)$  such that  $u$  is a co-author of a “scientific descendant” of  $v$ .

Formally, a *CRPQ*  $\gamma$  is defined as a tuple  $\bar{z} = (z_1, \dots, z_n)$  of *output variables*<sup>2</sup> together with a conjunction of *atoms* of the form  $\bigwedge_{j=1}^m x_j \xrightarrow{L_j} y_j$ , where each  $L_j$  is a regular language. The set of all variables occurring in  $\gamma$ , namely<sup>3</sup>  $\{z_1, \dots, z_n\} \cup \{x_1, y_1, \dots, x_m, y_m\}$ , is denoted by  $\text{vars}(\gamma)$ . Given a *database*  $G$ , we say that  $(u_1, \dots, u_n)$  *satisfies*  $\gamma$  on  $G$  if there is a mapping  $f: \text{vars}(\gamma) \rightarrow V(G)$  such that  $u_i = f(z_i)$  for all  $1 \leq i \leq n$ , and for each  $1 \leq j \leq m$ , there exists a path from  $f(x_j)$  to  $f(y_j)$  in  $G$ , labelled by a word from  $L_j$  (if the path is empty, the labelled word is  $\varepsilon$ ). The *evaluation* of  $\gamma$  on  $G$  is then the set of all tuples that *satisfy*  $\gamma$ . For example,  $(\text{author}_2, \text{author}_5)$  *satisfies*  $\gamma_2$  on the *graph database*  $G$  of Figure 1 via the function that maps  $x$  to  $\text{author}_2$ ,  $y$  to  $\text{author}_5$ ,  $z$  to  $\text{paper}_2$ , and  $z'$  to  $\text{author}_3$ .

<sup>2</sup> For technical reasons (see the definition of *expansion*) we allow for a variable to appear multiple times.

<sup>3</sup> We neither assume disjointness nor inclusion between  $\{z_1, \dots, z_n\}$  and  $\{x_1, y_1, \dots, x_m, y_m\}$

The language of **CRPQ** can be extended to navigate edges in both directions. Consider the database  $G^\pm$  obtained from  $G$  by adding, for every edge  $x \xrightarrow{a} y$  in  $G$ , an extra edge  $y \xrightarrow{a^-} x$ . We obtain a graph database on the alphabet  $\mathbb{A}^\pm = \mathbb{A} \dot{\cup} \mathbb{A}^-$  where  $\mathbb{A}^- = \{a^- \mid a \in \mathbb{A}\}$ . We then define the syntax of a *CRPQ with two-way navigation*, or **C2RPQ**, as a **CRPQ** on the alphabet  $\mathbb{A}^\pm$ . Its *evaluation* is defined as the *evaluation* of the **CRPQ** on  $G^\pm$ . For instance, the *evaluation* of the **C2RPQ**  $\gamma_3(x, y) = x \xrightarrow{(\text{wrote} \cdot \text{wrote}^-)^*} y$  on the graph database of Figure 1 returns all pairs of individuals linked by a chain of co-authorship. It includes  $(\text{author}_1, \text{author}_3)$  or  $(\text{author}_1, \text{author}_1)$  but not  $(\text{author}_1, \text{author}_4)$ . If a query has no *output variables* we call it *Boolean*, and its *evaluation* can either be the set  $\{\emptyset\}$ , in which case we say that  $G$  *satisfies* the query, or the empty set  $\{\}$ . For example,  $G$  *satisfies*  $\gamma_4() = x \xrightarrow{\text{wrote}} y$  if, and only if, the database contains one author together with the paper they wrote. We denote the set of *atoms* of a **C2RPQ**  $\gamma$  by  $\text{Atoms}(\gamma)$ , and by  $\|\gamma\|$  we denote its number of atoms, *i.e.*,  $|\text{Atoms}(\gamma)|$ .

Finally, a *union of CQs (UCQs)* [resp. *union of CRPQs (UCRPQs)*, resp. *union of C2RPQs (UC2RPQs)*] is defined as a finite set of **CQs** [resp. **CRPQs**, resp. **C2RPQs**], whose tuples of *output variables* have all the same arity. The *evaluation* of a union is defined as the union of its *evaluations*, for instance:

$$\Gamma_5 = \gamma_5^1(x, y) \vee \gamma_5^2(x, y) \text{ where } \gamma_5^1(x, y) = x \xrightarrow{\text{wrote}} y \text{ and } \gamma_5^2(x, y) = x \xleftarrow{\text{advised}} z \wedge z \xrightarrow{\text{wrote}} y$$

*evaluates* to the set of pairs  $(x, y)$  such that  $y$  is a paper written by either  $x$  or their advisor. *Infinitary unions* are defined analogously, except that we allow for potentially infinite unions.

For a more detailed introduction to **CRPQs**, we refer the reader to [10]. For a more general introduction to different query languages for graph databases—including **CRPQs**—see [4], and for a more practical approach, see [1].

Given two **UC2RPQ**  $\Gamma$  and  $\Gamma'$ , we say that  $\Gamma$  is *contained* in  $\Gamma'$ , denoted by  $\Gamma \subseteq \Gamma'$  if for every graph database  $G$ , for every tuple  $\bar{u}$  of  $G$ , if  $\bar{u}$  *satisfies*  $\Gamma$  on  $G$ , then so does  $\Gamma'$ . The *containment problem* for **UC2RPQs** is the problem of, given two **UC2RPQs**  $\Gamma$  and  $\Gamma'$ , to decide if  $\Gamma \subseteq \Gamma'$ . When  $\Gamma$  is *contained* in  $\Gamma'$  and vice versa, we say that  $\Gamma$  and  $\Gamma'$  are *semantically equivalent*, denoted by  $\Gamma \equiv \Gamma'$ . The *evaluation problem* for **UC2RPQ** is the problem of, given a **C2RPQ**  $\gamma$ , a graph database  $G$  and a tuple  $\bar{u}$  of elements of  $G$ , whether  $\bar{u}$  *satisfies*  $\gamma$  on  $G$ .

## Queries of small tree-width

It is known that the *evaluation problem* for **UC2RPQ** is NP-complete, just as for *conjunctive queries* [8]. However, queries whose underlying structure looks like a tree—formally, queries of bounded *tree-width*—can be *evaluated* in polynomial time.

*Tree-width* is a measure of how much a graph differs from a tree, introduced by Arnborg and Proskurowski [2]. Formally, a *tree decomposition* of a **C2RPQ**  $\gamma$  is a pair  $(T, \mathbf{v})$  where  $T$  is a tree and  $\mathbf{v} : V(T) \rightarrow \wp(\text{vars}(\gamma))$  is a function that associates to each node of  $T$ , called *bag*, a set of variables of  $\gamma$ . When  $x \in \mathbf{v}(b)$  we shall say that the bag  $b \in V(T)$  *contains* the variable  $x$ . Further, it must satisfy the following three properties:

- each variable  $x$  of  $\gamma$  is *contained* in at least one *bag* of  $T$ ;
  - for each atom  $x \xrightarrow{L} y$  of  $\gamma$ , there is at least one *bag* of  $T$  that *contains* both  $x$  and  $y$ ; and
  - for each variable  $x$  of  $\gamma$ , the set of bags of  $T$  *containing*  $x$  is a connected subset of  $V(T)$ .
- The *width* of  $(T, \mathbf{v})$  is the maximum of  $|\mathbf{v}(b)| - 1$  when  $b$  ranges over  $V(T)$ . The *tree-width* of  $\gamma$  is the minimum of the *width* of all *tree decompositions* of  $\gamma$ . We denote by  $\mathcal{T}w_k$  the set of all **C2RPQ** of *tree-width* at most  $k$ . The *tree-width* of a **UC2RPQ** is simply the maximum

of the **tree-width** of its **C2RPQs**. An example of **tree decomposition** of width 2 is given in Figure 3 on Page 11. For a gentle introduction to **tree-width**, see [14, §3.6].

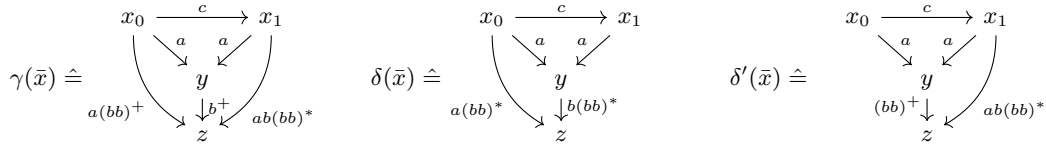
► **Proposition 1.1** (Folklore, see e.g. [15, Theorem IV.3]). *For each  $k \geq 1$ , the **evaluation problem** for **UC2RPQs** of **tree-width** at most  $k$  is in polynomial time.*

In practice, **graph databases** tend to be huge and often changing, while queries are in comparison very small. This motivates the following question, given some natural  $k \geq 1$ :

Given a **UC2RPQ**  $\Gamma$ , is it **equivalent** to a **UC2RPQ**  $\Gamma'$  of **tree-width** at most  $k$ ?  
That is, does it have **semantic tree-width** at most  $k$ ?

This problem is called the **semantic tree-width  $k$  problem**. Should it be decidable in a constructive way—that is, decidable, and if the answer is positive, we can compute a witnessing  $\Gamma'$  from  $\Gamma$ —, then one could, once and for all, compute  $\Gamma'$  from  $\Gamma$  and, whenever one wants to **evaluate**  $\Gamma$  on a database, **evaluate**  $\Gamma'$  instead.

► **Example 1.2.** Consider the following **CRPQs**, where  $\bar{x} = (x_0, x_1, y, z)$ :



The underlying graph of  $\gamma(\bar{x})$  being the directed 4-clique,  $\gamma(\bar{x})$  has **tree-width** 3. We claim that  $\gamma(\bar{x})$  is equivalent to the **UCRPQ**  $\delta(\bar{x}) \vee \delta'(\bar{x})$ , and hence has **semantic tree-width** 2.

Indeed, given a **graph database** satisfying  $\gamma(\bar{x})$  via some mapping  $\mu$ , it suffices to make a case disjunction on whether the number of  $b$ -labelled **atoms** in the path from  $\mu(y)$  to  $\mu(z)$  is even or odd. In the first case, the **atom**  $x_0 \xrightarrow{a(bb)^+} z$  becomes redundant since we can deduce the existence of such a path from the conjunction  $x_0 \xrightarrow{a} y \xrightarrow{(bb)^+} z$ , and hence the database satisfies  $\delta(\bar{x})$  via  $\mu$ . Symmetrically, in the second case, the **atom**  $x_1 \xrightarrow{b(bb)^*} z$  becomes redundant, and the database satisfies  $\delta'(\bar{x})$  via  $\mu$ . Thus,  $\gamma(\bar{x})$  is **contained**, and hence **equivalent** (the other **containment** being trivial), to the **UCRPQ**  $\delta(\bar{x}) \vee \delta'(\bar{x})$  of **tree-width** 2.

For **conjunctive queries**, the **semantic tree-width  $k$  problem** can be effectively decided quite easily—in fact, **CQs** enjoy the effective existence of unique minimal queries [8, Theorem 12] which happen to also minimize the tree-width. For **CRPQs** and **UC2RPQs**, the question is far more challenging, and it has only been solved for the case  $k = 1$  by Barceló, Romero, and Vardi [3, Theorem 6.1]. We solve the problem for every other  $k > 1$ , left open in [3, §7] [15, §VI-(3)]:

► **Theorem 1.3.** *For each  $k \geq 1$ , the **semantic tree-width  $k$  problem** is decidable. Moreover, it lies in 2EXPSpace and is EXPSpace-hard.*

Amusingly, our proof for  $k > 1$  cannot be stretched to capture the case  $k = 1$ : the two approaches seem to be intrinsically incompatible.

► **Remark 1.4.** To simplify proofs, we assume that all the regular languages are described via non-deterministic finite automata (NFA) instead of regular expressions, which does not affect any of our complexity bounds. However, for readability all our examples will be given in terms of regular expressions.

Moreover, we also show that for any class  $\mathcal{L}$  of regular languages over  $\mathbb{A}^\pm$  satisfying some mild hypothesis (“closure under sublanguages”), if  $\Gamma \in \text{UC2RPQ}(\mathcal{L})$  has semantic tree-width  $k > 1$ , then  $\Gamma$  is equivalent to a  $\text{UC2RPQ}(\mathcal{L})$  of tree-width at most  $k$ , where  $\text{UC2RPQ}(\mathcal{L})$  denotes the class of all  $\text{UC2RPQ}$ s whose atoms are all labelled by languages from  $\mathcal{L}$ . In other words, if a query can be defined with labels in  $\mathcal{L}$ , and if this query is equivalent to a query of small tree-width, then it is also equivalent to a query of small tree-width with labels in  $\mathcal{L}$ .

For a NFA  $\mathcal{A}$  and two states  $q, q'$  thereof, we denote by  $\mathcal{A}[q, q']$  the *sublanguage* of  $\mathcal{A}$  recognized when considering  $q$  as initial state and  $\{q'\}$  as set of final states. We say that  $\mathcal{L}$  is *closed under sublanguages* if (i) it contains every language of the form  $\{a\}$ , where  $a \in \mathbb{A}$  is any (positive) letter such that either  $a$  or  $a^-$  occur in a word of a language of  $\mathcal{L}$ , and (ii) for every language  $L \in \mathcal{L}$  there exists a NFA  $\mathcal{A}_L$  such that every sublanguage  $\mathcal{A}_L[q, q']$  distinct from  $\emptyset$  and  $\{\varepsilon\}$  belongs to  $\mathcal{L}$ .

To the best of our knowledge, all classes of regular expressions that have been considered in the realm of regular path queries (see, *e.g.*, [11, §1]) are *closed under sublanguages*. In particular, this is the case for the class  $\{\{a_1 + \dots + a_n\} \mid a_1, \dots, a_n \in \mathbb{A}\} \cup \{a^* \mid a \in \mathbb{A}\}$ , which will be our focus of study in Section 6. Moreover, even if some class  $\mathcal{L}$  is not *closed under sublanguages*, such as for example  $\{(aa)^*\}$ , then it is contained in a class *closed under sublanguages*— $\{a, a(aa)^*, (aa)^*\}$  in this example—, whose size<sup>4</sup> is polynomial in the size of the original class.

► **Theorem 1.5.** *Assume that  $\mathcal{L}$  is closed under sublanguages. For any query  $\Gamma \in \text{UC2RPQ}(\mathcal{L})$  and  $k > 1$ , the following are equivalent:*

1.  $\Gamma$  is equivalent to an infinitary union of conjunctive queries of tree-width at most  $k$ ;
2.  $\Gamma$  has semantic tree-width at most  $k$ ;
3.  $\Gamma$  is equivalent to a  $\text{UC2RPQ}(\mathcal{L})$  of tree-width at most  $k$ .

The implications (3)  $\Rightarrow$  (2)  $\Rightarrow$  (1) immediately follow from the definition of the semantic tree-width. On the other hand, the implications (1)  $\Rightarrow$  (2) and (2)  $\Rightarrow$  (3) are surprising, since they are both trivially false when  $k = 1$ , as we show next.

► **Remark 1.6.** (1)  $\not\Rightarrow$  (2) when  $k = 1$ : consider the  $\text{CRPQ}$   $\gamma(x, y) = x \xrightarrow{a^*} y \wedge y \xrightarrow{b} x$  of tree-width 1, and hence of semantic tree-width 1, and observe that it is not equivalent to any infinitary union of conjunctive queries of tree-width 1—this can be proven by considering, for example, the expansion  $x \xrightarrow{a} z \xrightarrow{a} y \wedge y \xrightarrow{b} x$  of  $\gamma(x, y)$  and applying Proposition 2.3.

(2)  $\not\Rightarrow$  (3) when  $k = 1$ : By [3, Proposition 6.4] the  $\text{CRPQ}$  of semantic tree-width 1  $\gamma(x) \doteq x \xleftarrow{a} z \xrightarrow{a} y \wedge x \xrightarrow{b} y \equiv x \xrightarrow{ba^*a} x$  is not equivalent to any  $\text{UCRPQ}$  of tree-width 1. Hence, the implication is false when  $\mathcal{L}$  is the class of regular languages over  $\mathbb{A}^\pm$  that do not use any letter of the form  $a^-$ . ◀

The proofs of both Theorems 1.3 and 1.5 rely on our key lemma (Lemma 3.7), which states essentially that every  $\text{UC2RPQ}$  has a computable “maximal under approximation” by a  $\text{UC2RPQ}$  of tree-width  $k$ . Formally, the key lemma has the following corollary:

► **Corollary 3.8.** *For each  $k > 1$  and for each class  $\mathcal{L}$  closed under sublanguages, for each query  $\Gamma \in \text{UC2RPQ}(\mathcal{L})$ , there exists  $\Gamma' \in \text{UC2RPQ}(\mathcal{L})$  of tree-width  $k$  such that  $\Gamma' \subseteq \Gamma$ , and for every  $\Delta \in \text{UC2RPQ}$ , if  $\Delta$  has tree-width  $k$  and  $\Delta \subseteq \Gamma$ , then  $\Delta \subseteq \Gamma'$ . Moreover,  $\Gamma'$  is computable from  $\Gamma$  in EXPSpace.*

<sup>4</sup> Defined as the sum of the number of states of the minimal automaton of the languages of  $\mathcal{L}$ .

The proof of our [key lemma](#) spans over Sections 3–5: in Section 3, we introduce necessary notions to formally state it, and deduce Theorems 1.3 and 1.5 from it; in Section 4 we introduce the central notion of [tagged tree decompositions](#) of [C2RPQs homomorphisms](#), and building on it, we finally describe the constructions used to prove the [key lemma](#) in Section 5.

Finally, in Section 6, given the high complexity of [semantic tree-width  \$k\$  problem](#), we focus on the case of [CRPQs](#) using some [simple regular expressions \(SRE\)](#), and show that the complexity of this problem is much lower:

► **Theorem 6.1.** *For  $k > 1$ , the [semantic tree-width  \$k\$  problem](#) for [UCRPQ\(SRE\)](#) is in  $\Pi_2^P$ .*

A discussion on differences with Barceló, Romero and Vardi’s contributions and open questions are left for Section 7.

## 2 Homomorphisms, refinements, and expansions

Before attacking the statement of our [key lemma](#) in Section 3, we first give a few elementary definitions on [C2RPQs](#) in this section. A [homomorphism](#)  $h$  from a [C2RPQ](#)  $\gamma(x_1, \dots, x_m)$  to a [C2RPQ](#)  $\gamma'(y_1, \dots, y_m)$  is a mapping from  $\text{vars}(\gamma)$  to  $\text{vars}(\gamma')$  such that  $h(x) \xrightarrow{L} h(y)$  is an atom of  $\gamma'$  for every atom  $x \xrightarrow{L} y$  of  $\gamma$ , and further  $h(x_i) = y_i$  for every  $i$ . Such a [homomorphism](#)  $h$  is [strong onto](#) if for every atom  $x' \xrightarrow{L} y'$  of  $\gamma'$  there is an atom  $x \xrightarrow{L} y$  of  $\gamma$  such that  $h(x) = x'$  and  $h(y) = y'$ . We write  $\gamma \xrightarrow{\text{hom}} \gamma'$  if there is a [homomorphism](#) from  $\gamma$  to  $\gamma'$ , and  $\gamma \xrightarrow{\text{hom}} \gamma'$  if there is a [strong onto homomorphism](#). It is easy to see that if  $\gamma \xrightarrow{\text{hom}} \gamma'$  then  $\gamma' \subseteq \gamma$ , and in the case where  $\gamma, \gamma'$  are [CQs](#) this is an “if and only if” [8].

**Refinements** An [atom  \$m\$ -refinement](#) of a [C2RPQ atom](#)  $\gamma(x, y) = x \xrightarrow{L} y$  where  $L$  is given by the NFA  $\mathcal{A}_L$  is any [C2RPQ](#) of the form

$$\rho(x, y) = x \xrightarrow{L_1} t_1 \xrightarrow{L_2} \dots \xrightarrow{L_{n-1}} t_{n-1} \xrightarrow{L_n} y \quad (1)$$

where  $1 \leq n \leq m$ ,  $t_1, \dots, t_{n-1}$  are fresh (existentially quantified) variables, and  $L_1, \dots, L_n$  are such that there exists a sequence  $(q_0, \dots, q_n)$  of states of  $\mathcal{A}_L$  such that  $q_0$  is initial,  $q_n$  is final, and for each  $i$ ,  $L_i$  is either of the form (i)  $\mathcal{A}_L[q_i, q_{i+1}]$ , or (ii)  $\{a\}$  if the letter  $a \in \mathbb{A}$  belongs to  $\mathcal{A}_L[q_i, q_{i+1}]$ , or (iii)  $\{a^{-1}\}$  if  $a^{-1} \in \mathbb{A}^-$  belongs to  $\mathcal{A}_L[q_i, q_{i+1}]$ . Additionally, if  $\varepsilon \in L$ , the equality atom “ $x = y$ ” is also an [atom  \$m\$ -refinement](#) (see Appendix B.1 for more details on these). Thus, an [atom  \$m\$ -refinement](#) can be either of the form (1) or “ $x = y$ ”. By convention,  $t \xrightarrow{a} t'$  is a shorthand for  $t' \xrightarrow{a} t$ . As a consequence, the underlying graph of an [atom  \$m\$ -refinement](#) of the form (1) is not necessarily a directed path. By definition, note that  $L_1 \cdots L_n \subseteq L$  and hence  $\rho \subseteq \gamma$  for any [atom  \$m\$ -refinement](#)  $\rho$  of  $\gamma$ . An [atom refinement](#) is an [atom  \$m\$ -refinement](#) for some  $m$ .

► **Definition 2.1.** *Given an [atom refinement](#)  $\rho = x \xrightarrow{L_1} t_1 \xrightarrow{L_2} \dots \xrightarrow{L_{n-1}} t_{n-1} \xrightarrow{L_n} y$  of  $\gamma = x \xrightarrow{L} y$  as in (1), define a [contraction](#) of  $\rho$  between  $t_i$  and  $t_j$ , where  $0 \leq i, j \leq n$  and  $j > i + 1$ , is any [C2RPQ](#) of the form:*

$$\rho' = x \xrightarrow{L_1} t_1 \xrightarrow{L_2} \dots \xrightarrow{L_i} t_i \xrightarrow{K} t_j \xrightarrow{L_{j+1}} \dots \xrightarrow{L_{n-1}} t_{n-1} \xrightarrow{L_n} y$$

such that  $K = \mathcal{A}_L[q_i, q_j]$ . Then every [contraction](#)  $\rho'$  of  $\rho$  is a [refinement](#) of  $\gamma$ , and  $\rho \subseteq \rho' \subseteq \gamma$ . Informally, we will abuse the notation and write  $[L_i \cdots L_j]$  to denote the language  $K$ —even if this language does not only depend on  $L_i \cdots L_j$ .



► **Example 2.2.** Let  $\gamma(x, y) = x \xrightarrow{(aa^-)^*} y$  be a **C2RPQ atom**, where  $(aa^-)^*$  is implicitly represented by its minimal automaton. Then  $\rho(x, y)$  is a **refinement** of **refinement length** seven of  $\gamma(x, y)$  and  $\rho'(x, y)$  is a **contraction** of  $\rho(x, y)$ , where:

$$\begin{aligned}\rho(x, y) &= x \xrightarrow{a} t_1 \xrightarrow{(a^-a)^*} t_2 \xrightarrow{(a^-a)^*} t_3 \xleftarrow{a} t_4 \xrightarrow{(aa^-)^*} t_5 \xrightarrow{(aa^-)^*a} t_6 \xleftarrow{a} y, \\ \rho'(x, y) &= x \xrightarrow{a} t_1 \xrightarrow{(a^-a)^*a^-} t_4 \xrightarrow{(aa^-)^*} y.\end{aligned}$$

On the other hand,  $\rho''(x, y) = x \xrightarrow{a} t_1 \xleftarrow{a} y$  is not a **contraction** of  $\rho(x, y)$ .

An **m-refinement** of a **C2RPQ**  $\gamma(\bar{x}) = \bigwedge_i x_i \xrightarrow{L_i} y_i$  is any query resulting from: 1) replacing every **atom** by one of its **m-refinements**, and 2) should some **m-refinements** be **equality atoms**, collapsing equal variables and getting rid of equalities, in a rather standard way—more details are given in Appendix B.1. A **refinement** is an **m-refinement** for some  $m$ . Note that any **atom m-refinements** is, by definition, also an **atom m'-refinements** when  $m < m'$ : as a consequence, in the **refinement** of a **C2RPQ** the **atom refinements** need not have the same length. For instance, both  $\rho(x, x) = x \xrightarrow{c} x$  and  $\rho'(x, y) = x \xrightarrow{a} t_1 \xrightarrow{a} y \xleftarrow{c} y$  are **refinements** of  $\gamma(x, y) = x \xrightarrow{a^*} y \xleftarrow{c} x$ . For a given **C2RPQ**  $\gamma$ , let  $\text{Ref}^{\leq m}(\gamma)$  be the set of all **m-refinements** of  $\gamma$ , and  $\text{Ref}(\gamma)$  be the set of all its **refinements**. Given a **refinement**  $\rho(\bar{x})$  of  $\gamma(\bar{x})$ , its **refinement length** is the least integer  $m$  such that  $\rho(\bar{x}) \in \text{Ref}^{\leq m}(\gamma)$ . Note that if the automaton representing a language  $L$  has more than one final state, for instance the minimal automaton for  $L = a^+ + b^+$ , then  $x \xrightarrow{L} y$  is not a **refinement** of itself. However, it will always be **equivalent** to a union of refinements: in this example,  $x \xrightarrow{a^+ + b^+} y$  is **equivalent** to the union of  $x \xrightarrow{a^+} y$  and  $x \xrightarrow{b^+} y$ , which are both **refinements** of the original **C2RPQ**.

**Expansions** Remember that a **C2RPQ** whose languages are  $\{a\}$  or  $\{a^-\}$  for  $a \in \mathbb{A}$  is in effect a **CQ**. The **expansions** of a **C2RPQ**  $\gamma$  is the set  $\text{Exp}(\gamma)$  of all **CQs** which are **refinements** of  $\gamma$ . In other words, an **expansion** of  $\gamma$  is any **CQ** obtained from  $\gamma$  by replacing each **atom**  $x \xrightarrow{L} y$  by a path  $x \xrightarrow{w} y$  for some word  $w \in L$ . For instance,  $\xi(x, y) = x \xrightarrow{a} t_1 \xleftarrow{a} t_2 \xrightarrow{a} t_3 \xleftarrow{a} y$  is an **expansion** of  $\rho(x, y) = x \xrightarrow{(aa^-)^*} y$ .

Any **C2RPQ** is equivalent to the infinitary union of its **expansions**. In light of this, the semantics for **UC2RPQ** can be rephrased as follows. Given a **UC2RPQ**  $\Gamma$  and a graph database  $G$ , the **evaluation** of  $\Gamma$  over  $G$ , denoted by  $\Gamma(G)$ , is the set of tuples  $\bar{v}$  of nodes for which there is  $\xi \in \text{Exp}(\Gamma)$  such that  $\xi \xrightarrow{\text{hom}} (G, \bar{v})$ . Similarly, **containment** of **UC2RPQs** can also be characterized in terms of expansions:

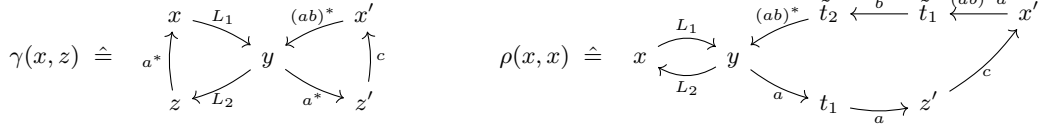
► **Proposition 2.3** (Folklore, see e.g. [12, Proposition 3.2] or [7, Theorem 2]). *Let  $\Gamma_1$  and  $\Gamma_2$  be **UC2RPQs**. Then the following are equivalent*

- $\Gamma_1 \subseteq \Gamma_2$ ;
- for every  $\xi_1 \in \text{Exp}(\Gamma_1)$ ,  $\xi_1 \subseteq \Gamma_2$ ;
- for every  $\xi_1 \in \text{Exp}(\Gamma_1)$  there is  $\xi_2 \in \text{Exp}(\Gamma_2)$  such that  $\xi_2 \xrightarrow{\text{hom}} \xi_1$ .

Hence, every **expansion** of  $\gamma$  is also a **refinement** of  $\gamma$ . Moreover, if  $\rho$  is a **refinement** of  $\gamma$ , then  $\rho \subseteq \gamma$ , and  $\gamma$  is **semantically equivalent** to the infinitary union of all its **refinements**, and to the infinitary union of all its **expansions**.

Our approach to proving Theorems 1.3 and 1.5 and the **key lemma** heavily rely on **refinements**. One crucial property that these objects satisfy is that they preserve **tree-width**  $k$ , unless  $k = 1$ . This is the main reason why our approach cannot capture the case  $k = 1$ , solved by Barceló, Romero and Vardi [3].

▷ **Fact 2.4.** Let  $k > 1$  and let  $\gamma$  be a **C2RPQ** of **tree-width** at most  $k$ . Then any **refinement** of  $\gamma$  has **tree-width** at most  $k$ .



■ **Figure 2** On the left-hand side, a CRPQ  $\gamma$  of tree-width 2. On the right-hand side, a refinement  $\rho$  of  $\gamma$ , whose refinement length is three, which is also of tree-width 2.

This fact is illustrated on an example in Figure 2.

Proof sketch. The underlying graph of a refinement of  $\gamma$  is obtained from the underlying graph of  $\gamma$  by either contracting some edges (when dealing with equality atoms), or by replacing a single edge by a path of edges (where the non-extremal nodes are new nodes). Both operations preserve the property “having tree-width at most  $k$ ” when  $k > 1$ . Details are given in Appendix B.2.  $\triangleleft$

For  $k = 1$ , the property fails: for instance the CRPQ  $\gamma(x) = x \xrightarrow{a^*} x$  has tree-width at most 1 (in fact it has tree-width 0), but its refinement  $\rho(x) = x \xrightarrow{a^*} t_1 \xrightarrow{a^*} t_2 \xrightarrow{a^*} x$  has tree-width two.

### 3 Maximal under-approximations of bounded tree-width

In this section, we state our key technical result, Lemma 3.7. Essentially, we follow the same structure as Theorem 1.5: given a C2RPQ  $\gamma$  and an integer  $k > 1$ , we start by consider its maximal under-approximation by infinitary unions of conjunctive queries of tree-width  $k$  (Definition 3.1), and then show that this query can in fact be expressed as a UC2RPQ of tree-width  $k$  whose atoms are obtained by taking sublanguages from  $\gamma$  (Lemma 3.7).

For the definitions of this section, let us fix any class  $\mathcal{C}$  of C2RPQ queries.

► **Definition 3.1** (Maximal under-approximation). *Let  $\gamma$  be a C2RPQ. The maximal under-approximation of  $\gamma$  by infinitary unions of  $\mathcal{C}$ -queries, is  $\text{App}_{\mathcal{C}}(\gamma) \triangleq \{\alpha \in \mathcal{C} \mid \alpha \subseteq \gamma\}$ .*

► **Remark 3.2.** Observe that  $\text{App}_{\mathcal{C}}(\gamma)$  is an infinitary union of  $\mathcal{C}$ -queries, that  $\text{App}_{\mathcal{C}}(\gamma) \subseteq \gamma$ , and that for every infinitary union of  $\mathcal{C}$ -queries  $\Delta$ , if  $\Delta \subseteq \gamma$ , then  $\Delta \subseteq \text{App}_{\mathcal{C}}(\gamma)$  (i.e., it is the unique maximal under-approximation). Similarly, the maximal under-approximation of a UC2RPQ is simply the union of the maximal under-approximations of the C2RPQs thereof.

Unfortunately, the fact that a query  $\alpha$  is part of this union,  $\alpha \in \text{App}_{\mathcal{C}}(\gamma)$ , does not yield any useful information on the shape of  $\alpha$ —we merely know that  $\alpha \subseteq \gamma$ . The rest of this subsection is dedicated to introducing another infinitary union of  $\mathcal{C}$ -queries, namely  $\text{App}_{\mathcal{C}}^*(\gamma) \subseteq \text{App}_{\mathcal{C}}(\gamma)$ , in which queries  $\alpha \in \text{App}_{\mathcal{C}}^*(\gamma)$  come together with a witness—a homomorphism—of their containment in  $\gamma$ .

► **Definition 3.3.** *The maximal under-approximation of  $\gamma$  by infinitary unions of homomorphically-smaller  $\mathcal{C}$ -queries is*

$$\text{App}_{\mathcal{C}}^*(\gamma) \triangleq \{\alpha \in \mathcal{C} \mid \exists \rho \in \text{Ref}(\gamma), \exists f: \rho \xrightarrow{\text{hom}} \alpha\}. \quad (2)$$

For instance, let  $\mathcal{C}$  be the class of all CRPQs and let  $\gamma(x, y) \triangleq x \xrightarrow{a^*} y \xrightarrow{a} z \xrightarrow{a^*} x \wedge y \xrightarrow{a^*} x$  be the CRPQ which asks for all pairs of nodes that belong to the same non-empty cycle of  $a$ ’s<sup>5</sup>. Then  $\gamma'(x, y) = x \xrightarrow{a^*} y \xrightarrow{a} z \xrightarrow{a^*} x$  is an element of  $\text{App}_{\mathcal{C}}^*(\gamma)$  since

<sup>5</sup> There exists CRPQs with fewer variables that expresses the same property, but this is irrelevant here.



there is a **strong onto homomorphism** from the refinement

$$\rho(x, y) = (x \xrightarrow{a^*} y \xrightarrow{a} z \xrightarrow{a^*} x \wedge y \xrightarrow{a} z' \xrightarrow{a^*} x) \in \text{Ref}(\gamma(x, y))$$

to  $\gamma'(x, y)$ , mapping  $z$  and  $z'$  to  $z$ .

Clearly,  $\text{App}_{\mathcal{C}}^*(\gamma)$ —whose queries are informally called *approximations*—is included, and thus **contained**, in  $\text{App}_{\mathcal{C}}(\gamma)$ , since  $\rho \subseteq \gamma$  and  $\alpha \subseteq \rho$  in (2). In fact, under some assumptions on  $\mathcal{C}$ , the converse **containment** also holds.

► **Observation 3.4.** *If  $\mathcal{C}$  is closed under expansions, then for any C2RPQ  $\gamma$ , we have  $\text{App}_{\mathcal{C}}(\gamma) \equiv \text{App}_{\mathcal{C}}^*(\gamma)$ .*

**Proof.** This follows immediately from Proposition 2.3—note that in the definition of  $\text{App}_{\mathcal{C}}^*(\gamma)$  we work with **strong onto homomorphisms**, but we can always restrict **homomorphisms** to their image to make them **strong onto**, without changing the expressiveness of the query.  $\square$

Observe then, by Fact 2.4, that the class  $\mathcal{T}_{w_k}$  of all C2RPQs of **tree-width** at most  $k$  is closed under **refinements** and hence under **expansions**, provided that  $k$  is greater or equal to 2. As an immediate consequence, we have:

► **Corollary 3.5.** *For  $k \geq 2$ , for all C2RPQ  $\gamma$ ,  $\text{App}_{\mathcal{T}_{w_k}}(\gamma) \equiv \text{App}_{\mathcal{T}_{w_k}}^*(\gamma)$ .*

► **Example 3.6** (Counter-example for  $k = 1$ ). Consider the followings queries:

$$\gamma(x, y) \triangleq \begin{array}{ccc} & z & \xrightarrow{(ba)^*} z' \\ a \uparrow & & \downarrow b \\ x & \xrightarrow{(ab)^+} & y \\ & \xleftarrow{c} & \end{array} \quad \text{and} \quad \delta(x, y) \triangleq x \xrightleftharpoons[c]{(ab)^+} y.$$

We claim that  $\text{App}_{\mathcal{T}_{w_1}}(\gamma) \not\equiv \text{App}_{\mathcal{T}_{w_1}}^*(\gamma)$  since  $\delta(x, y) \in \text{App}_{\mathcal{T}_{w_1}}(\gamma)$  but  $\delta(x, y) \notin \text{App}_{\mathcal{T}_{w_1}}^*(\gamma)$ . Details are given in Appendix C.1

By definition,  $\text{App}_{\mathcal{T}_{w_k}}^*(\gamma)$  is an infinitary union of C2RPQs. We show that, in fact,  $\text{App}_{\mathcal{T}_{w_k}}^*(\gamma)$  is always equivalent to a *finite* union of C2RPQs. This is done by bounding the length of the refinements occurring in the definition of  $\text{App}_{\mathcal{T}_{w_k}}^*(\gamma)$ . For a natural  $m$ , let  $\text{App}_{\mathcal{C}}^{*, \leq m}(\Gamma) \triangleq \{\alpha \in \mathcal{C} \mid \exists \rho \in \text{Ref}^{\leq m}(\Gamma), \exists f: \rho \xrightarrow{\text{hom}} \alpha\}$ . Our main technical lemma is then the following:

► **Lemma 3.7 (Key lemma).** *For  $k \geq 1$  and C2RPQ  $\gamma$ , we have  $\text{App}_{\mathcal{T}_{w_k}}(\gamma) \equiv \text{App}_{\mathcal{T}_{w_k}}^{*, \leq \ell}(\gamma)$ , where  $\ell = \Theta(\|\gamma\|^2 \cdot (k+1)^{\|\gamma\|+1})$ .*

► **Corollary 3.8.** *For each  $k > 1$  and for each class  $\mathcal{L}$  closed under sublanguages, for each query  $\Gamma \in \text{UC2RPQ}(\mathcal{L})$ , there exists  $\Gamma' \in \text{UC2RPQ}(\mathcal{L})$  of **tree-width**  $k$  such that  $\Gamma' \subseteq \Gamma$ , and for every  $\Delta \in \text{UC2RPQ}$ , if  $\Delta$  has **tree-width**  $k$  and  $\Delta \subseteq \Gamma$ , then  $\Delta \subseteq \Gamma'$ . Moreover,  $\Gamma'$  is computable from  $\Gamma$  in EXPSpace.*

Using Lemma 3.7 as a black box—which will be proven in Section 5—, we can now give a proof of Theorems 1.3 and 1.5. The upper bound of Theorem 1.3 follows directly from Corollary 3.8: to test whether a query  $\Gamma$  is of **semantic tree-width**  $k$ , it suffices to test the **containment**  $\Gamma \subseteq \Gamma'$ , where  $\Gamma'$  is the **maximal under-approximation** given by Corollary 3.8. The containment problem being in EXPSpace [12, 7], we obtain:

See the proof of  
Lemma 3.9 in  
page 20.

► **Lemma 3.9.** *For  $k \geq 1$ , the semantic tree-width  $k$  problem for UC2RPQ is in 2EXPSpace.*

An EXPSpace lower bound follows by a straightforward adaptation from the EXPSpace lower bound for the case  $k = 1$  [3, Proposition 6.2].

See the proof of  
Lemma 3.10 in  
page 21.

► **Lemma 3.10.** *The semantic tree-width  $k$  problem is EXPSpace-hard, even if restricted to Boolean CRPQs.*

We can now rely on the equivalence  $\text{App}_{\mathcal{T}_{w_k}}(\gamma) \equiv \text{App}_{\mathcal{T}_{w_k}}^{*, \leq \ell}(\gamma)$  to prove Theorem 1.5.

**Proof of Theorem 1.5.** The implications  $(3) \Rightarrow (2) \Rightarrow (1)$  are straightforward: they follow directly from Fact 2.4. For  $(1) \Rightarrow (3)$ , note that (1) implies that  $\Gamma \equiv \text{App}_{\mathcal{T}_{w_k}}(\Gamma)$ , and by Lemma 3.7,  $\text{App}_{\mathcal{T}_{w_k}}(\Gamma) \equiv \Delta \triangleq \bigvee_{\gamma \in \Gamma} \text{App}_{\mathcal{T}_{w_k}}^{*, \leq \ell}(\gamma)$ , so  $\Gamma$  is equivalent to the latter. Since queries of  $\Delta$  are obtained as homomorphic images of refinements of  $\Gamma$ , all of which are labelled by sublanguages of  $\mathcal{L}$ , and since  $\mathcal{L}$  is closed under sublanguages, it follows that hence  $\Gamma$  is equivalent to a UC2RPQ( $\mathcal{L}$ ) of tree-width  $k$ . ◀

We are left with the proof of Lemma 3.7, which will take the next two sections. Since  $\text{App}_{\mathcal{T}_{w_k}}(\gamma) \equiv \text{App}_{\mathcal{T}_{w_k}}^{*, \leq \ell}(\gamma)$  by Corollary 3.5 and since  $\text{App}_{\mathcal{T}_{w_k}}^{*, \leq \ell}(\gamma)$  is a subset of  $\text{App}_{\mathcal{T}_{w_k}}(\gamma)$ , we only need to show that  $\text{App}_{\mathcal{T}_{w_k}}(\gamma) \subseteq \text{App}_{\mathcal{T}_{w_k}}^{*, \leq \ell}(\gamma)$ . Formally, this means that for all  $\alpha \in \mathcal{T}_{w_k}$ , if there exists  $\rho \in \text{Ref}(\gamma)$  and  $f: \rho \xrightarrow{\text{hom}} \alpha$ , then there exists  $\alpha' \in \mathcal{T}_{w_k}$  such that  $\alpha \subseteq \alpha'$  and there exists  $\rho' \in \text{Ref}^{\leq \ell}(\gamma)$  and  $f': \rho' \xrightarrow{\text{hom}} \alpha'$ . We prove this by “massaging” the homomorphism  $f: \rho \xrightarrow{\text{hom}} \alpha$ , by looking where each atom refinement of  $\gamma$  is sent on  $\alpha$  relative to a “well-behaved” tree decomposition of  $\alpha$  of width  $k$ . Next, we introduce the notion of tagged tree decomposition to have a precise handle on this information.

## 4 Intermezzo: tagged tree decompositions

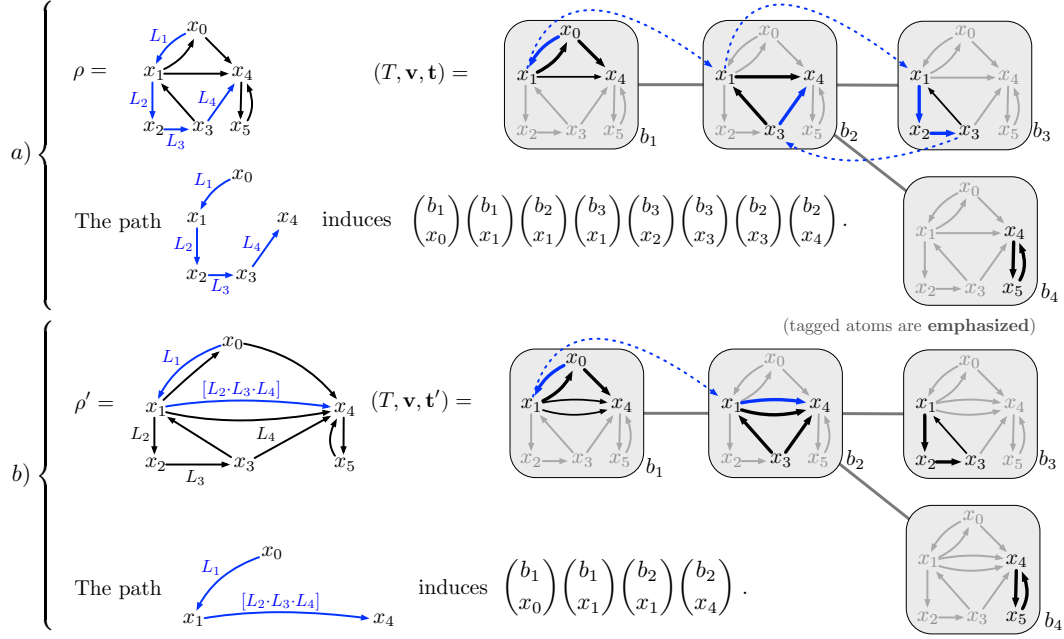
► **Definition 4.1.** *Let  $f: \rho \rightarrow \alpha$  be a homomorphism between two C2RPQs. A tagged tree decomposition of  $f$  is a triple  $(T, \mathbf{v}, \mathbf{t})$  where  $(T, \mathbf{v})$  is a tree decomposition of  $\alpha$ , and  $\mathbf{t}$  is a mapping  $\mathbf{t}: \text{Atoms}(\gamma) \rightarrow V(T)$ , called tagging, such that  $\mathbf{v}(\mathbf{t}(e))$  contains both  $f(x)$  and  $f(y)$  for each atom  $e = x \xrightarrow{\lambda} y \in \text{Atoms}(\gamma)$ .*

In other words,  $\mathbf{t}$  gives, for each atom of  $\gamma$ , a witnessing bag that contains it, in the sense that it contains the image by  $f$  of the atom source and target. By definition, given a tree decomposition  $(T, \mathbf{v})$  of  $\alpha$  and a homomorphism  $f: \rho \rightarrow \alpha$ , there is always one way (usually many) of extending  $(T, \mathbf{v})$  into a tagged tree decomposition of  $f$ .

► **Fact 4.2.** Let  $(T, \mathbf{v}, \mathbf{t})$  be a tagged tree decomposition of some homomorphism  $f: \rho \rightarrow \alpha$ . Let  $T'$  be the smallest connected subset of  $T$  induced by the image of  $\mathbf{t}$ . Then  $(T', \mathbf{v}|_{T'}, \mathbf{t})$  is still a tagged tree decomposition of  $f$ , whose width is at most the width of  $(T, \mathbf{v}, \mathbf{t})$ .

We extend the notion of tagging to paths: a formal definition (Definition D.1) can be found in Appendix D, and the notion is illustrated in Figure 3. In the context of a (nice) tagged tree decomposition  $(T, \mathbf{v}, \mathbf{t})$  of  $f: \rho \xrightarrow{\text{hom}} \alpha$ , given a path  $\pi$  of  $\rho$ , say  $x_0 \xrightarrow{\lambda_1} x_1 \xrightarrow{\lambda_2} \dots \xrightarrow{\lambda_n} x_n$  (in blue in Figure 3a), the path induced by  $\pi$ , denoted  $\mathbf{t}[\pi]$ , is informally defined as the following “path” in  $T \times \alpha$ , seen as a sequence of pairs from  $V(T) \times \text{vars}(\alpha)$ :

- it starts with the bag  $\mathbf{t}(x_0 \xrightarrow{\lambda_1} x_1)$  of  $T$  and the variable  $f(x_0)$  of  $\alpha$ ; and it continues with  $(\mathbf{t}(x_0 \xrightarrow{\lambda_1} x_1), f(x_1))$  (corresponding to the first blue edge in  $b_1$  of Figure 3a);
- it then follows the shortest path in  $T$  (unique, since it is a tree) that goes to the bag  $\mathbf{t}(x_1 \xrightarrow{\lambda_2} x_2)$ , while staying in  $f(x_1)$  in  $\alpha$  (in Figure 3a, it follows the blue path:  $(b_2, x_2), (b_3, x_2)$ ) and it traverses the atom  $x_1 \xrightarrow{\lambda_2} x_2$  (i.e., we go to  $(b_3, x_3)$ );



■ **Figure 3** a) Path induced by some path in  $\rho$ , in a tagged tree decomposition of  $f: \rho \xrightarrow{\text{hom}} \alpha$  in the case where  $\alpha = \rho$  and  $f$  is the identity homomorphism  $\text{id}_\rho: \rho \xrightarrow{\text{hom}} \rho$ .  
b) Suppose that the path in  $\rho$  is the image of an atom refinement of  $\gamma$ . Then the cycle in the induced path can be avoided by adding an atom  $x_1 \xrightarrow{[L_2 \cdot L_3 \cdot L_4]} x_4$  to  $\rho$ , obtaining the path of  $\rho'$ .

- it continues in the same way for all other atoms of the path, ending up with the bag  $\mathbf{t}(x_{n-1} \xrightarrow{\lambda_n} x_n)$  and the variable  $f(x_n)$  of  $\alpha$ .

By construction, note that the constructed sequence  $(b_i, z_i)_i$  is such that  $z_i \in \mathbf{v}(b_i)$ . Moreover, given a bag  $b$  of  $T$  and a variable  $z$  of  $\alpha$ , we say that  $\mathbf{t}[\pi]$  *leaves*  $b$  at  $z$  when  $z = z_i$  and  $b = b_i$  for some  $i$ , and either  $b_{i+1}$  is undefined, or distinct from  $b$ . For example, in Figure 3a,  $\tau[\pi]$  *leaves*  $b_1$  at  $x_1$  and *leaves*  $b_2$  at  $x_1$  and at  $x_4$ . We say that an induced path is *cyclic* if it contains three positions  $i < j < i'$  such that  $i$  and  $i'$  contain the same bag but  $j$  contains a different bag (the one in Figure 3a is *cyclic* because it passes twice by the bag  $b_2$ ).

For technical reasons—the proof of Claim 5.2—, we need to use the classical notion of *nice tree decomposition* (see e.g. [13, Definition 13.1.4, page 149]), which is a tree decomposition  $(T, \mathbf{v})$  such that any bag  $b$  that is not a leaf either: 1) has exactly two children  $b_1, b_2$  such that  $\mathbf{v}(b_1) = \mathbf{v}(b) = \mathbf{v}(b_2)$ , or 2) has exactly one child  $b'$ , and  $\mathbf{v}(b')$  is obtained from  $\mathbf{v}(b)$  by either adding a single vertex, or by removing a single vertex. A C2RPQ has *tree-width*  $k$  if and only if it has a *nice tree decomposition* of *width* at most  $k$  [13, Lemma 13.1.2, page 149]. In the context of a *nice tree decomposition* of *width*  $k$ , a *full bag* is any bag of size  $k + 1$ . Note that any *non-branching path*—i.e. a path whose non-extremal bags have degree 2—in a *nice tree decomposition* with  $n$  bags must have at least  $\lfloor n/2 \rfloor$  bags which are not full. Finally, we define a *nice tagged tree decomposition* of  $f: \rho \rightarrow \alpha$  to be a tagged tree decomposition of  $f$  that is also a *nice tree decomposition* of  $\alpha$ .

## 5 Proof of the key lemma

We can now start to describe the constructions used to prove Lemma 3.7. We call a *trio* to a triple  $(\alpha, \rho, f)$  such that  $\alpha \in \mathcal{F}\omega_k$ ,  $\rho \in \text{Ref}(\gamma)$  and  $f$  is a *strong onto homomorphism* from  $\gamma$

to  $\alpha$ . For clarity, we will denote such a **trio** by simply “ $f: \rho \xrightarrow{\text{hom}} \alpha$ ”. Using this terminology, in order to prove Lemma 3.7, we must show that for every **trio**  $f: \rho \xrightarrow{\text{hom}} \alpha$ , there exists another **trio**  $f': \rho' \xrightarrow{\text{hom}} \alpha'$  such that  $\alpha \subseteq \alpha'$  and  $\rho' \in \text{Ref}^{\leq \ell}(\gamma)$ . Our first construction, which will ultimately allow us to bound the size of **atom refinements**, shows that we can assume *w.l.o.g.* that they induce acyclic paths in a nice tagged tree decomposition of  $f$ .

See the proof of  
Claim 5.1 in  
page 24.

▷ **Claim 5.1.** For any **trio**  $f: \rho \xrightarrow{\text{hom}} \alpha$ , there exists a **trio**  $f': \rho' \xrightarrow{\text{hom}} \alpha'$  and a nice tagged tree decomposition  $(T', \mathbf{v}', \mathbf{t}')$  of width at most  $k$  of  $f'$  such that  $\alpha \subseteq \alpha'$ ,  $\|\rho'\| \leq \|\rho\|$  and every **atom refinement** of  $\rho'$  induces an acyclic path in the tree  $T'$ , in which case we say that  $(T', \mathbf{v}', \mathbf{t}')$  is *locally acyclic*.

The construction behind Claim 5.1 is quite simple, and is described and proven in Section 5, and illustrated in Figure 3.

Informal proof of Claim 5.1. Start with a **trio**  $f: \rho \xrightarrow{\text{hom}} \alpha$ , and let  $(T, \mathbf{v}, \mathbf{t})$  be a nice tagged tree decomposition of  $f$ . Consider an **atom refinement**  $\pi \hat{=} z_0 \xrightarrow{L_1} z_1 \xrightarrow{L_2} \dots \xrightarrow{L_n} z_n$  in  $\rho$  of some atom  $x \xrightarrow{L} y$  (with  $z_0 \hat{=} x$  and  $z_n \hat{=} y$ ), and assume that it induces a cyclic path in  $T$ , as in Figure 3a. It means that some variables  $z_i$  and  $z_j$  are mapped by  $f$  to the same bag of  $T$ , somewhere along the path induced by  $\pi$ . It suffices then to **contract**  $\rho$  by replacing the atoms  $z_i \xrightarrow{L_{i+1}} \dots \xrightarrow{L_j} z_j$  by a single atom  $z_i \xrightarrow{[L_{i+1} \dots L_j]} z_j$  (in Figure 3,  $z_i = x_1$  and  $z_j = x_4$ ). We thus obtain a new **refinement**  $\rho'$  of  $\gamma$ . Then define  $\alpha'$  be simply adding an atom  $f(z_i) \xrightarrow{L_{i+1} \dots L_j} f(z_j)$ , see Figure 3b. The definitions of  $f'$  and  $(T', \mathbf{v}', \mathbf{t}')$  are then straightforward—potentially,  $\alpha'$  should be restricted to the image of  $f': \rho' \rightarrow \alpha'$  so that  $f'$  is still **strong onto**. Crucially,  $\alpha \subseteq \alpha'$ , and  $\alpha'$  still has **tree-width** at most  $k$  since we picked  $f(z_i)$  and  $f(z_j)$  so that they belonged to the same bag of  $T$ . ◁

Ultimately, Claim 5.1 will allow us to give a bound on the number of leaves of a nice tagged tree decomposition of a **trio**. The following claim—which is significantly more technical than the foregoing—will give us a bound on the height of a **decomposition**.

▷ **Claim 5.2.** Let  $f: \rho \xrightarrow{\text{hom}} \alpha$  be a **trio** and  $(T, \mathbf{v}, \mathbf{t})$  be a locally acyclic nice tagged tree decomposition of width at most  $k$  of  $f$ . Then there is a **trio**  $f': \rho' \xrightarrow{\text{hom}} \alpha'$  and a nice tagged tree decomposition  $(T', \mathbf{v}', \mathbf{t}')$  of width at most  $k$  of  $f'$  such that:

- $\alpha \subseteq \alpha'$ ,
- $(T', \mathbf{v}', \mathbf{t}')$  is locally acyclic w.r.t.  $f'$ , and
- the size of the longest **non-branching path** in  $T$  is at most  $\Theta(\|\gamma\| \cdot (k+1)^{\|\gamma\|+1})$ .

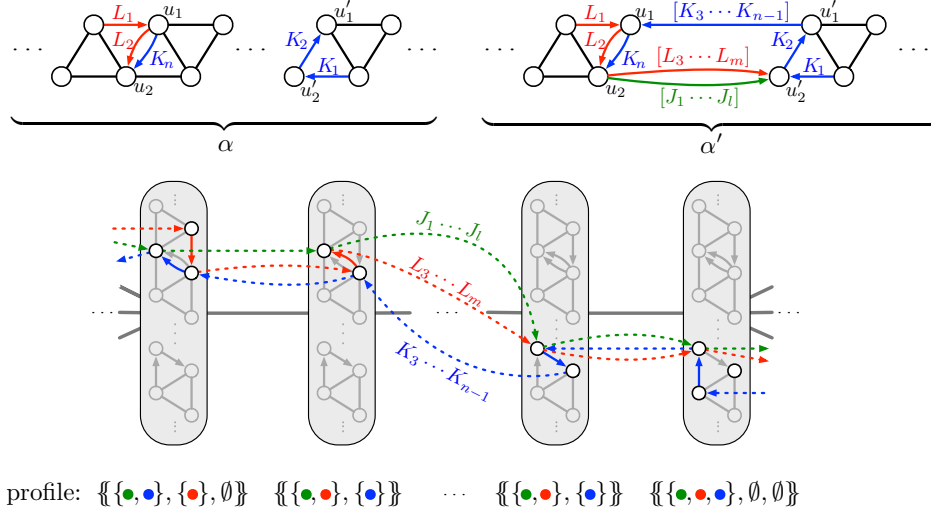
To prove Claim 5.2, we will try to find, in a long **non-branching path**, some kind of shortcut. The piece of information that is relevant to finding this shortcut is what we call the **profile** of a **bag**.

► **Definition 5.3.** Given a **trio**  $f: \rho \xrightarrow{\text{hom}} \alpha$  and a nice tagged tree decomposition  $(T, \mathbf{v}, \mathbf{t})$  of  $f$ , for each **bag**  $b$  of  $T$ , we say that:

- $b$  is “**atomic**” if there is at least one atom  $e \in \mathbf{t}^{-1}(b)$  and at least one variable  $x$  of  $e$  such that  $x \in \text{vars}(\gamma)$ , i.e., the atom  $e$  is not in the ‘middle’ part of an **atom refinement**;
- otherwise, when  $b$  is **non-atomic**, we assign to each variable  $z \in \mathbf{v}(b) \subseteq V(\alpha)$  a **type**

$$\text{type}_z \hat{=} \{x \xrightarrow{L} y \text{ atom of } \gamma \mid \text{the path induced by the atom refinement of } x \xrightarrow{L} y \text{ in } \rho \text{ leaves } b \text{ at } z\};$$

then the **profile** of  $b$  is the multiset of the types of  $z$  where  $z$  ranges over  $\mathbf{v}(b)$ .



■ **Figure 4** A long non-branching path in the tree decomposition of width 2 of an approximation  $\alpha$ . There are two non-full bags in the path with the same profile, and thus the query  $\alpha$  can be simplified to  $\alpha'$  by applying contractions to the atom refinements involved.

The rest of the proof consists in two parts: first, we show that if two non-atomic bags  $b$  and  $b'$  occurring in some non-branching path of  $T$  have the same profile, then we can essentially replace the path between  $b$  and  $b'$  by a path of constant length (Subclaim 5.4): while this construction is quite elementary, it motivates the intricate definition of the profile of a bag; then, we show that in every non-branching path, if it is sufficiently long, then we can find  $b$  and  $b'$  satisfying the aforementioned property: this part simply relies on a basic combinatorial argument (Fact E.1).

▷ **Subclaim 5.4.** Suppose there are two bags  $b$  and  $b'$  such that: (i) they contain at most  $k$  nodes (*i.e.*, not full bags), (ii) they have the same profile, (iii) there is a non-branching path in  $T$  between these bags, and (iv) no bags of the path between  $b$  and  $b'$  (both included) are atomic. Then, there exists a trio  $f': \rho' \xrightarrow{\text{hom}} \alpha'$  and a nice tagged tree decomposition of  $f'$  of width at most  $k$  that can be obtained by replacing the non-branching path between  $b$  and  $b'$  in the nice tagged tree-decomposition of  $f: \rho \xrightarrow{\text{hom}} \alpha$  by another non-branching path of at most  $2k + 1$  bags, such that  $\alpha \subseteq \alpha'$ .

See the proof of Subclaim 5.4 in page 25.

The basic idea behind Subclaim 5.4 is that we the definition of profile was carefully design so that we could contract every refinements between  $b$  and  $b'$ , while preserving every desirable properties on the trio. The construction is illustrated in Figure 4, and both an informal proof and a formal proof can be found in Appendix E.2. We can now describe key steps in the proof of Claim 5.2. A formal proof can also be found in Appendix E.2.

**Proof sketch of Claim 5.2.** We claim that, starting from a trio  $f: \rho \xrightarrow{\text{hom}} \alpha$  and a locally acyclic nice tagged tree decomposition of  $f$ , if we can find a long non-branching path, then an elementary argument (Fact E.1) yields the existence of two bags  $b$  and  $b'$  on this path, satisfying the assumptions of Subclaim 5.4, and sufficiently far apart that the construction described in Subclaim 5.4 strictly shortens the path between  $b$  and  $b'$ . Overall, the iterative application of this construction, which preserves both the niceness and the local acyclicity of the tagged tree decomposition, and only produces bigger approximations (in the sense of

containment), yields a trio  $f': \rho' \xrightarrow{\text{hom}} \alpha'$  with a locally acyclic nice tagged tree decomposition of width at most  $k$ , whose non-branching paths are all “small”, and such that  $\alpha \subseteq \alpha'$ .  $\triangleleft$

Finally, our main lemma follows from Claims 5.1 and 5.2.

**Proof of Lemma 3.7.** In order to show  $\text{App}_{\mathcal{T}_{w_k}}(\gamma) \subseteq \text{App}_{\mathcal{T}_{w_k}}^{*, \leq \ell}(\gamma)$ —the other containment being trivial—, pick a trio  $f: \rho \rightarrow \alpha$ . Applying Claim 5.1 and then Claim 5.2 yields the existence of a trio  $f': \rho' \rightarrow \alpha'$  together with a nice tagged tree decomposition  $(T', \mathbf{v}', \mathbf{t}')$  of  $f'$  such that  $\alpha \subseteq \alpha'$  and  $(T', \mathbf{v}', \mathbf{t}')$  is locally acyclic, and any non-branching path in  $T'$  has length at most  $\Theta(\|\gamma\| \cdot (k+1)^{\|\gamma\|})$ .

Moreover, we can assume *w.l.o.g.*, by applying Fact 4.2, that every leaf of  $T'$  is tagged by at least one atom of  $\rho'$ . The local acyclicity of  $T'$  implies that if  $b$  be a leaf of  $T'$ , and  $\pi \triangleq x \xrightarrow{L_1} t_1 \xrightarrow{L_2} \dots \xrightarrow{L_{n-1}} t_{n-1} \xrightarrow{L_n} y$  is an atom refinement in  $\rho'$  of some atom  $x \xrightarrow{L} y$  of  $\gamma$ , then if  $b$  is tagged by one atom of  $\pi$  this atom must either be  $z_0 \xrightarrow{L_1} z_1$  or  $z_{n-1} \xrightarrow{L_n} z_n$  by local acyclicity—see e.g. Figure 3 for a visual proof. The number of such atoms in  $\rho'$  being bounded by  $2\|\gamma\|$ , we conclude that  $T'$  has at most  $2\|\gamma\|$  leaves.

Then, observe that a tree with at most  $p$  leaves and whose non-branching paths have length at most  $q$  is of height at most<sup>6</sup>  $p \cdot q - 1$ . We conclude that the height of  $T'$  is  $\Theta(\|\gamma\|^2 \cdot (k+1)^{\|\gamma\|})$ . Using again the local acyclicity of  $T'$ , observe that the refinement length of  $\rho'$  is at most twice the height of  $T'$ , and hence  $\rho' \in \text{Ref}^{\leq \ell}(\gamma)$  where  $\ell = \Theta(\|\gamma\|^2 \cdot (k+1)^{\|\gamma\|})$ . In other words,  $\alpha' \in \text{App}_{\mathcal{T}_{w_k}}^{*, \leq \ell}(\gamma)$ . Hence, we have shown that for all  $\alpha \in \text{App}_{\mathcal{T}_{w_k}}(\gamma)$ , there exists  $\alpha' \in \text{App}_{\mathcal{T}_{w_k}}^{*, \leq \ell}(\gamma)$  such that  $\alpha \subseteq \alpha'$ .  $\blacktriangleleft$

## 6 Queries over simple regular expressions

A *simple regular expression*, or *SRE*, is a regular expression the form  $a^*$  for some letter  $a \in \mathbb{A}$  or of the form  $a_1 + \dots + a_m$  for some  $a_1, \dots, a_m \in \mathbb{A}$ .

Let  $\text{UCRPQ}(\text{SRE})$  be the set of all UCRPQ whose languages are expressed via SRE expressions. Observe that  $\text{UCRPQ}(\text{SRE})$  is semantically closed under concatenation, that is, concatenations of SRE expressions can be also expressed in the language. For example,  $\gamma(x, y) = x \xrightarrow{a^* \cdot (a+b) \cdot b^*} y$  is equivalent to  $\gamma'(x, y) = x \xrightarrow{a^*} z \wedge z \xrightarrow{a+b} z' \wedge z' \xrightarrow{b^*} y$ . One interest of  $\text{UCRPQ}(\text{SRE})$  comes from the fact that it is used widely in practice, as recent studies on SPARQL query logs on Wikidata, DBpedia and other sources show that this kind of regular expressions cover a majority of the queries investigated, *e.g.*, 75% of the “property paths” (C2RPQ atoms) of the corpus of 1.5M queries of [6, Table 15]. An additional interest comes from the fact that the containment problem for UCRPQ(SRE) is much better behaved than for general UCRPQs, since it is in  $\Pi_2^P$  [11, Corollary 5.2], that is, just one level up the polynomial hierarchy compared to the CQ containment problem, which is in NP [8], and in sharp contrast with the costly EXPSPACE-complete CRPQ containment problem [7, 12].

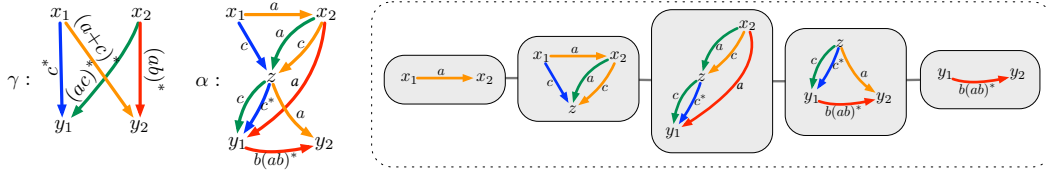
We devote this section to showing the following result.

► **Theorem 6.1.** *For  $k > 1$ , the semantic tree-width  $k$  problem for UCRPQ(SRE) is in  $\Pi_2^P$ .*

Observe that simple regular expressions are closed under sublanguages. Hence, in the light of Theorem 1.5, the maximal under-approximation of a UCRPQ(SRE) query by infinitary unions of CQs of tree-width  $k$  is always equivalent to a UCRPQ(SRE) query of tree-width  $k$ .

<sup>6</sup> The length of a path being its number of nodes, and with the convention that the height of a single node is zero.





■ **Figure 5** The query  $\mathbf{P}_l(\{x_1, x_2\}, \{y_1, y_2\}, \gamma)$  (where  $l = 2$ ) on the left contains the approximation  $\alpha$ , witnessed by the path decomposition on the right.

We will explain how the construction of the maximal under-approximation of the previous section can be exploited to improve the complexity from 2EXPSpace down to  $\Pi_2^P$ .

## 6.1 Summary queries

We will first show that the maximal under-approximation of tree-width  $k$  of a UC2RPQ can be expressed as a union of polynomial sized “summary” queries. Each summary query represents a union of exponentially-bounded C2RPQs sharing some *common structure*. These are normal UC2RPQ queries extended with some special kind of atoms, called “path- $l$  approximations”. Intuitively, a path- $l$  approximation is a maximal under-approximation of tree-width  $l$  of queries of the form  $\bigwedge_i x_i \xrightarrow{L_i} y_i$  such that  $x_i \neq y_j$  for all  $i, j$ . Path- $l$  approximations may require an exponential size when represented as UCRPQs. Formally, a *path- $l$  approximation* is a query of the form “ $\mathbf{P}_l(X, Y, \gamma)$ ” where: (i)  $X, Y$ , are two disjoint sets of variables of size at most  $l$ , (ii)  $\gamma(\bar{z})$  is a conjunction of atoms  $\bigwedge_{1 \leq i \leq n} A_i(x_i, y_i)$  where  $\bar{z}$  contains all variables of  $X \cup Y$ , (iii) each  $A_i$  is a C2RPQ atom of the form  $x_i \xrightarrow{L} y_i$  or  $y_i \xrightarrow{L} x_i$  such that  $x_i$  is in  $X$  and  $y_i$  is in  $Y$ . We give its semantics in terms of infinitary unions of CQs. A query like the one before is defined to be equivalent to the (infinitary) union of all queries  $\alpha(\bar{z}) \in \text{App}_{\mathcal{F}_{w_k}}(\gamma)$  that admit a path decomposition of width  $l$  having the bag  $X$  at the root and  $Y$  at the leaf—where a *path decomposition* is defined to be any tree decomposition whose underlying tree is a path. See Figure 5 for an example.

We now simply define a *k-summary query* as a C2RPQ extended with path- $l$  approximation atoms for any  $l \leq k$ , with the expected semantics. The important property of *summary queries* is that they are exponentially more succinct than the UC2RPQ counterpart for expressing maximal under-approximations, as the next lemma shows.

► **Proposition 6.2.** *For every class  $\mathcal{L}$  closed under sublanguages, and for every C2RPQ( $\mathcal{L}$ )  $\gamma$ ,  $\text{App}_{\mathcal{F}_{w_k}}(\gamma)$  can be expressed as a union of polynomial-sized  $k$ -summary queries having only C2RPQ( $\mathcal{L}$ ) atoms. Further, one can test in NP if a summary query is part of this union. We call  $\text{App}_{\mathcal{F}_{w_k}}^{\text{zip}}(\gamma)$  to any such a union of summary queries.*

See the proof of Proposition 6.2 in page 27.

**Informal proof.** As corollary of the proof of Lemma 3.7, we can assume to have  $\text{App}_{\mathcal{F}_{w_k}}(\gamma)$  expressed as a union of C2RPQ( $\mathcal{L}$ ) with a nice tree decomposition of width  $k$  with a linear number of leaves, and hence it suffices to replace non-branching paths with path- $l$  approximations. Concretely, for any such C2RPQ  $\alpha$  having a witnessing tree decomposition with a long non-branching path, the tree must contain a sub-path whose every bag is non-atomic, and such that it starts and ends in bags of size at most  $k$  (by the niceness property). Such a non-atomic non-branching path can be “compressed” by replacing the subquery corresponding to the path with a corresponding path- $l$  approximation query. The resulting summary query will contain  $\alpha$  and in turn be contained in  $\gamma$ . Simultaneously applying such replacement to all non-atomic non-branching paths of maximal length then yields a polynomial sized summary query.

Further, given a  $k$ -summary query  $\sigma$ , one can test in NP whether there exists an element of  $\text{App}_{\mathcal{J}_{w_k}}^*(\gamma)$  that leads to such summary query by the process just mentioned. This is done by first checking that  $\sigma$  is of the “right shape” (essentially a query of tree-width  $k$  when disregarding the path- $l$  approximation atoms), and that it is contained in  $\gamma$  via a polynomial refinement  $\rho \in \text{Ref}(\gamma)$  and a strong onto homomorphism to the query resulting from replacing  $\mathbf{P}_l(X, Y, \delta)$  atoms with  $\delta$  in  $\sigma$ . ◀

## 6.2 Semantic tree-width problem

With the previous results in place, we now show that the semantic tree-width  $k$  problem is in  $\Pi_2^p$  for  $\text{UCRPQ}(\text{SRE})$ , for every  $k > 1$ .

► **Theorem 6.1.** *For  $k > 1$ , the semantic tree-width  $k$  problem for  $\text{UCRPQ}(\text{SRE})$  is in  $\Pi_2^p$ .*

**Proof.** It suffices to show the statement for any  $\text{CRPQ}(\text{SRE}) \gamma$ . Remember that  $\gamma$  is of semantic tree-width  $k$  if, and only if,  $\gamma \subseteq \text{App}_{\mathcal{J}_{w_k}}^{\text{zip}}(\gamma)$ . The first ingredient to this proof is the fact that this containment has a polynomial counter-example property:

See the proof of  
Claim 6.3 in  
page 28.

▷ **Claim 6.3.** If  $\gamma \not\subseteq \text{App}_{\mathcal{J}_{w_k}}^{\text{zip}}(\gamma)$  then there is a polynomial-sized expansion  $\xi$  of  $\gamma$  such that  $\xi \not\subseteq \text{App}_{\mathcal{J}_{w_k}}^{\text{zip}}(\gamma)$ .

This is because any path  $x_0 \xrightarrow{a} \dots \xrightarrow{a} x_m$  in an expansion  $\xi$  of  $\gamma$  such that  $m > \|\gamma\|$  and  $\xi \subseteq \text{App}_{\mathcal{J}_{w_k}}^{\text{zip}}(\gamma)$  can be ‘pumped’ to an even longer path of any length greater than  $m$  obtaining another expansion  $\xi'$  such that  $\xi' \subseteq \text{App}_{\mathcal{J}_{w_k}}^{\text{zip}}(\gamma)$ . This implies that a minimal counterexample must be of polynomial size.

The second ingredient is that testing whether CQ is a counterexample is in NP.

See the proof of  
Claim 6.4 in  
page 28.

▷ **Claim 6.4.** The problem of testing, given a CQ  $\gamma$ , whether  $\gamma \subseteq \text{App}_{\mathcal{J}_{w_k}}^{\text{zip}}(\gamma)$ , is in NP.

**Informal proof.** We first guess a polynomial-sized  $k$ -summary query  $\delta_{\text{zip}}$  and test in NP that it is part of  $\text{App}_{\mathcal{J}_{w_k}}^{\text{zip}}(\gamma)$  by Proposition 6.2. We now guess a valuation  $\mu : \text{vars}(\delta_{\text{zip}}) \rightarrow \text{vars}(\gamma)$  and test that it is a homomorphism  $\xi \xrightarrow{\text{hom}} \gamma$ , where  $\xi$  is an expansion of the CRPQ resulting from discarding the path- $l$  approximation atoms of  $\delta_{\text{zip}}$ , which can be done in NL. It remains to check that each atom  $\mathbf{P}_l(X, Y, \hat{\delta})$  of  $\delta$  has an expansion  $\hat{\xi}$  such that  $h : \hat{\xi} \xrightarrow{\text{hom}} \gamma$  for a homomorphism such that  $h(x) = \mu(x)$  for all  $x \in X \cup Y$ . This can be done via an NL algorithm using  $l + 1$  pointers to traverse the width- $l$  nice path decomposition of  $\hat{\xi}$ , simultaneously guessing  $\hat{\xi}$ , the expansion of  $\gamma$  which homomorphically maps to  $\hat{\xi}$ , and the valuation of  $h$ . ◀

As a consequence of the two claims, we obtain a  $\Sigma_2^p$  algorithm for non-containment of  $\gamma \subseteq \text{App}_{\mathcal{J}_{w_k}}^{\text{zip}}(\gamma)$ . We first guess an expansion  $\xi$  of  $\gamma$  of polynomial size, and we then test  $\xi \not\subseteq \text{App}_{\mathcal{J}_{w_k}}^{\text{zip}}(\gamma)$  in coNP. This gives a  $\Pi_2^p$  algorithm for the semantic tree-width  $k$  problem. ◀

## 7 Discussion

We have studied the definability and approximation of  $\text{UC2RPQ}$  queries by queries of bounded tree-width and shown that the maximal under-approximation in terms of an infinitary union of conjunctive queries of tree-width  $k > 1$  can be always effectively expressed as a  $\text{UC2RPQ}$  of tree-width  $k$  (Corollary 3.8). However, while the semantic tree-width 1 problem as shown to be EXSPACE-complete [3, Theorem 6.1, Proposition 6.2], we have left a gap between our lower and upper bounds in Theorem 1.3.

▷ Question 7.1. For  $k > 1$ , is the semantic tree-width  $k$  problem EXPSPACE-complete?

We also do not know whether the  $\Pi_2^P$  bound on the semantic tree-width  $k$  problem for UCRPQ(SRE) has a matching lower bound. The known lower bound for the UCRPQ(SRE) containment problem [11, Theorem 5.1] does not seem to be useful to be used in a reduction, since it necessitates queries of arbitrary high tree-width.

It is worth stressing that in [3] two-way navigation plays a crucial part to prove the existence of the maximal under-approximation by a UC2RPQ of tree-width 1 [3, Theorem 5.2], but this feature plays no role whatsoever in our proof—see Theorem 1.5 and Remark 1.6. Moreover,  $\mathcal{F}w_k$  queries enjoy the very nice property of being closed under refinement when  $k > 1$  but not when  $k = 1$ —see Example 3.6—which forces Barceló, Romero, and Vardi to introduce the notion of “pseudoacyclicity” [3, §5.2.1], namely the greatest subclass of  $\mathcal{F}w_1$  closed under refinement, while we can directly work with the rather comfortable  $\text{App}_{\mathcal{F}w_k}^*(\gamma)$ . On the other hand, graphs of tree-width  $k > 1$  being combinatorially less trivial than graphs of tree-width 1, our proof must carefully handle this information, using tagged tree decompositions of Section 4. Similarly to [3, Theorem 6.3] for the case  $k = 1$ , our results implies that for each  $k > 1$  the evaluation problem for UC2RPQs  $\Gamma$  of semantic tree-width  $k$  is fixed-parameter tractable (FPT) in the size of the query, *i.e.* in  $O(|G|^{k+1} \cdot f(|\Gamma|))$  for a computable function  $f$ , where  $G$  is the database given as input. This improves the dependence on the size of the database, namely  $O(|G|^{2k+2} \cdot f(|\Gamma|))$ , proven by Romero, Barceló and Vardi [15, Corollary IV.12].

▷ Question 7.2 (Also mentioned in [15, §IV-(4)]). Does every r.e. class of CRPQs with FPT evaluation has bounded semantic tree-width?

## References

- 1 Renzo Angles, Marcelo Arenas, Pablo Barceló, Aidan Hogan, Juan Reutter, and Domagoj Vrgoč. Foundations of modern query languages for graph databases. *ACM Comput. Surv.*, 50(5), sep 2017. doi:10.1145/3104031.
- 2 Stefan Arnborg, Derek G Corneil, and Andrzej Proskurowski. Complexity of finding embeddings in a k-tree. *SIAM Journal on Algebraic Discrete Methods*, 8(2):277–284, 1987.
- 3 Pablo Barceló, Miguel Romero, and Moshe Y. Vardi. Semantic acyclicity on graph databases. *SIAM Journal on computing*, 45(4):1339–1376, 2016. doi:10.1137/15M1034714.
- 4 Pablo Barceló Baeza. Querying graph databases. In *Proceedings of the 32nd ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems*, PODS ’13, page 175–188, New York, NY, USA, 2013. Association for Computing Machinery. doi:10.1145/2463664.2465216.
- 5 Hans L. Bodlaender. A partial k-arboretum of graphs with bounded treewidth. *Theoretical Computer Science*, 209(1):1–45, 1998. doi:https://doi.org/10.1016/S0304-3975(97)00228-4.
- 6 Angela Bonifati, Wim Martens, and Thomas Timm. An analytical study of large SPARQL query logs. *VLDB Journal*, 29(2):655–679, 2020. doi:10.1007/s00778-019-00558-9.
- 7 Diego Calvanese, Giuseppe De Giacomo, Maurizio Lenzerini, and Moshe Y. Vardi. Containment of conjunctive regular path queries with inverse. In *Principles of Knowledge Representation and Reasoning (KR)*, pages 176–185, 2000.
- 8 Ashok K. Chandra and Philip M. Merlin. Optimal implementation of conjunctive queries in relational data bases. In John E. Hopcroft, Emily P. Friedman, and Michael A. Harrison, editors, *Proceedings of the 9th Annual ACM Symposium on Theory of Computing, May 4-6, 1977, Boulder, Colorado, USA*, pages 77–90. ACM, 1977. doi:10.1145/800105.803397.
- 9 Diego Figueira. Containment of UC2RPQ: the hard and easy cases. In *International Conference on Database Theory (ICDT)*, Leibniz International Proceedings in Informatics (LIPIcs). Leibniz-Zentrum für Informatik, 2020.

- 10    Diego Figueira. Foundations of graph path query languages - course notes for the reasoning web summer school 2021. In *Reasoning Web. Declarative Artificial Intelligence - 17th International Summer School 2021, Leuven, Belgium, September 8-15, 2021, Tutorial Lectures*, volume 13100 of *Lecture Notes in Computer Science*, pages 1–21. Springer, 2021. doi:10.1007/978-3-030-95481-9\_1.
- 11    Diego Figueira, Adwait Godbole, S. Krishna, Wim Martens, Matthias Niewerth, and Tina Trautner. Containment of simple conjunctive regular path queries. In *Principles of Knowledge Representation and Reasoning (KR)*, 2020. URL: <https://hal.archives-ouvertes.fr/hal-02505244>.
- 12    Daniela Florescu, Alon Levy, and Dan Suciu. Query containment for conjunctive queries with regular expressions. In *ACM Symposium on Principles of Database Systems (PODS)*, pages 139–148. ACM Press, 1998. doi:10.1145/275487.275503.
- 13    Ton Kloks. *Treewidth: computations and approximations*. Lecture Notes in Computer Science. Springer, 1994. doi:10.1007/BFb0045375.
- 14    Jaroslav Nešetřil and Patrice Ossona de Mendez. *Sparsity - Graphs, Structures, and Algorithms*, volume 28 of *Algorithms and combinatorics*. Springer, 2012. doi:10.1007/978-3-642-27875-4.
- 15    Miguel Romero, Pablo Barceló, and Moshe Y. Vardi. The homomorphism problem for regular graph patterns. In *Annual Symposium on Logic in Computer Science (LICS)*, pages 1–12. IEEE Computer Society Press, 2017. doi:10.1109/LICS.2017.8005106.

## A Classes of languages closed under sublanguages

The following classes are [closed under sublanguages](#):

- $\mathcal{L}_a \triangleq \{\{a\} \mid a \in \mathbb{A}\}$ ,
- $\mathcal{L}_{A,a^*} \triangleq \{\{a_1 + \dots + a_n\} \mid a_1, \dots, a_n \in \mathbb{A}\} \cup \{a^* \mid a \in \mathbb{A}\}$ ,
- $\mathcal{L}_{A,A^*} \triangleq \{\{a_1 + \dots + a_n\} \mid a_1, \dots, a_n \in \mathbb{A}\} \cup \{(a_1 + \dots + a_n)^* \mid a_1, \dots, a_n \in \mathbb{A}\}$ ,
- $\mathcal{L}_{\text{mod } n} \triangleq \{\{a\}\} \cup \{a^i(a^n)^* \mid 0 \leq i < n\}$ ,
- the closure under concatenation of any class [closed under sublanguages](#),
- an arbitrary union of classes [closed under sublanguages](#),
- as a consequence any class of the form

$$\mathcal{L}_a \cup \left\{ \bigcup_{q \in Q, q' \in Q'} \mathcal{A}_i[q, q'] \mid Q, Q' \text{ sets of states of } \mathcal{A}_i, i \in I \right\}$$

where  $(\mathcal{A}_i)_{i \in I}$  is an arbitrary family of automata, and

- in particular, the classes of star-free languages (*a.k.a.* aperiodic languages, *a.k.a.* first-order definable languages), of piecewise testable languages, of commutative languages, etc.

## B Missing details to Section 2

### B.1 Equality atoms

We also work with [CRPQs](#) with *equality atoms*, which are queries of the form  $\gamma(\bar{x}) = \delta \wedge I$ , where  $\delta$  is a [CRPQ](#) (without equality atoms) and  $I$  is a conjunction of [equality atoms](#) of the form  $x = y$ . Again, we denote by  $\text{vars}(\gamma)$  the set of variables appearing in the (equality and non-equality) atoms of  $\gamma$ . We define the binary relation  $=_\gamma$  over  $\text{vars}(\gamma)$  to be the reflexive-symmetric-transitive closure of the binary relation  $\{(x, y) : x = y \text{ is an equality atom in } \gamma\}$ . In other words, we have  $x =_\gamma y$  if the equality  $x = y$  is forced by the [equality atoms](#) of  $\gamma$ . Note that every [CRPQ](#) with [equality atoms](#)  $\gamma(\bar{x}) = \delta \wedge I$  is equivalent to a [CRPQ](#) without [equality atoms](#)  $\gamma^\approx$ , which is obtained from  $\gamma$  by collapsing each equivalence class of the relation  $=_\gamma$  into a single variable. This transformation gives us a *canonical* renaming from  $\text{vars}(\gamma)$  to  $\text{vars}(\gamma^\approx)$ .

### B.2 Proof of Fact 2.4

▷ **Fact 2.4.** Let  $k > 1$  and let  $\gamma$  be a [C2RPQ](#) of [tree-width](#) at most  $k$ . Then any [refinement](#) of  $\gamma$  has [tree-width](#) at most  $k$ .

**Proof of Fact 2.4.** The underlying graph of a [refinement](#) of  $\gamma$  is obtained from the underlying graph of  $\gamma$  by either contracting some edges (when dealing with [equality atoms](#)), or by replacing a single edge by a path of edges (where the non-extremal nodes are new nodes).

This first operation preserves [tree-width](#) at most  $k$  (even if  $k = 1$ ), see *e.g.* [5, Lemma 16]. The second operation preserves [tree-width](#) at most  $k$ , provided that  $k > 1$ : if a graph  $G'$  is obtained from a graph  $G$  by replacing an edge  $x_0 \rightarrow x_n$  by a path  $x_0 \rightarrow x_1 \rightarrow \dots \rightarrow x_n$ , then from a [tree decomposition](#) of  $G$  it suffices to pick a [bag](#) containing both  $x_0$  and  $x_n$ , and add a branch to the tree, rooted at this [bag](#), and containing bags with nodes

$$\begin{aligned} &\{x_0, x_n\}, \{x_0, x_1, x_n\}, \{x_1, x_n\}, \{x_1, x_2, x_n\}, \\ &\dots, \\ &\{x_{n-2}, x_n\}, \{x_{n-2}, x_{n-1}, x_n\}, \{x_{n-1}, x_n\}, \{x_n\}. \end{aligned}$$

All bags contain at most three nodes, so we obtain **tree decomposition** of  $G'$  whose **width** is the maximum between 2 and the **width** of the original **tree decomposition** of  $G$ .  $\triangleleft$

## C Missing details to Section 3

### C.1 Example 3.6

Recall that we have the followings queries:

$$\gamma(x, y) \hat{=} \begin{array}{ccc} z & \xrightarrow{(ba)^*} & z' \\ \uparrow a & & \downarrow b \\ x & \xrightarrow{(ab)^+} & y \\ & \xleftarrow{c} & \end{array} \quad \text{and} \quad \delta(x, y) \hat{=} x \xrightleftharpoons[c]{(ab)^+} y.$$

and we want to show that  $\delta(x, y) \in \text{App}_{\mathcal{T}_{w_1}}(\gamma)$  but  $\delta(x, y) \not\subseteq \text{App}_{\mathcal{T}_{w_1}}^*(\gamma)$ .

First note that  $\delta(x, y) \in \mathcal{T}_{w_1}$  and  $\delta(x, y) \subseteq \gamma(x, y)$ —we even have  $\delta(x, y) \equiv \gamma(x, y)$ —so  $\delta(x, y) \in \text{App}_{\mathcal{T}_{w_1}}(\gamma)$ . However,  $\delta(x, y) \not\subseteq \text{App}_{\mathcal{T}_{w_1}}^*(\gamma)$ : informally, this is because to obtain a query in  $\text{App}_{\mathcal{T}_{w_1}}^*(\gamma)$ , we must start from a **refinement** of  $\rho(x, y)$ , and, after applying some **strong onto homomorphism**, obtain a query of **tree-width** 1. This last condition forces us to either have  $f(x) = f(y)$ , or to send the variables  $z$  and  $z'$  in the **refinement** of  $\rho$  on  $\{f(x), f(y)\}$ , meaning that we obtain a **C2RPQ** with a self-loop. We conclude by observing an **expansion** of  $\delta(x, y)$  cannot be **contained** in such a query.

### C.2 EXPSPACE upper-bound for the semantic tree-width $k$ problem

First stated in  
page 10.

► **Lemma 3.9.** *For  $k \geq 1$ , the **semantic tree-width  $k$  problem** for **UC2RPQ** is in 2EXPSPACE.*

**Proof of Lemma 3.9.** If  $k = 1$ , use the algorithm [3, Theorem 6.1] of Barceló, Romero and Vardi, which runs in EXPSPACE.  $\text{App}_{\mathcal{T}_{w_k}}(\Gamma) \hat{=} \bigvee_{\gamma \in \Gamma} \text{App}_{\mathcal{T}_{w_k}}(\gamma)$  being the **maximal under-approximation** of  $\Gamma$  by **infinitary unions** of **C2RPQs** of **tree-width  $k$** , we conclude by Lemma 3.7 that  $\bigvee_{\gamma \in \Gamma} \text{App}_{\mathcal{T}_{w_k}}^{*, \leq \ell_\gamma}(\gamma)$  is the **maximal under-approximation** of  $\Gamma$  by a **UC2RPQ** of **tree-width  $k$** , and hence  $\Gamma$  has **semantic tree-width  $k$**  if and only if

$$\Gamma \equiv \bigvee_{\gamma \in \Gamma} \text{App}_{\mathcal{T}_{w_k}}^{*, \leq \ell_\gamma}(\gamma). \quad (3)$$

Since **semantical equivalence** is decidable in EXPSPACE in the size of the queries—see [12, Theorem 4.8] or [7, Theorem 5]—and since the right-hand side of Equation (3) has exponential size in  $\|\Gamma\|$ —indeed, recall that  $\ell_\gamma = \Theta(\|\gamma\|^2 \cdot (k+1)^{\|\gamma\|+1})$ —the conclusion follows.  $\blacktriangleleft$

### C.3 EXPSPACE-hardness for the semantic tree-width $k$ problem

Our proof of Lemma 3.10 is a direct adaptation from the proof of [3, Proposition 6.2], which prove the EXPSPACE-hardness of the **semantic tree-width 1 problem**. We say that a **C2RPQ** is **connected** when its underlying undirected graph is connected. We first give a small usefull fact:

► **Fact C.1** (Implicit in [3, Proof of Proposition 6.2]).

1. Let  $G, G'$  be two **databases** and  $\gamma$  be a **connected Boolean C2RPQ**. If the disjoint sum  $G + G'$  satisfies  $\gamma$ , then either  $G$  satisfies  $\gamma$  or  $G'$  satisfies  $\gamma$ .



2. Let  $\gamma, \gamma'$  be **Boolean C2RPQs**. If  $\gamma'$  is **connected** and  $\gamma \subseteq \gamma'$ , then there exists a subquery  $\gamma'$  of  $\gamma$ , obtained as connected component of  $\gamma$ , such that  $\gamma' \subseteq \gamma$ .

Notice first that if  $\gamma$  and  $\gamma'$  are **CQs** then the proof of Fact C.1 follows directly from the equivalence of  $\gamma \subseteq \gamma'$  (resp.  $G$  satisfies  $\gamma$ ) and the existence of a **homomorphism** from  $\gamma'$  to  $\gamma$  (resp.  $\gamma$  to  $G$ ).

**Proof.** The first point follows directly from the second one, so let us prove the latter. Write  $\gamma \hat{=} \gamma_1 \wedge \dots \wedge \gamma_n$  where  $\gamma_1, \dots, \gamma_n$  are connected components of  $\gamma$ , and assume by contradiction that for all  $i$ ,  $\gamma_i \not\subseteq \gamma'$ . Then there exists a database  $G_i$  such that  $G_i$  satisfies  $\gamma_i$  but not  $\gamma'$ . Consider the disjoint union  $G = G_1 + \dots + G_n$ .

On the one hand, since the  $\gamma_i$ 's have disjoint variables and  $\gamma_i$  satisfies  $G_i$  for each  $i$ , then  $G$  satisfies  $\gamma$ . On the other hand,  $G$  cannot satisfy  $\gamma'$  since  $\gamma'$  is **connected** and none of the  $G_i$  satisfies  $\gamma'$ —otherwise, the **homomorphism** from  $\gamma'$  to  $G$  would yield a **homomorphism** from  $\gamma'$  to some  $G_i$ . This contradicts the **containment** of  $\gamma$  in  $\gamma'$ . ◀

We can then prove Lemma 3.10.

► **Lemma 3.10.** *The semantic tree-width  $k$  problem is EXPSPACE-hard, even if restricted to Boolean CRPQs.*

First stated in page 10.

**Proof of Lemma 3.10.** Fix  $k \geq 1$ . The proof contains two parts:

1. first, we reduce the decision of problem of, given two **Boolean CRPQs**  $\gamma$  and  $\gamma'$ , where  $\gamma$  has **tree-width**  $k$ ,  $\gamma'$  is **connected** and does *not* have **semantic tree-width**  $k$ , deciding whether  $\gamma \subseteq \gamma'$ , to the **semantic tree-width  $k$  problem**;
2. then, we prove that the former problem is EXPSPACE-hard.

(1): We reduce the instance  $(\gamma, \gamma')$  to the instance  $\gamma \wedge \gamma'$ . We simply have to check that  $\gamma \subseteq \gamma'$  if and only if  $\gamma \wedge \gamma'$  has **semantic tree-width**  $k$ . The left-to-right implication is straightforward since  $\gamma \subseteq \gamma'$  implies that  $\gamma \wedge \gamma' \equiv \gamma$  and  $\gamma$  was assumed to have **tree-width**  $k$ . For the converse implication, if  $\gamma \wedge \gamma' \equiv \delta$  where  $\delta$  is a **UC2RPQ** of **tree-width**  $k$  then write  $\delta = \bigvee_{i=1}^n \delta_i$  where the  $\delta_i$ 's are **C2RPQs** and let  $\delta_{i,1}, \dots, \delta_{i,k_i}$  be the connected components of  $\delta_i$ .

Since for each  $i$  we have  $\delta_i \subseteq \delta \equiv \gamma \wedge \gamma' \subseteq \gamma'$ , by Fact C.1, there exists  $j_i$  such that  $\delta_{i,j_i} \subseteq \gamma'$ . Let  $\delta' \hat{=} \bigvee_{i=1}^n \delta_{i,j_i}$  so that, by construction  $\delta' \subseteq \gamma'$ . However, note that  $\delta'$  has **tree-width** at most  $k$  but  $\gamma'$  was assumed not to have **semantic tree-width**  $k$ , hence  $\delta' \not\subseteq \gamma'$ , so there exists  $G'$  such that:

$$G' \text{ satisfies } \gamma' \quad \text{and} \quad G' \text{ does not satisfy } \delta'. \quad (4)$$

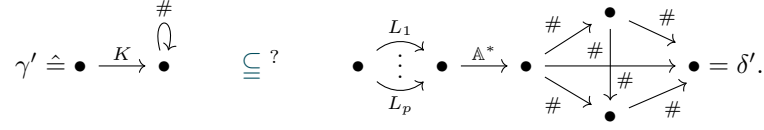
We now prove that  $\gamma \subseteq \gamma'$ . Let  $G$  be a **database satisfying**  $\gamma$ . Then  $G + G'$  satisfies  $\gamma \wedge \gamma'$  since  $G$  satisfies  $\gamma$ ,  $G'$  satisfies  $\gamma'$  and  $\gamma$  and  $\gamma'$  have disjoint variables. As a consequence,  $G + G'$  satisfies  $\delta$  and hence  $\delta'$ , so there exists  $i$  such that  $G + G'$  satisfies  $\delta_{i,j_i}$ . Since  $\delta_{i,j_i}$  is **connected**, either  $G$  satisfies  $\delta_{i,j_i}$  or  $G'$  satisfies  $\delta_{i,j_i}$ . By Equation (4), the latter cannot hold, so  $G$  satisfies  $\delta_{i,j_i}$  and hence  $\gamma'$ .

Therefore, we've shown that for each **database**  $G$  that **satisfies**  $\gamma$ , then  $G$  satisfies  $\gamma'$ , i.e.  $\gamma \subseteq \gamma'$ . Overall,  $\gamma \wedge \gamma'$  has **semantic tree-width**  $k$  if and only if  $\gamma \subseteq \gamma'$ .

(2): It was shown in [9, Lemma 8] that the **containment** of **CRPQs** was still EXPSPACE-hard when restricted to inputs of the form:

$$\gamma = \bullet \xrightarrow{K} \bullet \quad \subseteq ? \quad \bullet \begin{array}{c} \xrightarrow{L_1} \\ \vdots \\ \xrightarrow{L_p} \end{array} \bullet = \delta,$$

where  $K, L_1, \dots, L_p$  are regular languages over  $\mathbb{A}$ . Reduce it to the following problem:



where the right-hand side of  $\delta'$  is a directed  $(k+2)$ -clique and where  $\#$  is a new symbol, i.e.  $\# \notin \mathbb{A}$ .

We claim that  $\gamma \subseteq \delta$  if and only if  $\gamma' \subseteq \delta'$ . The direct implication is direct and the converse implication simply relies on the fact that  $\# \notin \mathbb{A}$ . Then, observe that  $\gamma'$  has [tree-width](#)  $1 \leq k$ , and that  $\delta'$  is connected but do not have [semantic tree-width](#) at most  $k$ .

To prove the last point, consider a [UC2RPQ](#)  $\Delta''$  that is equivalent to  $\delta'$ . Pick any [expansion](#)  $\xi'_1$  of  $\delta'$ . Since  $\Delta'' \subseteq \delta'$ , there exists an [expansion](#)  $\xi''$  of  $\Delta''$  such that there is a [homomorphism](#) from  $\xi'_1$  to  $\xi''$ . Dually, since  $\delta' \subseteq \Delta''$ , there exists an [expansion](#)  $\xi'_2$  of  $\delta'$  such that there is a [homomorphism](#) from  $\xi''$  to  $\xi'_2$ . Overall, we have [homomorphisms](#)  $\xi'_1 \rightarrow \xi'' \rightarrow \xi'_2$ . Since  $\xi'_1$  and  $\xi'_2$  are both [expansions](#) of  $\delta'$ , they contain a  $\#$ -labelled directed  $(k+2)$ -clique, and the  $\#$ -letter appears nowhere else. Should the [homomorphism](#)  $\xi'_1 \rightarrow \xi''$  not be injective,  $\xi''$  would contain a  $\#$ -labelled self-loop, and hence, the [homomorphism](#)  $\xi'' \rightarrow \xi'_2$  would yield a  $\#$ -self loop in  $\xi'_2$ , which does not exist! Hence, the [homomorphism](#) from  $\xi'_1$  to  $\xi''$  is injective on the  $(k+2)$ -clique. As a result,  $\xi''$  contains a  $(k+2)$ -clique and has [tree-width](#) at least  $k+1$ . We conclude that  $\delta''$  has [tree-width](#) at least  $k+1$ .

Hence, we have shown that  $\gamma \subseteq \delta$  if and only if  $\gamma' \subseteq \delta'$  where  $\gamma'$  has [tree-width](#) at most  $k$ , where  $\delta'$  is [connected](#) and has [semantic tree-width](#) at least  $k+1$ . Since our reduction can be implemented in polynomial time, we conclude that the problem of Lemma 3.10 (2) is EXPSpace-hard. ◀

## D Path induced in a tagged tree decomposition

Recall that we are given an [homomorphism](#)  $f: \rho \rightarrow \alpha$  between [C2RPQs](#), and a [tagged tree decomposition](#)  $(T, \mathbf{v}, \mathbf{t})$  of  $f$ . Consider a path  $\pi$  in  $\rho$ , say  $x_0 \xrightarrow{\lambda_1} x_1 \xrightarrow{\lambda_2} \dots \xrightarrow{\lambda_n} x_n$ .

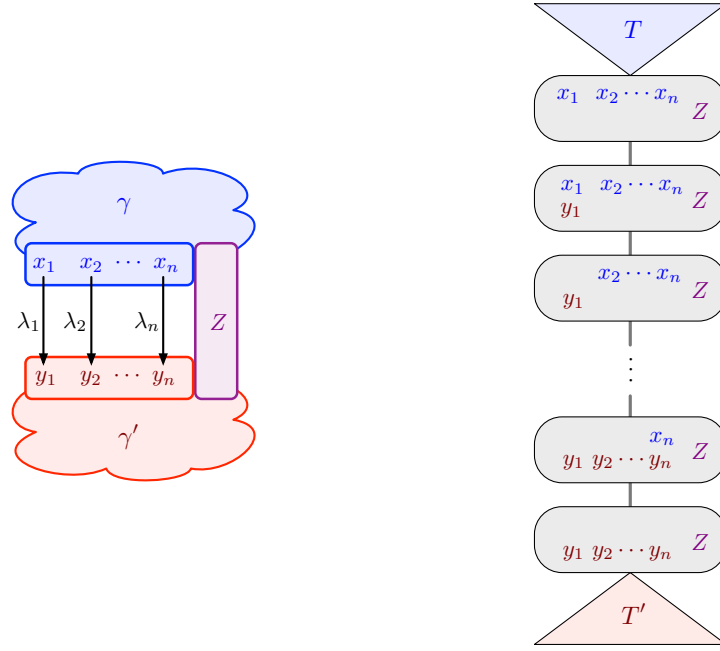
► **Definition D.1** (Path induced in a tagged tree decomposition). *The path induced by  $\pi$  is the unique sequence*

$$\left\langle \left( \begin{smallmatrix} b_1 \\ f(x_0) \end{smallmatrix} \right), \dots, \left( \begin{smallmatrix} b_{i_0} \\ f(x_0) \end{smallmatrix} \right), \left( \begin{smallmatrix} b_{i_0+1} \\ f(x_1) \end{smallmatrix} \right), \dots, \left( \begin{smallmatrix} b_{i_1} \\ f(x_1) \end{smallmatrix} \right), \dots, \left( \begin{smallmatrix} b_{i_{n-1}+1} \\ f(x_n) \end{smallmatrix} \right), \dots, \left( \begin{smallmatrix} b_{i_n} \\ f(x_n) \end{smallmatrix} \right) \right\rangle$$

of elements of  $V(T) \times \text{vars}(\alpha)$  such that for all  $k \in \llbracket 1, n \rrbracket$ :

- $i_1, \dots, i_n$  are strictly positive and  $i_1 = i_n = 1$ ;
- $x_k \in \mathbf{v}(b_j)$  for  $i_{k-1} < j \leq i_k$ , with the convention  $i_{-1} \hat{=} 0$ ;
- $b_{i_k} = b_{i_k+1} = \mathbf{t}(x_k \xrightarrow{\lambda_{k+1}} x_{k+1})$ ;
- $\langle b_{i_{k-1}+1}, \dots, b_{i_k} \rangle$  is the unique simple path from  $b_{i_{k-1}}$  to  $b_{i_k}$  in  $T$ .

In Section 4, we introduced the notion of [nice tree decompositions](#) of width  $k$ . One interesting property of such a decomposition is that in any of its [non-branching path](#), at least half of the [bags](#) contain at most  $k$  variables—recall that in a decomposition of width  $k$ , [bags](#) are allowed to contain at most  $k+1$  variables. Such [bags](#) are interesting for us because of the following property:



■ **Figure 6** The query  $\gamma \wedge \gamma' \wedge \bigwedge_{i=1}^n x_i \xrightarrow{\lambda_i} y_i$  (left-hand side) and one of its nice tagged tree decomposition of width at most  $k$  (right-hand side).

► **Proposition D.2.** Let  $\gamma, \gamma'$  be C2RPQs, and  $(T, \mathbf{v}, \mathbf{t})$  [resp.  $(T', \mathbf{v}', \mathbf{t}')$ ] be a nice tagged tree decomposition of width  $k$  of  $\gamma$  [resp.  $\gamma'$ ]. Let  $b, b'$  be leaves of  $T$  and  $T'$ , such that:

$$\mathbf{v}(b) = \{x_1, \dots, x_n\} \cup Z \quad \text{and} \quad \mathbf{v}(b') = \{y_1, \dots, y_n\} \cup Z,$$

where  $n + |Z| \leq k$ . Then  $\gamma \wedge \gamma' \wedge \bigwedge_{i=1}^n x_i \xrightarrow{\lambda_i} y_i$  has a nice tagged tree decomposition of width  $k$  in which the length of the non-branching paths is smaller than the sum of the longest non-branching paths of  $T$  and of  $T'$ , plus  $2k$ .

The proof of Appendix E.2 is elementary and illustrated in Figure 6.

**Proof.** The idea is to connect  $T$  with  $T'$  with  $2n \leq 2k$  bags: start from  $b_0 \triangleq b$ , which contains  $\{x_1, \dots, x_n\} \cup Z$ . Then create the following bags:

- $\mathbf{v}(b_1) \triangleq \{x_1, x_2, \dots, x_n\} \cup \{y_1\} \cup Z = \mathbf{v}(b_0) \cup \{y_1\}$ ,
- $\mathbf{v}(b_2) \triangleq \{x_2, \dots, x_n\} \cup \{y_1\} \cup Z = \mathbf{v}(b_1) \setminus \{x_1\}$ ,
- $\mathbf{v}(b_{2i-1}) \triangleq \{x_i, \dots, x_n\} \cup \{y_1, \dots, y_i\} \cup Z = \mathbf{v}(b_{2i-2}) \cup \{y_i\}$ ,
- $\mathbf{v}(b_{2i}) \triangleq \{x_{i+1}, \dots, x_n\} \cup \{y_1, \dots, y_i\} \cup Z = \mathbf{v}(b_{2i-1}) \setminus \{x_i\}$

for  $1 \leq i \leq n$ , and observe that  $\mathbf{v}(b_{2n}) = \mathbf{v}(b')$ . Then, tag the newly created atom  $x_i \xrightarrow{\lambda_i} y_i$  in the bag  $b_{2i-1}$ . ◀

## E Missing details to Section 5

### E.1 Missing proof to Claim 5.1

▷ **Claim 5.1.** For any trio  $f: \rho \xrightarrow{\text{hom}} \alpha$ , there exists a trio  $f': \rho' \xrightarrow{\text{hom}} \alpha'$  and a nice tagged tree decomposition  $(T', \mathbf{v}', \mathbf{t}')$  of width at most  $k$  of  $f'$  such that  $\alpha \subseteq \alpha'$ ,  $\|\rho'\| \leq \|\rho\|$  and every atom refinement of  $\rho'$  induces an acyclic path in the tree  $T'$ , in which case we say that  $(T', \mathbf{v}', \mathbf{t}')$  is locally acyclic.

First stated in page 12.

Proof of Claim 5.1. Let  $\pi$  be an **atom refinement** in  $\rho$ , say

$$\pi = x \xrightarrow{L_1} t_1 \xrightarrow{L_2} \dots \xrightarrow{L_{n-1}} t_{n-1} \xrightarrow{L_n} y,$$

that induce a **cyclic path** in  $T$ .

Our proof is illustrated in Figure 3b. In order to build the **trio**  $f: \rho' \xrightarrow{\text{hom}} \alpha'$  and a **tagged tree decomposition**  $T'$  of  $f'$  of width at most  $k$ , we will mainly use the fact that if two vertices  $(u, v)$  of some graph  $G$  belong to the same bag of a **tree decomposition**  $(T, \mathbf{v})$  of  $G$ , then  $(T, \mathbf{v})$  is still also a **tree decomposition** of the graph obtained by adding an edge from  $u$  to  $v$ .

Since the **induced path**

$$\mathbf{t}[\pi] \triangleq \left\langle \left( \begin{smallmatrix} b_1 \\ f(x) \end{smallmatrix} \right), \dots, \left( \begin{smallmatrix} b_{i_0} \\ f(x) \end{smallmatrix} \right), \left( \begin{smallmatrix} b_{i_0+1} \\ f(t_1) \end{smallmatrix} \right), \dots, \left( \begin{smallmatrix} b_{i_1} \\ f(t_1) \end{smallmatrix} \right), \dots, \left( \begin{smallmatrix} b_{i_{n-1}+1} \\ f(y) \end{smallmatrix} \right), \dots, \left( \begin{smallmatrix} b_{i_n} \\ f(y) \end{smallmatrix} \right) \right\rangle$$

is not **acyclic**, there exists  $j+1 < j'$  such that  $b_j = b_{j'}$ . Let  $n(j)$  (resp.  $n(j')$ ) denote the unique index such that  $i_{n(j)-1} < j \leq i_{n(j)}$  (resp.  $i_{n(j')-1} < j' \leq i_{n(j')}$ ). In particular, we have  $f(t_{n(j)}) \in \mathbf{v}(b_j)$  and  $f(t_{n(j')}) \in \mathbf{v}(b_{j'})$ . We claim that  $n(j) < n(j')$ .

We can then define

$$\pi' \triangleq t_0 \xrightarrow{L_1} \dots \xrightarrow{L_{n(j)}} t_{n(j)} \xrightarrow{K} t_{n(j')} \xrightarrow{L_{n(j')+1}} \dots \xrightarrow{L_n} t_n,$$

where  $K \triangleq [L_{n(j)+1} \dots L_{n(j')}]$  and let  $\rho'$  be the query obtained from  $\rho$  by replacing  $\pi$  with  $\pi'$ . Then, define  $\alpha'$  to be the query obtained from  $\alpha$  by adding an atom

$$f(z_{n(j)}) \xrightarrow{K} f(z_{n(j')}),$$

so that by construction, we have  $\alpha \subseteq \alpha'$ , that  $\rho' \in \text{Ref}(\gamma)$  with  $\|\rho'\| \leq \|\rho\|$  and  $f$  induces a morphism  $f': \rho' \rightarrow \alpha'$ .

We must then build a **tagged tree decomposition**  $(T', \mathbf{v}', \mathbf{t}')$  of  $f'$ . Simply start from the **tagged tree decomposition**  $(T, \mathbf{v}, \mathbf{t})$  of  $f$ , restrict  $\mathbf{t}$  to the atoms  $\text{Atoms}(\rho') \setminus \{z_{n(j)} \xrightarrow{K} z_{n(j')}\} \subseteq \text{Atoms}(\rho)$ , and tag the atom  $z_{n(j)} \xrightarrow{K} z_{n(j')}$  to the bag  $b_j = b_{j'}$ . This **tree decomposition** has the same **width** as  $T$ .

Observe then that the **path induced** by  $\pi'$  in  $(T', \mathbf{v}', \mathbf{t}')$  is simply

$$\mathbf{t}'[\pi'] = \left\langle \left( \begin{smallmatrix} b_1 \\ f(x) \end{smallmatrix} \right), \dots, \left( \begin{smallmatrix} b_j \\ f(t_{n(j)}) \end{smallmatrix} \right), \left( \begin{smallmatrix} b_{j'} \\ f(t_{n(j')}) \end{smallmatrix} \right), \dots, \left( \begin{smallmatrix} b_{i_n} \\ f(y) \end{smallmatrix} \right) \right\rangle$$

and thus  $\mathbf{t}'[\pi']$  is strictly shorter than  $\mathbf{t}[\pi]$  since  $j+1 < j'$ , by definition of these indices. Finally, observe that (1) we restrict  $\alpha'$  to be the image of  $f': \rho' \rightarrow \alpha'$ , in order to obtain a **strong onto homomorphism** and (2) if  $(T, \mathbf{v}, \mathbf{t})$  is **nice** then so is  $(T', \mathbf{v}', \mathbf{t}')$ .

Overall, we built  $f': \rho' \xrightarrow{\text{hom}} \alpha'$  together with a **nice tagged tree decomposition**  $(T', \mathbf{v}', \mathbf{t}')$  of **width** at most  $k$  where  $\alpha \subseteq \alpha'$ , and  $\rho' \in \text{Ref}(\gamma)$  is such that  $\|\rho'\| \leq \|\rho\|$ , and for each atom of  $\gamma$ , the **refinement** of this atom in  $\rho$  is exactly the same as the **refinement** of this atom in  $\rho'$  except possibly for one atom, for which the path induced in  $T'$  by its **refinement** in  $\rho'$  is strictly shorter than the **path induced** in  $T$  by its **refinement** in  $\rho$ . After iterating this construction as many times as needed, we obtain a **trio** as in the conclusion of Claim 5.1, which concludes our proof.  $\triangleleft$

## E.2 Proof of Claim 5.2 and related propositions

$\triangleright$  Subclaim 5.4. Suppose there are two **bags**  $b$  and  $b'$  such that: (i) they contain at most  $k$  nodes (*i.e.*, not **full bags**), (ii) they have the same **profile**, (iii) there is a **non-branching path** in  $T$  between these bags, and (iv) no **bags** of the path between  $b$  and  $b'$  (both included) are

atomic. Then, there exists a trio  $f': \rho' \xrightarrow{\text{hom}} \alpha'$  and a nice tagged tree decomposition of  $f'$  of width at most  $k$  that can be obtained by replacing the non-branching path between  $b$  and  $b'$  in the nice tagged tree-decomposition of  $f: \rho \xrightarrow{\text{hom}} \alpha$  by another non-branching path of at most  $2k + 1$  bags, such that  $\alpha \subseteq \alpha'$ .

We then give both an informal and a formal proof of Subclaim 5.4.

Informal proof of Subclaim 5.4. If  $b$  and  $b'$  have the same profile, then in particular they have the same cardinality, which is at most  $k$  by assumption. For the sake of simplicity, assume that it is exactly  $k$ , and let  $\mathbf{v}(b) = \{x_1, \dots, x_n, z_{n+1}, \dots, z_k\}$  and  $\mathbf{v}(b') = \{y_1, \dots, y_n, z_{n+1}, \dots, z_k\}$  be such that the  $x_i, y_i$  and  $z_i$ 's are distinct and  $\text{type}_{x_i} = \text{type}_{y_i}$  for all  $1 \leq i \leq n$ . Essentially, we can then contract every atom refinement in  $\rho$  of some atom occurring in a set of the form  $\text{type}_{x_i} = \text{type}_{y_i}$  for some  $i$ , as depicted in Figure 4. At this point, bags strictly comprised between  $b$  and  $b'$  are discarded, and so are variables of  $\alpha$  that do not occur anywhere else. We are left with two halves of a nice tagged tree decomposition that we need to merge, which can easily be done by using . The construction makes use of some crucial ingredients to guarantee its correctness:

- first, an atom  $y \xrightarrow{L} y'$  of  $\gamma$  cannot occur in two different types, allowing us to do the contraction of each atom refinement independently—this property is guaranteed by the definition of the profile;
- second, this contraction forces us to add new atoms in  $\alpha$  (to preserve the existence of a homomorphism from the refinement to the approximation) from some variables of  $\mathbf{v}(b)$  to some variables of  $\mathbf{v}(b')$ , but we never add an atom from  $z_i$  to  $z'_j$  for some  $i \neq j$ , which allows us to preserve the tree-width of the approximation.  $\triangleleft$

Proof of Subclaim 5.4. Let

$$\mathbf{v}(b) = \{x_1, \dots, x_n, z_{n+1}, \dots, z_k\} \text{ and } \mathbf{v}(b') = \{y_1, \dots, y_n, z_{n+1}, \dots, z_k\}$$

be as in the informal proof. Note that given an atom  $x \xrightarrow{L} y$  of  $\gamma$ , there is at most one variable of  $\alpha$  such that  $x \xrightarrow{L} y$  is in the type of this variable, because we assumed that refinements were sent on acyclic paths in  $T$ .

For every atom  $x \xrightarrow{L} y$  of  $\gamma$ , let  $\pi(x \xrightarrow{L} y) \triangleq t_0 \xrightarrow{L_1} \dots \xrightarrow{L_n} t_n$  be its refinement in  $q$ . If  $x \xrightarrow{L} y$  is not in some type of the profile of  $b$  (or equivalently, of  $b'$ ), leave it as it is. Otherwise, let  $i$  (resp.  $j$ ) be the unique<sup>7</sup> index such that  $\mathbf{t}[t_0 \xrightarrow{L_1} \dots \xrightarrow{L_n} t_n]$  leaves  $b$  at  $f(x_i)$  (resp. leaves  $b'$  at  $f(x_j)$ ).

Define  $\pi'(x \xrightarrow{L} y)$  as

$$t_0 \xrightarrow{L_1} \dots \xrightarrow{L_i} t_i \xrightarrow{[L_{i+1} \dots L_j]} t_j \xrightarrow{L_{j+1}} \dots \xrightarrow{L_n} t_n$$

when  $i \leq j$  and otherwise the definition is symmetric. Then, let  $\rho'$  be the refinement of  $\gamma$  obtained by simultaneously substituting  $\pi(x \xrightarrow{L} y)$  with  $\pi'(x \xrightarrow{L} y)$  in  $\rho$ , for every atom  $x \xrightarrow{L} y$  of  $\gamma$ .

Finally, let  $\alpha'$  be the query obtained by first adding the atoms

$$f(t_i) \xrightarrow{[L_{i+1} \dots L_j]} f(t_j),$$

and observe that  $f: \rho \rightarrow \alpha$  induces a homomorphism  $f': \rho' \rightarrow \alpha'$ . Moreover, by construction,  $\alpha \subseteq \alpha'$  (in fact,  $\alpha \equiv \alpha'$ ).

<sup>7</sup> by acyclicity

Then, build a **tagged tree decomposition**  $(T', \mathbf{v}', \mathbf{t}')$  of  $q'(\bar{x})$  by starting from  $(T, \mathbf{v}, \mathbf{t})$  and replacing the **non-branching path** in  $T$  from  $b$  to  $b'$  by a path  $(b, b_1, \tilde{b}_1, b_2, \tilde{b}_2, \dots, b_{k-1}, \tilde{b}_{k-1}, b_k, b')$  where:  $b_1, \dots, b_k$  and  $\tilde{b}_1, \dots, \tilde{b}_{k-1}$  are new bags such that:

- $\mathbf{v}'(b_i) \triangleq \{u'_1, \dots, u'_i, u_i, \dots, u_k\}$ ,
- $\mathbf{v}'(\tilde{b}_i) = \{u'_1, \dots, u'_i, u_{i+1}, \dots, u_k\}$ ,
- and  $\mathbf{t}'(f(x_i) \xrightarrow{[L_{i+1} \dots L_j]} f(x_j)) = b_m$ , where  $m$  is the unique integer such that the **type** of  $u_i$  (or equivalently, the **type** of  $u'_i$ ) contains the atom  $x \xrightarrow{L} y$ .

First, as we assumed that no bag in  $T$  between  $b$  and  $b'$  is **atomic**,  $\mathbf{t}'$  is still a **tagging** of the vertices of  $\alpha'$ . Then, since  $T$  has **width** at most  $k$  and each new **bag** contains at most  $k+1$  vertices, it follows that  $T'$  still has **width** at most  $k$ .

To conclude the proof, note that, as usual, we restrict  $\alpha'$  to be the image of  $f': \rho' \rightarrow \alpha'$  in order to obtain a **strong onto homomorphism**.  $\triangleleft$

In order to prove Claim 5.2, we need to be able to find bags  $b$  and  $b'$  as in the statement of Subclaim 5.4.

▷ **Fact E.1.** Let  $n, d, t \in \mathbb{N}$ . Let  $P$  be a set with at most  $n$  elements, and  $\tilde{P}$  be the disjoint union of  $P$  and  $\{\text{trap}, \text{avoid}\}$ . For every natural number  $m \geq 2(t+1)d(n+1) + 2t$ , for every sequence  $(p_i)_{0 \leq i < m} \in \tilde{P}^m$  containing at most  $t$  elements equal to **trap**, if at most half of the elements of the sequence are equal to **avoid**, then there exists  $i < i'$  such that  $x_i = p_{i'} \neq \text{avoid}$ ,  $i' - i \geq d$  and  $p_j \neq \text{trap}$  for every  $i \leq j \leq i'$ .

*Proof.* First extract from  $(x_i)_{0 \leq i < m}$  the subsequence of elements distinct from **avoid**, of length at least  $\lceil \frac{m}{2} \rceil \geq (t+1)d(n+1) + t$ . Then extract from it contiguous subsequences that avoid the **trap** element. Since there is at most  $t+1$  subsequences like this, one of them must have size at least  $d(n+1)$ . Denote by  $(y_i)_{0 \leq i < d(n+1)}$  the prefix of such a subsequence. Applying the pigeon-hole principle to  $(y_{i \cdot d})_{0 \leq i < n+1}$  yields the desired result.  $\triangleleft$

We can now prove Claim 5.2.

▷ **Claim 5.2.** Let  $f: \rho \xrightarrow{\text{hom}} \alpha$  be a **trio** and  $(T, \mathbf{v}, \mathbf{t})$  be a **locally acyclic nice tagged tree decomposition** of **width** at most  $k$  of  $f$ . Then there is a **trio**  $f': \rho' \xrightarrow{\text{hom}} \alpha'$  and a **nice tagged tree decomposition**  $(T', \mathbf{v}', \mathbf{t}')$  of **width** at most  $k$  of  $f'$  such that:

- $\alpha \subseteq \alpha'$ ,
- $(T', \mathbf{v}', \mathbf{t}')$  is **locally acyclic** w.r.t.  $f'$ , and
- the size of the longest **non-branching path** in  $T$  is at most  $\Theta(\|\gamma\| \cdot (k+1)^{\|\gamma\|+1})$ .

**Proof of Claim 5.2.** Let  $f: \rho \xrightarrow{\text{hom}} \alpha$  be a **trio**, and  $(T, \mathbf{v}, \mathbf{t})$  be a **locally acyclic nice tagged tree decomposition** of  $f$ . If there is a **non-branching path**  $(b_i)_{0 \leq i < m}$  in  $T$  of length at least  $m$ , let  $(p_i)_{0 \leq i < m}$  be the sequence defined by letting:

$$p_i \triangleq \begin{cases} \text{trap} & \text{if } b_i \text{ is atomic,} \\ \text{avoid} & \text{if } b_i \text{ contains } k+1 \text{ variables,} \\ \text{profile of } b_i & \text{otherwise.} \end{cases}$$

Applying Fact E.1 for  $n = (k+1)^{\|\gamma\|}$ ,  $d = 2k+1$ ,  $t = \|\gamma\|$  yields, under the assumption that

$$m \geq m_0 \triangleq 2(\|\gamma\| + 1)(2k+1)((k+1)^{\|\gamma\|} + 1) + 2k+1,$$

the existence of indices  $i < i'$  such that  $i' - i \geq 2k+1$ , and  $b_i$  and  $b_{i'}$  have the same profile, contain at most  $k$  variables, and every bag  $b_j$  for  $i \leq j \leq i'$  is **non-atomic**—note that the



hypothesis of Fact E.1 are satisfied since at most  $t = \|\gamma\|$  bags of  $(b_i)_{0 \leq i < m}$  are atomic, and assuming *w.l.o.g.* that no two consecutive bags of  $(b_i)_{0 \leq i < m}$  are identical, since the tagged tree decomposition  $(T, \mathbf{v}, \mathbf{t})$  of width  $k$  is nice, at most half of the bags contain  $k + 1$  variables. The assumption  $i' - i \geq 2k + 1$  means that the path from  $b_i$  to  $b_{i'}$  has length at least  $2k + 2$ , and thus applying Subclaim 5.4 will strictly shorten this path—this it has length at most  $2k + 1$ . Note that Subclaim 5.4 preserves the niceness of the tagged tree decomposition, its local acyclicity, and that the size of this tree decomposition is strictly smaller (in number of nodes) than the original tree decomposition. By iteratively applying this construction, we obtain a trio  $f': \rho' \xrightarrow{\text{hom}} \alpha'$  together with a locally acyclic nice tagged tree decomposition  $T'$  of width at most  $k$ , such that  $\alpha \subseteq \alpha'$  and every non-branching path of  $T'$  has length at most  $m_0 - 1 = \Theta(\|\gamma\| \cdot (k + 1)^{\|\gamma\| + 1})$ . ◀

## F Missing proofs to Section 6

► **Proposition 6.2.** *For every class  $\mathcal{L}$  closed under sublanguages, and for every  $\text{C2RPQ}(\mathcal{L})$   $\gamma$ ,  $\text{App}_{\mathcal{J}_{w_k}}(\gamma)$  can be expressed as a union of polynomial-sized  $k$ -summary queries having only  $\text{C2RPQ}(\mathcal{L})$  atoms. Further, one can test in NP if a summary query is part of this union. We call  $\text{App}_{\mathcal{J}_{w_k}}^{\text{zip}}(\gamma)$  to any such a union of summary queries.*

First stated in page 15.

**Proof of Proposition 6.2.** As corollary of the proof of Lemma 3.7, we can assume to have  $\text{App}_{\mathcal{J}_{w_k}}(\gamma)$  expressed as a union of  $\text{C2RPQ}(\mathcal{L})$  with a nice tree decomposition of width  $k$  with a linear number of leaves, and hence it suffices to replace non-branching paths with path- $l$  approximations. Concretely, any non-branching path of length greater than  $3\|\gamma\|$  must contain a sub-path whose every bag is non-atomic, and such that it starts and ends in bags of size at most  $k$  (by the niceness property). Such a non-atomic non-branching path can be “shrunk” using a path- $l$  approximation query.

Indeed, in the tree decomposition of any query  $\alpha \in \text{App}_{\mathcal{J}_{w_k}}^*(\gamma)$  witnessed by a trio  $f: \rho \xrightarrow{\text{hom}} \alpha$ , a long non-branching path from a node with a bag  $X$  to a node with a bag  $Y$  in the tree decomposition corresponds to a sub-query  $\delta(X', Y')$  of  $\alpha$  where  $X' = X \setminus Y$  and  $Y' = Y \setminus X$ . Such  $\delta$  is of the form  $\bigwedge_i \delta_i(x_i, y_i)$  where for each  $i$ ,  $x_i \in X'$ ,  $y_i \in Y'$ , and  $\delta_i$  is part of the refinement of an atom  $\rho_i$  of  $\delta$  in  $\rho$ , either of the form (a)  $\delta_i(x_i, y_i) = x_i \xrightarrow{L_1} z_1 \xrightarrow{L_2} \dots z_{n-1} \xrightarrow{L_n} y_i$  or of the form (b)  $\delta_i(x_i, y_i) = y_i \xrightarrow{L_1} z_1 \xrightarrow{L_2} \dots z_{n-1} \xrightarrow{L_n} x_i$  (where expansion atoms of the form  $z_i \xrightarrow{a} z_{i+1}$  and  $z_i \xleftarrow{a} z_{i+1}$  are represented as  $z_i \xrightarrow{\{a\}} z_{i+1}$  and  $z_i \xleftarrow{\{a\}} z_{i+1}$ , respectively). We can then safely replace  $\delta(X', Y')$  with a path- $l$  approximation query  $\mathbf{P}_l(X', Y', \delta')$  where  $l = k - |X \cap Y|$ . Here,  $\delta'$  contains an atom  $x' \xrightarrow{L'} y'$  for every  $\delta_i$  of the form (a) and an atom of the form  $y' \xrightarrow{L'} x'$  for every  $\delta_i$  of the form (b), where  $L'$  is a sublanguage of  $L$  such that  $L_1 \dots L_n \subseteq L'$ . This corresponds to the contraction of the atom refinement  $\rho_i$  between  $x_i$  and  $y_i$ . The summary query resulting from this replacement contains  $\alpha$  and in turn is contained in  $\gamma$ . Applying this replacement to all maximal non-atomic non-branching paths yields a summary query which is polynomially bounded in size.

Further, given a  $k$ -summary query  $\sigma$ , one can test in NP whether there exists an element of  $\text{App}_{\mathcal{J}_{w_k}}^*(\gamma)$  that leads to such summary query via the aforementioned procedure. For this, we need to first check that  $\sigma$  is of the “right shape” by checking that each connected component  $\sigma_i$  of  $\sigma$  obtained when disregarding the path- $l$  approximation atoms has a tree decomposition  $(T_i, \mathbf{v}_i)$  of width  $k$  such that

1. for every path- $l$  approximation atom  $A = \mathbf{P}_l(X, Y, \delta)$  of  $\sigma$  there exist components  $i, j$  and vertices  $x_A$  of  $T_i$  and  $y_A$  of  $T_j$  such that  $x_A$  contains the bag  $X$  and  $y_A$  contains the bag  $Y$ ;

2. the graph resulting from the disjoint union of all the trees  $T_i$ 's, plus the edges  $\{x_A, y_A\}$  for every *path- $l$  approximation* atom  $A$  of  $\sigma$  is a tree.

Finally, we guess some polynomial-sized *refinement*  $\rho$  of  $\gamma$  and we check that there is a *strong onto homomorphism* from  $\rho$  to  $\sigma'$ , where  $\sigma'$  is the *C2RPQ* resulting from replacing every  $\mathbf{P}_l(X', Y', \delta)$  atom with  $\delta$ .  $\blacktriangleleft$

First stated in  
page 16.

$\triangleright$  **Claim 6.3.** If  $\gamma \not\subseteq \text{App}_{\mathcal{T}\omega_k}^{\text{zip}}(\gamma)$  then there is a polynomial-sized *expansion*  $\xi$  of  $\gamma$  such that  $\xi \not\subseteq \text{App}_{\mathcal{T}\omega_k}^{\text{zip}}(\gamma)$ .

**Proof of Claim 6.3.** Let us call any atom with a language of the form  $a^*$  a *recursive atom*, and any other atom a *non-recursive atom*. Let  $n$  be the number of *non-recursive atoms* of  $\gamma$ . Hence, any refinement  $\rho \in \text{Ref}(\gamma)$  has  $n$  atoms deriving from *non-recursive atom refinements*, all the remaining ones derive from *recursive atom refinements*.

We will work with the infinitary union of *conjunctive queries*  $\mathcal{U} \triangleq \text{App}_{\mathcal{T}\omega_k}^*(\gamma) \cap \text{CQ}$ . Note that  $\mathcal{U} = \{\alpha \in \text{CQ} \mid \alpha \in \mathcal{T}\omega_k \text{ and there is } \xi \in \text{Exp}(\gamma) \text{ s.t. } \xi \xrightarrow{\text{hom}} \alpha\}$ . It is easy to see that  $\mathcal{U} \equiv \text{App}_{\mathcal{T}\omega_k}(\gamma)$  (and thus  $\mathcal{U} \equiv \text{App}_{\mathcal{T}\omega_k}^{\text{zip}}(\gamma)$ ). By Proposition 2.3  $\gamma \not\subseteq \mathcal{U}$  if, and only if, there is some *expansion*  $\xi$  of  $\gamma$  such that  $\xi \not\subseteq \mathcal{U}$ . In turn, this happens if, and only if, there is no  $\delta \in \mathcal{U}$  such that  $\delta \xrightarrow{\text{hom}} \xi$ .

Take any such counterexample  $\xi$  of minimal size (in number of atoms). We show that for any path of  $\xi$  of the shape  $\pi = x_0 \xrightarrow{a} x_1 \xrightarrow{a} x_2 \cdots x_{m-1} \xrightarrow{a} x_m$  for some  $a \in \mathbb{A}$  such that no  $x_i$  ( $0 < i < m$ ) is part of any atom not present in  $\pi^8$ , we have  $m \leq n + 1$ . Hence, since  $\xi$  is an *expansion* of a *CRPQ(SRE)*, this means that the size of  $\xi$  is linearly bounded in the size of  $\gamma$ .

By means of contradiction, if  $m > n + 1$  consider the *expansion*  $\xi'$  resulting from “shrinking” the path  $\pi$  to a path  $\pi'$  of length  $n + 1$ . Hence,  $\xi'$  is smaller than  $\xi$ , and since  $\xi$  is minimal,  $\xi'$  cannot be a counter-example. Thus, there is some  $\delta \in \mathcal{U}$  such that  $\delta \xrightarrow{\text{hom}} \xi'$ ; and by definition of  $\mathcal{U}$  we have  $f : \xi'' \xrightarrow{\text{hom}} \delta$  for some  $\xi'' \in \text{Exp}(\gamma)$ . Consider the composition  $\xi'' \xrightarrow{\text{hom}} \delta \xrightarrow{\text{hom}} \xi'$  and let us call it  $g : \xi'' \xrightarrow{\text{hom}} \xi'$ . By definition of  $n$  there must be at least one atom  $x_i \xrightarrow{a} x_{i+1}$  of the shrunk path  $\pi'$  of  $\xi'$  which either (i) does not have an  $f$ -preimage, or (ii) all its  $g$ -preimages proceed from atoms  $z \xrightarrow{a} z'$  of  $\xi''$  which proceed from *expansions* of *recursive atoms* of  $\gamma$ . We show that, in both cases, we can replace  $x_i \xrightarrow{a} x_{i+1}$  with a path of  $a$ 's of any arbitrary length  $l$ , obtaining a *conjunctive query*  $\xi'_{+l}$  which is not a counter-examples. In particular, for  $l = m - n - 1$ , we have  $\xi'_l = \xi$ , and this would contradict the fact that  $\xi$  is a counter-example. In the first case (i), we actually obtain that  $\delta \xrightarrow{\text{hom}} \xi'_{+l}$ . In the second case (ii), we have to replace each atom  $z \xrightarrow{a} z'$  in the  $g$ -preimage of  $x_i \xrightarrow{a} x_{i+1}$  in  $\xi''$  by an  $a$ -path of length  $l$ , obtaining some *expansion*  $\xi''_{+l}$  of  $\gamma$ . We also replace each atom in the  $f$ -preimage of  $x_i \xrightarrow{a} x_{i+1}$  by an  $a$ -path of length  $l$  obtaining some  $\delta_{+l}$  such that  $\xi''_{+l} \xrightarrow{\text{hom}} \delta_{+l} \xrightarrow{\text{hom}} \xi'_{+l}$ . Further,  $\delta_{+l} \in \mathcal{T}\omega_k$  since  $\mathcal{T}\omega_k$  is closed under subdivision (for any  $k > 1$ ). In both cases this shows that  $\xi'_{+l}$  is *not* a counter-example. Therefore, there exists a counter-example of polynomial (quadratic) size whenever  $\gamma \not\subseteq \text{App}_{\mathcal{T}\omega_k}^{\text{zip}}(\gamma)$ .  $\blacktriangleleft$

First stated in  
page 16.

$\triangleright$  **Claim 6.4.** The problem of testing, given a *CQ*  $\gamma$ , whether  $\gamma \subseteq \text{App}_{\mathcal{T}\omega_k}^{\text{zip}}(\gamma)$ , is in NP.

**Proof of Claim 6.4.** We first guess a polynomial-sized  $k$ -summary query  $\delta_{\text{zip}}$  and test in NP that it is part of  $\text{App}_{\mathcal{T}\omega_k}^{\text{zip}}(\gamma)$  by Proposition 6.2. Let us call  $\Delta$  be the equivalent *UCRPQ(SRE)* query, given by Proposition 6.2 *cum* Lemma 3.7. We have to check that there is some *expansion*  $\xi$  of  $\Delta$  such that there is a *homomorphism*  $\xi \xrightarrow{\text{hom}} \gamma$ . We first guess a

<sup>8</sup> This means that each  $x_i$  ( $0 < i < m$ ) is incident to exactly two atoms in  $\xi$ .

valuation  $\mu : \text{vars}(\delta_{\text{zip}}) \rightarrow \text{vars}(\gamma)$  for all the variables of the **summary query**  $\delta_{\text{zip}}$ . Now it remains to check that:

1. For every **CRPQ** atom  $x \xrightarrow{a^*} y$  of  $\delta_{\text{zip}}$  there is an  $a$ -path in  $\gamma$  from  $\mu(x)$  to  $\mu(y)$ .
2. Every **path- $l$  approximation**  $\mathbf{P}_l(X, Y, \bigwedge_{1 \leq i \leq n} A_i(x_i, y_i))$  of  $\delta_{\text{zip}}$  contains a **CQ**  $\delta_{\text{path}}$  admitting a **path decomposition** of **width**  $l$  which starts with the bag  $X$  and ends with  $Y$ . And further, there is a homomorphism  $h : \delta_{\text{path}} \xrightarrow{\text{hom}} \gamma$  which coincides with  $\mu$  on variables  $X \cup Y$ .

Observe that these two properties hold true if, and only if, there is some **expansion**  $\xi$  of  $\Delta$  such that  $\xi \xrightarrow{\text{hom}} \gamma$ . It is clear the first point can be achieved in NL since it is a simple reachability query. The second point can also be achieved in NL, since the **nice path decomposition** of **width**  $l$  can be guessed on-the-fly using  $l + 1$  pointers to the variables of  $\gamma$ . An NL algorithm can advance down the **path decomposition** while simultaneously

1. guessing the **conjunctive query**  $\delta_{\text{path}}$  via its **nice path decomposition** of **width**  $k$ ,
2. checking that there is a partial **homomorphism** to  $\gamma$ ,
3. ensuring that the **CQ**  $\delta_{\text{path}}$  being built is an element of  $\mathbf{P}_l(X, Y, \bigwedge_{1 \leq i \leq n} A_i(x_i, y_i))$ , which requires to also guess a homomorphism  $\rho \xrightarrow{\text{hom}} \delta_{\text{path}}$  from a refinement  $\rho$  of  $\bigwedge_{1 \leq i \leq n} A_i(x_i, y_i)$ .

Further, a simple test can ensure that the first and last bags of the decomposition coincide with the guessed assignment  $\mu$ . Since the number of pointers (bounded by  $k + 1$ ) is fixed, this algorithm is in NL.  $\blacktriangleleft$