

Conjunctive Regular Path Queries with String Variables

Markus L. Schmid¹

¹Humboldt-Universität zu Berlin, Unter den Linden 6, D-10099,
Berlin, Germany, MLSchmid@MLSchmid.de

December 20, 2019

Abstract

We introduce the class CXRPQ of conjunctive xregex path queries, which are obtained from conjunctive regular path queries (CRPQs) by adding string variables (also called backreferences) as found in practical implementations of regular expressions. CXRPQs can be considered user-friendly, since they combine two concepts that are well-established in practice: pattern-based graph queries and regular expressions with backreferences. Due to the string variables, CXRPQs can express inter-path dependencies, which are not expressible by CRPQs. The evaluation complexity of CXRPQs, if not further restricted, is PSpace-hard in data-complexity. We identify three natural fragments with more acceptable evaluation complexity: their data-complexity is in NL, while their combined complexity varies between ExpSpace, PSpace and NP. In terms of expressive power, we compare the CXRPQ-fragments with CRPQs and unions of CRPQs, and with extended conjunctive regular path queries (ECRPQs) and unions of ECRPQs.

1 Introduction

The popularity of graph databases (commonly abstracted as directed, edge-labelled graphs) is due to their applicability in a variety of settings where the underlying data is naturally represented as graphs, e. g., Semantic Web and social networks, biological data, chemical structure analysis, pattern recognition, network traffic, crime detection, object oriented data. The problem of querying graph-structured data has been studied over the last three decades and still receives a lot of attention. For more background information, we refer to the introductions of the recent papers [33, 18, 8, 9], and to the survey papers [4, 6, 43, 5].

Many query languages for graph databases (for practical systems as well as those studied in academia) follow an elegant and natural declarative approach: a query is described by a *graph pattern*, i. e., a graph $G = (V, E)$ with edge labels that represent some path-specifications. The evaluation of such a query consists in *matching* it to the graph database $\mathcal{D} = (V_{\mathcal{D}}, E_{\mathcal{D}})$, i. e., finding a mapping $h : V \rightarrow V_{\mathcal{D}}$, such that, for every $(x, s, y) \in E$, in \mathcal{D} there is a path from $h(x)$ to

$h(y)$ whose edge labels satisfy the path-specification s . In the literature, such query languages are also called *pattern-based*. Let us now briefly summarise where this concept can be found in theory and practice.

The most simple graph-patterns (called *basic* in [4]) have just fixed relations (i. e., edge-labels) from the graph database as their edge labels. A natural extension are *wildcards*, which can match any edge-label of the database (e. g., as described in [18]), or *label variables*, which are like wildcards, but different occurrences of the same variable must match the same label (see, e. g., [9]). It is common to extend such basic graph patterns with relational features like, e. g., projection, union, and difference (see [4]). In order to implement *navigational features* that can describe more complex connectivities between nodes via longer paths instead of only single arcs, we need more complicated path specifications.

Navigational features are popular, since they allow to query the *topology* of the data and, if transitivity can be described, exceed the power of the basic relational query languages. Using regular expressions as path specifications is the most common way of implementing navigational features. The *regular path queries* (RPQs) given by *single-edge* graph patterns ($\{x, y\}, \{(x, s, y)\}$), where s is a regular expression, can be considered the simplest navigational graph patterns. General graph patterns labelled by regular expressions are called *conjunctive regular path queries* (CRPQs).

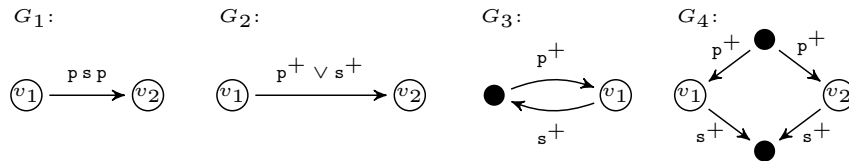


Figure 1: Simple graph patterns.

For example, consider a graph database with nodes representing persons, arcs (u, p, v) meaning “ u is a (biological) parent of v ” and arcs (u, s, v) meaning “ v is u ’s PhD-supervisor”. We consider the graph patterns from Figure 1 (labelled nodes are considered as free variables of the query). Then G_1 describes pairs (v_1, v_2) , where v_1 ’s child has been supervised by v_2 ’s parent; G_2 describes pairs (v_1, v_2) , where v_1 is a biological ancestor or an academical descendant of v_2 ; G_3 describes v_1 that have a biological ancestor that is also their academical ancestor; G_4 describes pairs (v_1, v_2) , where v_1 and v_2 are biologically related as well as academically. Note that G_1, G_2 represent RPQs, while G_3, G_4 represent CRPQs.

The classes of RPQs and CRPQs (and modifications of them) have been intensively studied in the literature (see, e. g., [17, 13, 14, 16, 10, 24, 38, 2, 1, 19]). The former can be evaluated efficiently (see, e. g., [16]), while evaluation for the latter it is NP-complete in combined-complexity, but NL-complete in data-complexity (see [8]).

Despite their long-standing investigation, these basic classes still pose several challenges that are currently studied. For example, [34, 36, 35] provide an in-depth analysis of the complexity of RPQs for different path semantics. So far in this introduction, we implicitly assumed *arbitrary path* semantics, but since there are potentially infinitely many such paths, query languages that also

retrieve paths often restrict this by considering simple paths or trails. However, such semantics make the evaluation of RPQs much more difficult (see [34, 36, 35] for details). Much effort has also been spent on extending RPQs and CRPQs to the setting where the data-elements stored at nodes of the graph database can also be queried (see [33, 32]). In [9], the authors represent partially defined graph data by graph patterns and then query them with CRPQs (among others). In the very recent paper [7], the authors study the boundedness problem for unions of CRPQs (i. e., the problem of finding an equivalent union of (relational) conjunctive queries).

Also in the practical world, pattern-based query languages for graph databases play a central role. Most prominently, the *W3C Recommendation* for *SPARQL 1.1* “is based around graph pattern matching” (as stated in Section 5 of [30]), and *Neo4J Cypher* also uses graph patterns as a core functionality (see [41]). Moreover, the graph computing framework *Apache TinkerPopTM* contains the graph database query language *Gremlin* [42], which is more based on the navigational graph traversal aspect, but nevertheless supports pattern-based query mechanisms. Note that [4] surveys the main features of these three languages.

1.1 Main Goal of this Work

CRPQs are not expressive enough for many natural querying tasks (see, e. g., the introduction of [8]). The most obvious shortcoming is that we cannot express any *inter-path* dependencies, i. e., relations between the paths of the database that are matched by the arcs of the graph pattern, except that they must start or end with the same node. The main goal of this work is to extend CRPQs in order to properly increase their expressive power. In particular, we want to meet the following objectives:

1. The increased expressive power should be reasonable, i. e., it should cover natural and relevant querying tasks.
2. The extensions should be user-friendly, i. e., the obtained query language should be intuitive.
3. The evaluation complexity should still be acceptable.

The main idea is to allow *string variables* in the edge labels of the graph patterns. For example, in G_1 of Figure 2, the $x\{a \vee b\}$ label sets variable x to some word matched by regular expression $a \vee b$ and the occurrence of x on the other edge label then refers to the value of x . Hence, G_1 describes all pairs (v_1, v_2) such that v_1 has a direct a -predecessor that has v_2 as a transitive successor with respect to a or c , or v_1 has a direct b -predecessor that has v_2 as a transitive successor with respect to b or c . Similarly, G_2 of Figure 2 describes triangles (v_1, v_2, v_3) with a complicated connectivity relation: v_1 reaches v_2 with aa or b , v_2 reaches v_3 with some path labelled only with symbols different than a and b (Σ is the set of edge labels), while v_3 reaches v_1 either in the same way as v_1 reaches v_2 , or in the same way as v_2 reaches v_3 .

The graph patterns G_1 and G_2 of Figure 2 could also be considered as CRPQs with some syntactic sugar. More precisely, both these graph patterns can be transformed into a union of CRPQs by simply “spelling out” all combinations that are possible for the variables x and y . However, we can easily build examples

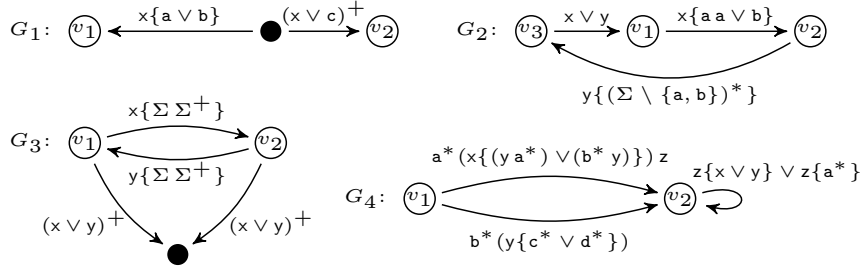


Figure 2: CRPQs with string variables.

that would translate into an exponential number of CRPQs, and, moreover, it can be argued that the necessity for explicitly listing *all* possible combinations that can conveniently be stated in a concise way, is exactly what a user should not be bothered with.

We move on to an example, where a simple application of string variables adds substantial expressive power to a CRPQ. Let us assume that the nodes in a graph database represent persons and arcs represent text messages sent by mobile phone (let Σ be the set of messages). The idea is that some individuals try to hide their direct communication by encoding their messages by sequences of simple text messages that are sent via intermediate senders and receivers. In particular, we want to discover individuals who are likely to be involved in such a hidden communication network. In this regard, G_3 of Figure 2 describes pairs (v_1, v_2) such that v_1 reaches v_2 (and v_2 reaches v_1) by a sequence x (y , respectively) of at least 2 messages, and there is some person that has been contacted by v_1 and v_2 by paths that are repetitions of these message-sequences. Note that, in this example, both the length of the message-sequences x and y , as well as the number of their repetitions in order to reach the mutual friend of v_1 and v_2 , are unbounded.

Finally, G_4 of Figure 2 shows a feature that has not been used in the previous examples: references to variables could also occur in the definitions of other variables, e.g., y is defined on one edge, and used in the definition of both x and z on the other two edges. Moreover, note that the same variable z has two definitions, which are mutually exclusive and therefore do not cause ambiguities.

This formalism obviously extends CRPQ; moreover, it is easy to see that it also covers wildcards and edge variables described above, as well as the fragment of *extended conjunctive regular path queries* (ECRPQs) [8] that only have equality-relations as non-unary relations (ECRPQs shall be discussed in more detail below).

1.2 Conjunctive Xregex Path Queries

Regular expressions of the kind used in the graph patterns of Figure 2 are actually a well-established concept, which, in the theoretical literature, is usually called *regex* or *xregex*, and string variables are often called *backreferences*. Xregex have been investigated in the formal language community [40, 39, 28, 15, 25], and, despite the fact that allowing them in regular expressions has many negative consequences (see [3, 22, 23, 21, 25]), regular expression libraries of

almost all modern programming languages (like, e. g., Java, PERL, Python and .NET) support backreferences (although they syntactically and even semantically slightly differ from each other (see the discussion in [28])), and they are even part of the POSIX standard [31]. The syntax of xregex is quite intuitive: on top of normal regular expressions, we can *define* variables by the construct $x\{\dots\}$ and *reference* them by occurrences of x (see Figure 2). Formally defining their semantics is more tricky (mainly due to nested variable definitions and undefined variables), but for simple xregex their meaning is intuitively clear.

Obviously, we can define various query classes by replacing the regular expressions in CRPQ by some more powerful language descriptors. However, this will not remedy the lack of a means to describe inter-path dependencies, and, furthermore, regular expressions seem powerful enough to describe the desired navigational features (in fact, the analyses carried out in [11, 12, 37] suggest that the regular expressions used in practical queries are even rather simple). What makes xregex interesting is that defining a variable on some edge and referencing it on another is a convenient way of describing inter-path dependencies, while, syntactically, staying in the realm of graph patterns with regular expressions.

To define such a query class, we first have to lift xregex to (multi-dimensional) *conjunctive xregex*, i. e., tuples $\bar{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_m)$ of xregex that can generate tuples $\bar{w} = (w_1, w_2, \dots, w_m)$ of words, but such that the w_i match the α_i in a *conjunctive* way with respect to the string variables (i. e., occurrences of the same variable x in different α_i and α_j must refer to the same string). Labeling graph patterns with xregex and interpreting the edge labels as conjunctive xregex yields our class of *conjunctive xregex path queries* (CXRPQs).

With respect to Objective 2 from above, we note that CXRPQs are purely pattern-based, i. e., the queries are just graph patterns with edge labels (as pointed out above, such queries are widely adapted in practice and we can therefore assume their user-friendliness), and the xregex used as path-specifications are a well-known practical tool (in fact, the use of string variables (backreferences) in regular expressions is a topic covered by standard textbooks on practical application of regular expressions (see, e. g., [29])). Regarding Objective 1, string variables add some mechanism to describe inter-path relationships and we already saw some illustrating examples. We will further argue in favour of their expressive power in the remainder of the introduction (see also Figure 5).

1.3 Barcelo et al.’s ECRPQ

An existing class of graph queries that is suitable for comparison with the class CXRPQ are the extended conjunctive regular path queries (ECRPQs) introduced in [8]. While CRPQs can be seen as graph patterns with unary regular relations on the edges (i. e., the regular expressions), the class ECRPQ permits regular relations of arbitrary arity, which allows to formulate a wide range of inter-path dependencies. ECRPQs have acceptable evaluation complexity: NL-complete in data-complexity and PSpace-complete in combined-complexity (see [8] for details).

In terms of expressive power, the class of CXRPQs completely cover the fragment ECRPQ^{er} of those ECRPQs that have only unary relations or equality relations (i. e., relations requiring certain paths to be equal). On the other hand, we can show that there are CXRPQ that are not expressible by ECRPQ^{er} . From

an intuitive point of view CXRPQ and ECRPQ are incomparable in the sense that ECRPQ can describe inter-path dependencies beyond simple equality, while CXRPQ can describe equality of arbitrarily many paths as, e. g., by $(x \vee y)^+$ in the query G_3 of Figure 2.

The class of ECRPQ is not purely pattern-based anymore, since also the relations must be given as regular expressions. In this regard, the authors of [8] mention that “[...] specifying regular relations with regular expressions is probably less natural than specifying regular languages (at least it is more cumbersome) [...]” and suggest that any practical standard would rather provide some reasonable regular relations as built-in predicates.

1.4 Technical Contributions

In terms of Objective 3, the best evaluation complexity that we can hope for is NL in data-complexity and NP in combined-complexity (since CRPQs have these lower bounds). Higher combined-complexity (e. g., PSpace) is still acceptable, as long as the optimum of NL in data-complexity is reached (this makes sense if we can assume our queries to be rather small in comparison to the data).

Unfortunately, CXRPQs have a surprisingly high evaluation complexity: even for the fixed xregex $\alpha_{ni} = \#x\{(a \vee b)^*\}(\#\#x)^*\#\#\#$, it is PSpace-hard to decide whether a given graph database contains a path labelled by a word from $\mathcal{L}(\alpha_{ni})$ (so Boolean evaluation is PSpace-hard in data complexity). This hardness result has nevertheless a silver lining: it directly points us to restrictions of CXRPQ that might lead to more tractable fragments. More precisely, for PSpace-hardness it seems vital that references for variable x are subject to the star-operator, and that the variable x can store words of unbounded length. Our main positive result will be that by restricting CXRPQs accordingly, we can tame their evaluation complexity and obtain more tractable fragments.

Let $\text{CXRPQ}^{\text{vsf}}$ be the class of *variable-star free* CXRPQs, i. e., no variable can be used under a star or plus (but normal symbols still can). For example, G_2, G_4 of Figure 2 are in $\text{CXRPQ}^{\text{vsf}}$. This restriction is enough to make the data-complexity drop from PSpace-hardness to the optimum of NL-completeness (although combined-complexity is ExpSpace). The upper bound is obtained by showing that $q \in \text{CXRPQ}^{\text{vsf}}$ can be transformed into equivalent q' in a certain *normal form*, which can be evaluated in nondeterministic space $O(|q'| \log(|\mathcal{D}|))$. However, $|q'| = O(2^{2^{|q|}})$ and we only get the single exponential space upper bound for combined-complexity by handling one exponential size blow-up with nondeterminism. A closer look at the normal form construction reveals that the exponential size blow-up is caused by chains of the following form: a reference of x occurs in the definition of y , a reference of y occurs in the definition of z , a reference of z occurs in the definition of u and so on (e. g., this happens with respect to x, y and z in G_4 of Figure 2). If we require that every variable definition only contains references of variables that themselves have definitions without variables, then we obtain the fragment $\text{CXRPQ}^{\text{vsf,fl}}$ which have normal forms of polynomial size and therefore the combined complexity drops to PSpace (i. e., the same complexity as ECRPQ).

The second successful approach is to add a constant upper-bound k on the *image size* of CXRPQs, i. e., the length of the words stored in variables. Let $\text{CXRPQ}^{\leq k}$ be the corresponding fragments. For every $\text{CXRPQ}^{\leq k}$, the evaluation complexity drops to the optimum of NL-completeness in data-complexity

and NP-completeness in combined-complexity (i. e., the same as for CRPQs). Unlike for CXPQ^{vsf} , bounding the image size does not impose any syntactical restrictions; it is rather a restriction of how a CXPQ can match a graph database. For example, we can treat G_3 of Figure 2 as a $\text{CXPQ}^{\leq 10}$, i. e., we only have a successful match if the paths between v_1 and v_2 are of length at most 10 (note that the paths from v_1 and v_2 to their mutual friend can have unbounded size). For G_1 of Figure 2, on the other hand, the image size of variables is necessarily bounded by 1 and therefore it does not matter whether we interpret it as CXPQ or $\text{CXPQ}^{\leq k}$ with $k \geq 1$. We stress the fact that this bound only applies to strings stored in variables; we can still specify paths of arbitrary length with regular expressions, i. e., $\text{CRPQ} \subseteq \text{CXPQ}^{\leq k}$. Moreover, evaluating $\text{CXPQ}^{\leq k}$ is not as simple as just replacing all variables by fixed words of length at most k and then evaluating a CRPQ, since we also have to take care of dependencies between variable definitions.

Finally, there are two more noteworthy results about $\text{CXPQ}^{\leq k}$ (a negative and a positive one). While in combined-complexity CRPQ can be evaluated in polynomial-time if the underlying graph pattern is acyclic (see [10, 8, 6]), $\text{CXPQ}^{\leq k}$ remain NP-hard in combined-complexity even for single-edge graph patterns (and $k = 1$). This also demonstrates the general difference of $\text{CXPQ}^{\leq k}$ and CRPQ. On the positive side, if instead of a constant upper bound, we allow images bounded logarithmically in the size of the database, then the NP upper bound in combined-complexity remains, while the data-complexity slightly increases to $O(\log^2(|\mathcal{D}|))$.

The question is whether these fragments are still interesting with respect to Objectives 1 and 2. We believe the answer is yes. First observe that the restrictions are quite natural: Not using the star-operator over variables is a rule not difficult for users to comply with, if they are familiar with regular expressions; it is also easily to be checked algorithmically, and the same holds for the additional restriction required by $\text{CXPQ}^{\text{vsf,fl}}$. The class $\text{CXPQ}^{\leq k}$ does not require any syntactical restriction; when interpreting the query result, the user only has to keep in mind that the paths corresponding to images of variables are bounded in length.

Regarding expressive power, all these fragments contain non-trivial examples of CXPQs. With respect to the examples from Figure 2, $G_4 \in \text{CXPQ}^{\text{vsf}}$ and $G_2 \in \text{CXPQ}^{\text{vsf,fl}}$; any CXPQ can be interpreted as $\text{CXPQ}^{\leq k}$ for any k . Both CXPQ^{vsf} and $\text{CXPQ}^{\text{vsf,fl}}$ still cover the fragment ECRPQ^{er} . It is tempting to misinterpret queries from $\text{CXPQ}^{\leq k}$ as CRPQ with mere syntactic sugar (since string variables range over finite sets of words). However, it can be proven that even $\text{CXPQ}^{\leq 1}$ contains queries that are not expressible by CRPQs.

2 Preliminaries

Let $\mathbb{N} = \{1, 2, 3, \dots\}$ and $[n] = \{1, 2, \dots, n\}$ for $n \in \mathbb{N}$. A^+ denotes the set of non-empty words over an alphabet A and $A^* = A^+ \cup \{\varepsilon\}$ (where ε is the empty word). For a word $w \in A^*$, $|w|$ denotes its length, and for $k \in \mathbb{N}$, $A^{\leq k} = \{w \in A^* \mid |w| \leq k\}$. For $w_1, w_2, \dots, w_n \in A^*$, we set $\prod_{i=1}^n w_i = w_1 w_2 \dots w_n$, and if $w = w_i$, for every $i \in [n]$, then we also write w^n instead of $\prod_{i=1}^n w_i$. For some alphabet A , a word $w \in A^*$ and a $b \in A$, $|w|_b$ is the number of occurrences of symbol b in w .

We fix a finite *terminal* alphabet Σ and an enumerable set \mathcal{X}_s of *string variables*, where $\mathcal{X}_s \cap \Sigma = \emptyset$. As a convention, we use symbols $\mathbf{a}, \mathbf{b}, \mathbf{c}, \mathbf{d}, \dots$ for elements from Σ , and $x, y, z, x_1, x_2, \dots, y_1, y_2, \dots$ for variables from \mathcal{X}_s . We consistently use sans-serif font for string variables to distinguish them from *node variables* to be introduced later. We use *regular expressions* and (*nondeterministic*) *finite automata* (NFA for short) as commonly defined in the literature (see Section 3 and the remainder of this section for more details).

2.1 Ref-Words

The following *ref-words* (first introduced in [40]) are convenient for defining the semantics of xregex (Section 3). They have also been used in [28] and for so-called document spanners in [27, 26, 20]. Ref-words will be vital in our definition of conjunctive xregex (Section 3.1), which are the basis of the class CXRPQ.

For every $x \in \mathcal{X}_s$, we use the pair of symbols $\mathbf{x}\triangleright$ and $\triangleleft\mathbf{x}$, which are interpreted as opening and closing parentheses.

Definition 1 (Ref-Words). *A subword-marked word (over terminal alphabet Σ and variables \mathcal{X}_s) is a word $w \in (\Sigma \cup \{\mathbf{x}\triangleright, \triangleleft\mathbf{x} \mid x \in \mathcal{X}_s\} \cup \mathcal{X}_s)^*$ that, for every $x \in \mathcal{X}_s$, contains the parentheses $\mathbf{x}\triangleright$ and $\triangleleft\mathbf{x}$ at most once and all these parentheses form a well-formed parenthesised expression. For every $x \in \mathcal{X}_s$, a subword $\mathbf{x}\triangleright v \triangleleft\mathbf{x}$ in w is called a definition (of variable x), and an occurrence of symbol x is called a reference (of variable x). For a subword-marked word w over Σ and \mathcal{X}_s , the binary relation \preceq_w over \mathcal{X}_s is defined by setting $x \preceq_w y$ if in w there is a definition of y that contains a reference or a definition of x . A ref-word (over terminal alphabet Σ and variables \mathcal{X}_s) is a subword-marked word over Σ and \mathcal{X}_s , such that the transitive closure of \preceq_w is acyclic.*

Ref-words are just words over alphabet $\Sigma \cup \mathcal{X}_s$, in which some subwords are uniquely marked by means of the parentheses $\mathbf{x}\triangleright$ and $\triangleleft\mathbf{x}$. Moreover, the marked subwords are not allowed to overlap, i. e., $\mathbf{x}\triangleright \mathbf{y}\triangleright \triangleleft\mathbf{x} \triangleleft\mathbf{y}$ must not occur as subsequence. For every variable $x \in \mathcal{X}_s$, all occurrences of x in a ref-word are interpreted as references to the definition of x . Since definitions may contain itself references or definitions of other variables, there are chains of references, e. g., the definition of x contains references of y , but the definition of y contains references of z and so on. Therefore, in order to make this implicit referencing process terminate, we have to require that it is acyclic, which is done by requiring the transitive closure of \preceq_w to be acyclic. For example, $\mathbf{a}\mathbf{x}\mathbf{b}\mathbf{x}\triangleright \mathbf{a}\mathbf{b}\triangleleft\mathbf{x}\mathbf{c}\mathbf{y}\triangleright \mathbf{x}\mathbf{a}\mathbf{a}\triangleleft\mathbf{y}\mathbf{y}$ is a valid ref-word, while $\mathbf{a}\mathbf{x}\mathbf{b}\mathbf{x}\triangleright \mathbf{a}\mathbf{b}\triangleleft\mathbf{x}\mathbf{c}\mathbf{y}\triangleright \mathbf{x}\mathbf{a}\mathbf{y}\triangleleft\mathbf{y}\mathbf{y}$, or $\mathbf{a}\mathbf{x}\mathbf{a}\mathbf{x}\triangleright \mathbf{a}\mathbf{y}\mathbf{b}\triangleleft\mathbf{x}\mathbf{c}\mathbf{y}\triangleright \mathbf{x}\mathbf{a}\triangleleft\mathbf{y}\mathbf{y}$ are not.

For ref-words over Σ and \mathcal{X}_s , it is therefore possible to successively substitute references by their definitions until we obtain a word over Σ . This can be formalised as follows.

Definition 2 (Deref-Function). *For a ref-word w over Σ and \mathcal{X}_s , $\text{deref}(w) \in \Sigma^*$ is obtained from w by the following procedure:*

1. Delete all occurrences of $x \in \mathcal{X}_s$ without definition in w .
2. Repeat until we have obtained a word over Σ :
 - (a) Let $\mathbf{x}\triangleright v_x \triangleleft\mathbf{x}$ be a definition such that $v_x \in \Sigma^*$.

(b) Replace $\succ v_x \prec^x$ and all occurrences of x in w by v_x .

Proposition 1. *The function $\text{deref}(w)$ is well-defined.*

Proof. If we reach Step 2 without any definition for a variable, then there can also be no reference of a variable, since we deleted all variable references that have no definition in Step 1, and for every variable definition deleted in Step 2b, we also deleted all corresponding references. Consequently, if we reach Step 2 without variable definitions, then we have obtained a word over Σ and terminate. If, on the other hand, there is at least one definition for some variable x when we reach Step 2, then, by the definition of a ref-word, there must be at least one definition that does not contain another definition. Hence, there must be at least one definition $\succ v_x \prec^x$ such that $v_x \in \Sigma^*$, as required by Step 2a. In Step 2b this definition of x and all its references are then replaced by v_x , which is a word over Σ . Moreover, it can be easily verified that $\text{deref}(w)$ does not depend on the actual choices of the definitions that are replaced in the iterations of Step 2a. \square

For a ref-word w , the procedure of Definition 2 that computes $\text{deref}(w)$ also uniquely allocates a subword v_x of $\text{deref}(w)$ to each variable x that has a definition in w (i.e., the subwords v_x defined in the iterations of Step 2a). In this way, a ref-word w over terminal alphabet Σ and variables \mathcal{X}_s describes a *variable mapping* $\text{vmap}_{w, \mathcal{X}_s} : \mathcal{X}_s \rightarrow \Sigma^*$, i.e., we set $\text{vmap}_{w, \mathcal{X}_s}(x) = v_x$ if x has a definition in w and we set $\text{vmap}_{w, \mathcal{X}_s}(x) = \varepsilon$ otherwise. The elements $\text{vmap}_{w, \mathcal{X}_s}(x)$ for $x \in \mathcal{X}_s$ are called *variable images*.

Note that even if x has a definition in w , $\text{vmap}_{w, \mathcal{X}_s}(x) = \varepsilon$ is possible due to a definition $\succ \prec^x$ or to a definition $\succ v_x \prec^x$, where v_x is a non empty word with only references of variables y with $\text{vmap}_{w, \mathcal{X}_s}(y) = \varepsilon$. Also note that any ref-word w over terminal alphabet Σ and variables \mathcal{X}_s is also a ref-word over any terminal alphabet $\Sigma' \supset \Sigma$ and variables $\mathcal{X}'_s \supset \mathcal{X}_s$ ($\text{vmap}_{w, \mathcal{X}'_s}$ then equals the extension of $\text{vmap}_{w, \mathcal{X}_s}$ by $\text{vmap}_{w, \mathcal{X}'_s}(x') = \varepsilon$ for every $x' \in \mathcal{X}'_s \setminus \mathcal{X}_s$). If the set of variables \mathcal{X}_s is clear from the context or negligible, we shall also denote the variable mapping by vmap_w and if there is some obvious implicit order on \mathcal{X}_s , e.g., given by indices as in the case $\mathcal{X}_s = \{x_1, x_2, \dots, x_m\}$, then we also write vmap_w as a tuple $(\text{vmap}_w(x_1), \text{vmap}_w(x_2), \dots, \text{vmap}_w(x_m))$.

A set L of ref-words is called *ref-language* and we extend the deref -function to ref-languages in the obvious way, i.e., $\text{deref}(L) = \{\text{deref}(w) \mid w \in L\}$. Let us illustrate ref-words with an example.

Example 1. *Let $\Sigma = \{a, b, c\}$ and let $x_1, x_2, x_3, x_4 \in \mathcal{X}_s$.*

$$w = \mathbf{ax_4a} \succ^{x_1} \mathbf{ab} \succ^{x_2} \mathbf{acc} \prec^{x_2} \mathbf{ax_2x_4} \prec^{x_1} \succ^{x_3} \mathbf{x_1ax_2} \prec^{x_3} \mathbf{x_3bx_1}.$$

The procedure of Definition 2 will first delete all occurrences of x_4 . Then $\succ^{x_2} \mathbf{acc} \prec^{x_2}$ and all references of x_2 are replaced by \mathbf{acc} . After this step, the definition for variable x_1 is $\succ^{x_1} \mathbf{abaccaacc} \prec^{x_1}$, so $\mathbf{abaccaacc}$ can be substituted for the definitions and references of x_1 . After replacing the last variable x_3 , we obtain $\text{deref}(w)$. Note that $\text{vmap}_w = (\mathbf{abaccaacc}, \mathbf{acc}, \mathbf{abaccaaccaacc}, \varepsilon)$.

2.2 Graph-Databases

A *graph-database* (over Σ) is a directed, edge labelled multigraph $\mathcal{D} = (V_{\mathcal{D}}, E_{\mathcal{D}})$, where $V_{\mathcal{D}}$ is the set of *vertices* (or *nodes*) and $E_{\mathcal{D}} \subseteq V_{\mathcal{D}} \times \Sigma \times V_{\mathcal{D}}$ is the set

of *edges* (or *arcs*). A path from $u \in V_{\mathcal{D}}$ to $v \in V_{\mathcal{D}}$ of length $k \geq 0$ is a sequence $p = (w_0, a_1, w_1, a_2, w_2, \dots, w_{k-1}, a_k, w_k)$ with $(w_{i-1}, a_i, w_i) \in E_{\mathcal{D}}$ for every $i \in [k]$. We say that p is *labelled* with the *word* $a_1 a_2 \dots a_k \in \Sigma^*$. According to this definition, for every $v \in V_{\mathcal{D}}$, (v) is a path from v to v of length 0 that is labelled by ε . Hence, every node of every graph-database has an ε -labelled path to itself (and these are the only ε -labelled paths in \mathcal{D}).

Nondeterministic finite automata (NFAs) are just graph databases, the nodes of which are called *states*, and that have a specified *start state* and a set of specified *final states*. Moreover, we allow the empty word as edge label as well (which is not the case for graph databases). The language $\mathcal{L}(M)$ of an NFA M is the set of all labels from paths from the start state to some final state.

In the following, \mathcal{X}_n is an enumerable set of *node-variables*; we shall use symbols $x, y, z, x_1, x_2, \dots, y_1, y_2, \dots$ for node variables (in contrast to the string variables \mathcal{X}_s in sans-serif font).

2.3 Conjunctive Path Queries

Let \mathfrak{R} be a class of language descriptors, and, for every $r \in \mathfrak{R}$, let $\mathcal{L}(r)$ denote the language represented by r . An \mathfrak{R} -*graph pattern* is a directed, edge-labelled graph $G = (V, E)$ with $V \subseteq \mathcal{X}_n$ and $E \subseteq V \times \mathfrak{R} \times V$; it is an \mathfrak{R} -graph pattern over alphabet Σ , if $\mathcal{L}(\alpha) \subseteq \Sigma^*$ for every $(x, \alpha, y) \in E$. For an \mathfrak{R} -graph pattern $G = (V, E)$ over Σ and a graph-database $\mathcal{D} = (V_{\mathcal{D}}, E_{\mathcal{D}})$ over Σ , a mapping $h : V \rightarrow V_{\mathcal{D}}$ is a *matching morphism* for G and \mathcal{D} if, for every $e = (x, \alpha, y) \in E$, \mathcal{D} contains a path from $h(x)$ to $h(y)$ that is labelled with a word $w_e \in \mathcal{L}(\alpha)$. The tuple $(w_e)_{e \in E}$ is a tuple of *matching words* (with respect to h). In particular, a matching morphism can have several different tuples of matching words.

A *conjunctive \mathfrak{R} -path query* (\mathfrak{R} -CPQ for short) is a query $q = \bar{z} \leftarrow G_q$, where $G_q = (V_q, E_q)$ is an \mathfrak{R} -graph pattern and $\bar{z} = (z_1, \dots, z_\ell)$ with $\{z_1, z_2, \dots, z_\ell\} \subseteq V_q$. We say that q is an \mathfrak{R} -CPQ over alphabet Σ if G_q is an \mathfrak{R} -graph pattern over Σ . The query q is a *single-edge* query, if $|E_q| = 1$.

For an \mathfrak{R} -CPQ $q = \bar{z} \leftarrow G_q$ over Σ with $\bar{z} = (z_1, z_2, \dots, z_\ell)$, a graph-database $\mathcal{D} = (V_{\mathcal{D}}, E_{\mathcal{D}})$ over Σ and a matching morphism h for G_q and \mathcal{D} (we also call h a *matching morphism for q and \mathcal{D}*), we define $q_h(\mathcal{D}) = (h(z_1), h(z_2), \dots, h(z_\ell))$ and we set $q(\mathcal{D}) = \{q_h(\mathcal{D}) \mid h \text{ is a matching morphism for } q \text{ and } \mathcal{D}\}$. The mapping $\mathcal{D} \mapsto q(\mathcal{D})$ from the set of graph-databases to the set of relations over Σ of arity ℓ that is defined by q shall be denoted by $\llbracket q \rrbracket$, and for any class A of conjunctive path queries, we set $\llbracket A \rrbracket = \{\llbracket q \rrbracket \mid q \in A\}$.

A *Boolean \mathfrak{R} -CPQ* has the form $\bar{z} \leftarrow G_q$ where \bar{z} is the empty tuple. In this case, we also denote q just by G_q instead of $(\) \leftarrow G_q$. For a Boolean \mathfrak{R} -CPQ q and a graph-database \mathcal{D} , we either have $q(\mathcal{D}) = \{(\)\}$ or $q(\mathcal{D}) = \emptyset$, which we shall also denote by $\mathcal{D} \models q$ and $\mathcal{D} \not\models q$, respectively. For Boolean queries q , the set $\llbracket q \rrbracket$ can also be interpreted as $\{\mathcal{D} \mid \mathcal{D} \models q\}$. Two \mathfrak{R} -CPQs q and q' are *equivalent*, denoted by $q \equiv q'$, if $\llbracket q \rrbracket = \llbracket q' \rrbracket$, i. e., $q(\mathcal{D}) = q'(\mathcal{D})$ for every graph-database \mathcal{D} (or, in the Boolean case, $\mathcal{D} \models q \Leftrightarrow \mathcal{D} \models q'$).

For a class Q of conjunctive path queries, Q -BOOL-EVAL is the problem to decide, for a given Boolean $q \in Q$ and a graph database \mathcal{D} , whether $\mathcal{D} \models q$. By Q -CHECK, we denote the problem to check $\bar{t} \in q(\mathcal{D})$ for a given $q \in Q$, a graph database \mathcal{D} and a tuple \bar{t} .

As common in database theory, the *combined-complexity* for an algorithm solving Q -BOOL-EVAL or Q -CHECK is the time or space needed by the algorithm

measured in both $|q|$ and $|\mathcal{D}|$, while for the *data-complexity* the query q is considered constant. For simplicity, we assume $|q| = O(|\mathcal{D}|)$ throughout the paper.

Conjunctive regular path queries (CRPQ) are \mathfrak{R} -CPQ where \mathfrak{R} is the class of regular expressions (which are defined in Section 3). For some of our results, we need the following result about CRPQs.

Lemma 1 ([8]). *CRPQ-BOOL-EVAL is NP-complete in combined complexity and NL-complete in data-complexity.*

3 Xregex

We next define the underlying language descriptors for our class of conjunctive path queries. Complete formalisations of the class of xregex can also be found elsewhere in the literature (e.g., [40, 28, 25]). However, since we will extend xregex to conjunctive xregex (the main building block for CXRPQ), we give a full definition.

Definition 3 (Xregex). *The set $\text{XRE}_{\Sigma, \mathcal{X}_s}$ of regular expressions with backreferences (over Σ and \mathcal{X}_s), also denoted by xregex, for short, is recursively defined as follows:*

1. $a \in \text{XRE}_{\Sigma, \mathcal{X}_s}$ and $\text{var}(a) = \emptyset$, for every $a \in \Sigma \cup \{\varepsilon\}$,
2. $x \in \text{XRE}_{\Sigma, \mathcal{X}_s}$ and $\text{var}(x) = \{x\}$, for every $x \in \mathcal{X}_s$,
3. $(\alpha \cdot \beta) \in \text{XRE}_{\Sigma, \mathcal{X}_s}$, $(\alpha \vee \beta) \in \text{XRE}_{\Sigma, \mathcal{X}_s}$, and $(\alpha)^+ \in \text{XRE}_{\Sigma, \mathcal{X}_s}$, for every $\alpha, \beta \in \text{XRE}_{\Sigma, \mathcal{X}_s}$;
furthermore, $\text{var}((\alpha \cdot \beta)) = \text{var}((\alpha \vee \beta)) = \text{var}(\alpha) \cup \text{var}(\beta)$ and $\text{var}((\alpha)^+) = \text{var}(\alpha)$,
4. $x\{\alpha\} \in \text{XRE}_{\Sigma, \mathcal{X}_s}$ and $\text{var}(x\{\alpha\}) = \text{var}(\alpha) \cup \{x\}$, for every $\alpha \in \text{XRE}_{\Sigma, \mathcal{X}_s}$ and $x \in \mathcal{X}_s \setminus \text{var}(\alpha)$.

For technical reasons, we also add \emptyset to $\text{XRE}_{\Sigma, \mathcal{X}_s}$. For $\alpha \in \text{XRE}_{\Sigma, \mathcal{X}_s}$, we use r^* as a shorthand form for $r^+ \vee \varepsilon$, and we usually omit the operator ‘ \cdot ’, i.e., we use juxtaposition. If this does not cause ambiguities, we often omit parenthesis. For example, x , $x\{ya\}$ and $x\{(y\{z\{a^* \vee bc\}a\}y)^+b\}x$ are all valid xregex, while neither $x\{ax\}b$ nor $x\{ax\{b^*\}a\}b$ is a valid xregex.

We call an occurrence of $x \in \mathcal{X}_s$ a *reference of variable* x and a subexpression $x\{\alpha\}$ a *definition of variable* x . The set $\text{XRE}_{\Sigma, \emptyset}$ is exactly the set of regular expressions over Σ , which shall be denoted by RE_{Σ} in the following. We also use the term *classical* regular expressions for a clearer distinction from xregex. If the underlying alphabet Σ or set \mathcal{X}_s of variables is negligible or clear from the context, we also drop these and simply write XRE and RE.

We next define the semantics of xregex, for which we heavily rely on the concept of ref-words (see Section 2). First, for any $\alpha \in \text{RE}$, let $\mathcal{L}(\alpha)$ be the regular language described by the (classical) regular expression α defined as usual: $\mathcal{L}(a) = \{a\}$, $\mathcal{L}(\alpha \cdot \beta) = \mathcal{L}(\alpha) \cdot \mathcal{L}(\beta)$, $\mathcal{L}(\alpha \vee \beta) = \mathcal{L}(\alpha) \cup \mathcal{L}(\beta)$, $\mathcal{L}(\alpha^+) = \mathcal{L}(\alpha)^+$. Now in order to define the language described by an xregex $\alpha \in \text{XRE}_{\Sigma, \mathcal{X}_s}$, we first define a classical regular expression α_{ref} over the alphabet

$\Sigma \cup \mathcal{X}_s \cup \{x_{\triangleright}, x_{\triangleleft} \mid x \in \mathcal{X}_s\}$ that is obtained from α by iteratively replacing all variable definitions $x\{\beta\}$ by $x_{\triangleright}\beta x_{\triangleleft}$. For example,

$$\begin{aligned}\alpha &= x\{(y\{z\{a^* \vee bc\}a\}y)^+b\}x, \\ \alpha_{\text{ref}} &= x_{\triangleright}(y_{\triangleright}z_{\triangleright}(a^* \vee bc)x_{\triangleleft}^z a x_{\triangleleft}^y y)^+b x_{\triangleleft}.\end{aligned}$$

We say that α is *sequential* if every $w \in \mathcal{L}(\alpha_{\text{ref}})$ contains for every $x \in \mathcal{X}_s$ at most one occurrence of x_{\triangleright} (if not explicitly stated otherwise, all our xregex are sequential). If an xregex α is sequential, then every $w \in \mathcal{L}(\alpha_{\text{ref}})$ must be a ref-word. Indeed, this directly follows from the fact that the definitions $x\{\dots\}$ are always subexpressions. Hence, for sequential xregex, $\mathcal{L}(\alpha_{\text{ref}})$ is a ref-language, which we shall denote by $\mathcal{L}_{\text{ref}}(\alpha)$. If a ref-word $v \in \mathcal{L}_{\text{ref}}(\alpha)$ contains a definition $x_{\triangleright}v_x x_{\triangleleft}$, then we say that the corresponding definition $x\{\gamma_x\}$ in α is *instantiated* (by v). In particular, we observe that sequential xregex can nevertheless have several definitions for the same variable x , but at most one of them is instantiated by any ref-word. Finally, the language described by α is defined as $\mathcal{L}(\alpha) = \text{deref}(\mathcal{L}_{\text{ref}}(\alpha))$. As a special case, we also define $\mathcal{L}(\emptyset) = \emptyset$. For $\alpha \in \text{XRE}_{\Sigma, \mathcal{X}_s}$ and $w \in \Sigma^*$, we say that w *matches* α with witness $u \in \mathcal{L}_{\text{ref}}(\alpha)$ and variable mapping vmap_u , if $\text{deref}(u) = w$.

Example 2. Let $\alpha \in \text{XRE}_{\{a,b\}, \mathcal{X}_s}$ with $x_1, x_2 \in \mathcal{X}_s$:

$$\begin{aligned}\alpha &= a^* x_1 \{a^* x_2 \{(a \vee b)^*\} b^* a^*\} x_2^* (a \vee b)^* x_1, \\ \alpha_{\text{ref}} &= a^* x_{1\triangleright} a^* x_{2\triangleright} (a \vee b)^* x_{2\triangleleft} b^* a^* x_{1\triangleleft} x_2^* (a \vee b)^* x_1, \\ u_1 &= aaaa x_{1\triangleright} x_{2\triangleright} ba x_{2\triangleleft} baa x_{1\triangleleft} x_2 x_2 bba x_1 \in \mathcal{L}_{\text{ref}}(\alpha), \\ u_2 &= aaa x_{1\triangleright} a x_{2\triangleright} bab x_{2\triangleleft} aa x_{1\triangleleft} x_2 abb x_1 \in \mathcal{L}_{\text{ref}}(\alpha), \\ w_\alpha &= \text{deref}(u_1) = \text{deref}(u_2) = a^4 (ba)^2 (ab)^3 (ba)^3 a \in \mathcal{L}(\alpha), \\ \text{vmap}_{u_1} &= (babaa, ba), \text{vmap}_{u_2} = (ababaa, bab).\end{aligned}$$

For $\gamma = x_1\{c^*(x_2\{a^*\} \vee x_3\{b^*\})\}cx_2cx_3bx_1$, $c^2a^2ca^2cbc^2a^2 \in \mathcal{L}(\gamma)$ is witnessed by ref-word $u_5 = x_{1\triangleright}c^2x_{2\triangleright}a^2x_{3\triangleright}c^2x_{2\triangleleft}a^2x_{3\triangleleft}cx_2cx_3bx_1$, which has the variable mapping $\text{vmap}_{u_5} = (c^2a^2, a^2, \varepsilon)$. Note that empty variable images can result from variables without definitions in the ref-word (as in this example), but also from definitions of variables that define the empty word, as, e.g., in the ref-word $x_{1\triangleright}x_{1\triangleleft}cx_1 \in \mathcal{L}_{\text{ref}}(x_1\{(a \vee b)^*\}cx_1)$.

For $\alpha \in \text{XRE}_{\Sigma, \mathcal{X}_s}$, we define the relation \preceq_α analogously how it is done for ref-words, i.e., $x \preceq_\alpha y$ if in α there is a definition of y that contains a reference or a definition of x . Even though the transitive closure of \preceq_v is acyclic for every $v \in \mathcal{L}_{\text{ref}}(\alpha)$, the transitive closure of \preceq_α is not necessarily acyclic. For example, $\alpha = x\{a^*\}y\{x\} \vee y\{a^*\}x\{y\}$ is an xregex, but the transitive closure of \preceq_α is not acyclic. We call xregex *acyclic* if \preceq_α is acyclic.

3.1 Conjunctive Xregex

Syntactically, conjunctive xregex are tuples of xregex. Their semantics, however, is more difficult and we spend some more time with intuitive explanations before giving a formal definition.

Definition 4 (Conjunctive Xregex). A tuple $\bar{\alpha} = (\alpha_1, \dots, \alpha_m) \in (\text{XRE}_{\Sigma, \mathcal{X}_s})^m$ is a conjunctive xregex of dimension m , if $\alpha_1\alpha_2 \dots \alpha_m$ is an acyclic xregex.

By $m\text{-CXRE}_{\Sigma, \mathcal{X}_s}$, we denote the set of conjunctive xregex of dimension m (over Σ and \mathcal{X}_s) and we set $\text{CXRE}_{\Sigma, \mathcal{X}_s} = \bigcup_{m \geq 1} m\text{-CXRE}_{\Sigma, \mathcal{X}_s}$. Note that $1\text{-CXRE}_{\Sigma, \mathcal{X}_s} = \text{XRE}_{\Sigma, \mathcal{X}_s}$. If the terminal alphabet Σ or set \mathcal{X}_s of variables are negligible or clear from the context, we also drop the corresponding subscripts. We also write $\bar{\alpha}[i]$ to denote the i^{th} element of some $\bar{\alpha} \in m\text{-CXRE}$.

Before defining next the semantics of conjunctive xregex, we first give some intuitive explanations. The central idea of a conjunctive xregex $\bar{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_m)$ is that the definition of some x in some α_i also serves as definition for possible references of x in some α_j with $i \neq j$ (which, due to the requirement that $\alpha_1 \alpha_2 \dots \alpha_m$ is an xregex, cannot contain a definition of x). Let us illustrate the situation with the concrete example $\bar{\gamma} = (\gamma_1, \gamma_2)$, with $\gamma_1 = (x\{\mathbf{a}^*\} \vee \mathbf{b}^*)y$ and $\gamma_2 = y\{\mathbf{xaxb}\}\mathbf{b}y^*$.

As an intermediate step, we first add dummy definitions for the undefined variables, i. e., we define $\langle \gamma_1 \rangle_{\text{int}} = y\{\Sigma^*\}\#\gamma_1$ and $\langle \gamma_2 \rangle_{\text{int}} = x\{\Sigma^*\}\#\gamma_2$, which are both xregex in which all variables have a definition. Now, we can treat $(\langle \gamma_1 \rangle_{\text{int}}, \langle \gamma_2 \rangle_{\text{int}})$ as a generator for pairs of ref-words (and therefore as pairs of words over Σ), but we only consider those pairs of ref-words that have the same variable mapping. For example, $u_1 = {}^y \triangleright \mathbf{a}^5 \mathbf{b} \triangleleft^y \# {}^x \triangleright \mathbf{aa} \triangleleft^x y \in \mathcal{L}_{\text{ref}}(\langle \gamma_1 \rangle_{\text{int}})$ and $u_2 = {}^x \triangleright \mathbf{aa} \triangleleft^x \# {}^y \triangleright \mathbf{xaxb} \triangleleft^y \mathbf{byy} \in \mathcal{L}_{\text{ref}}(\langle \gamma_2 \rangle_{\text{int}})$ and, furthermore, since u_1 and u_2 have the same variable mapping $(\mathbf{aa}, \mathbf{a}^5 \mathbf{b})$, they are witnesses for the conjunctive match $(w_1, w_2) = (\mathbf{aaa}^5 \mathbf{b}, \mathbf{a}^5 \mathbf{bb}(\mathbf{a}^5 \mathbf{b})^2)$ of $\bar{\gamma}$, since $\text{deref}(u_1) = \mathbf{a}^5 \mathbf{b} \# w_1$ and $\text{deref}(u_2) = \mathbf{aa} \# w_2$. On the other hand, $v_1 = {}^y \triangleright \mathbf{a} \triangleleft^y \# {}^x \triangleright \mathbf{a} \triangleleft^x y$ and $v_2 = {}^x \triangleright \mathbf{a} \triangleleft^x \# {}^y \triangleright \mathbf{xaxb} \triangleleft^y \mathbf{by}$ are also ref-words from $\mathcal{L}_{\text{ref}}(\langle \alpha_1 \rangle_{\text{int}})$ and $\mathcal{L}_{\text{ref}}(\langle \alpha_2 \rangle_{\text{int}})$, respectively, but, even though $\text{deref}(v_1) = \mathbf{a} \# \mathbf{aa}$ and $\text{deref}(v_2) = \mathbf{a} \# \mathbf{a}^3 \mathbf{bba}^3 \mathbf{b}$, the tuple $(\mathbf{aa}, \mathbf{a}^3 \mathbf{bba}^3 \mathbf{b})$ is not a conjunctive match for $\bar{\alpha}$, since $\text{vmap}_{v_1}(y) = \mathbf{a} \neq \mathbf{a}^3 \mathbf{b} = \text{vmap}_{v_2}(y)$. We shall now generalise this idea to obtain a sound definition of the semantics of conjunctive xregex.

For any xregex $\gamma \in \text{XRE}_{\Sigma, \mathcal{X}_s}$, let $\nabla_\gamma = \prod_{x \in A} x\{\Sigma^*\}$, where $A \subseteq \mathcal{X}_s$ is the set of variables that have no definition in γ (the order of the definitions $x\{\Sigma^*\}$ with $x \in A$ in ∇_γ shall be irrelevant), and we also define $\langle \gamma \rangle_{\text{int}} = \nabla_\gamma \# \gamma$, where $\#$ is a new symbol with $\# \notin \Sigma$. Finally, a tuple $\bar{w} = (w_1, w_2, \dots, w_m) \in (\Sigma^*)^m$ is a (*conjunctive*) *match* for $\bar{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_m) \in (\text{XRE}_{\Sigma, \mathcal{X}_s})^m$ with variable mapping ψ if, for every $i \in [m]$, there is a ref-word $v_i = v'_i \# v''_i \in \mathcal{L}_{\text{ref}}(\langle \alpha_i \rangle_{\text{int}})$ with $\text{deref}(v_i) = \text{deref}(v'_i) \# w_i$ and $\text{vmap}_{v_i} = \psi$.

Conjunctive xregex behave similar to xregex in terms of how ε can be allocated to a variable. It might happen that α_i contains the definition of some variable x (and therefore ∇_{α_i} does not contain $x\{\Sigma^*\}$), but the corresponding ref-word $v_i \in \mathcal{L}_{\text{ref}}(\langle \alpha_i \rangle_{\text{int}})$ does not contain a definition of x . This means that $\text{vmap}_{v_i}(x) = \varepsilon$ and therefore all other v_j with $i \neq j$ must also allocate ε to x in their definitions $x\{\Sigma^*\}$ contained in ∇_{α_j} . On the other hand, it is possible that $\text{vmap}_{v_i}(x) = \varepsilon$ even though v_i contains a definition of x .

For an $\bar{\alpha} \in \text{CXRE}$, $\mathcal{L}(\bar{\alpha})$ is the *set of conjunctive matches* for $\bar{\alpha}$. We say that conjunctive xregex $\bar{\alpha}$ and $\bar{\beta}$ are *equivalent* if $\mathcal{L}(\bar{\alpha}) = \mathcal{L}(\bar{\beta})$ (note that $\bar{\alpha}$ and $\bar{\beta}$ having the same dimension is necessary for this). We stress the fact that the order of the elements α_i of a conjunctive xregex $\bar{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_m)$ is vital.

A special class of conjunctive xregex is $m\text{-CXRE}_{\Sigma, \emptyset}$, i. e., the class of all m -dimensional tuples of classical regular expressions. For every $\bar{\alpha} \in m\text{-CXRE}_{\Sigma, \emptyset}$, $\mathcal{L}(\bar{\alpha}) = \mathcal{L}(\bar{\alpha}[1]) \times \mathcal{L}(\bar{\alpha}[2]) \times \dots \times \mathcal{L}(\bar{\alpha}[m])$.

Example 3. Consider the xregex

$$\begin{aligned}\alpha_1 &= x_2\{x_1 \vee a^*\}b, & \alpha_2 &= x_1\{(a \vee b)^*\}x_3\{c^*\}bx_3, \\ \alpha_3 &= x_2^*a^*x_1, & \alpha_4 &= x_4\{a^*\}bx_4x_1\{x_2a\}.\end{aligned}$$

Then (α_2, α_4) is not a conjunctive xregex since $\alpha_2\alpha_4$ is not sequential, but both (α_3, α_4) and $(\alpha_1, \alpha_2, \alpha_3)$ are conjunctive xregex.

Obviously, $w_1 = \mathbf{aab}$, $w_2 = \mathbf{bbacbc}$ and $w_3 = \mathbf{aa}$ are matches for α_1 , α_2 and α_3 , respectively, with variable mappings $(\varepsilon, \mathbf{aa}, \varepsilon)$, $(\mathbf{bba}, \varepsilon, \mathbf{c})$ and $(\varepsilon, \varepsilon, \varepsilon)$, respectively. However, for (w_1, w_2, w_3) to be a conjunctive match for $(\alpha_1, \alpha_2, \alpha_3)$, there must be words u_1, u_2, u_3 such that the following w'_i match the $\langle \alpha_i \rangle_{\text{int}}$ all with the same variable mapping (u_1, u_2, u_3) , which can be easily seen to be impossible.

$$\begin{aligned}w'_1 &= u_1u_3\#w_1 & \langle \alpha_1 \rangle_{\text{int}} &= x_1\{\Sigma^*\}x_3\{\Sigma^*\}\#\alpha_1 \\ w'_2 &= u_2\#w_2 & \langle \alpha_2 \rangle_{\text{int}} &= x_2\{\Sigma^*\}\#\alpha_2 \\ w'_3 &= u_1u_2u_3\#w_3 & \langle \alpha_3 \rangle_{\text{int}} &= x_1\{\Sigma^*\}x_2\{\Sigma^*\}x_3\{\Sigma^*\}\#\alpha_3\end{aligned}$$

However, it can be verified that $(\mathbf{abb}, \mathbf{abccbcc}, \mathbf{ababaaab})$ is a conjunctive match for $(\alpha_1, \alpha_2, \alpha_3)$ with variable mapping $(\mathbf{ab}, \mathbf{ab}, \mathbf{cc})$.

The following lemma is straightforward, but nevertheless helpful in the following proofs. It also points out how the semantics of a conjunctive xregex $\bar{\alpha}$ depends on the ref-languages of the individual components $\bar{\alpha}[i]$.

Lemma 2. Let $\bar{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_m) \in m\text{-CXRE}_{\Sigma, \mathcal{X}_s}$ and let $\beta_1, \beta_2, \dots, \beta_m \in \text{XRE}_{\Sigma, \mathcal{X}_s}$ such that, for every $i \in [m]$, $\mathcal{L}_{\text{ref}}(\alpha_i) = \mathcal{L}_{\text{ref}}(\beta_i)$. Then $\beta = (\beta_1, \beta_2, \dots, \beta_m)$ is a conjunctive xregex with $\mathcal{L}(\bar{\alpha}) = \mathcal{L}(\beta)$.

Proof. We first note that $\mathcal{L}_{\text{ref}}(\alpha_i) = \mathcal{L}_{\text{ref}}(\beta_i)$ for every $i \in [m]$ also implies that $\mathcal{L}_{\text{ref}}(\alpha_1\alpha_2\dots\alpha_m) = \mathcal{L}_{\text{ref}}(\beta_1\beta_2\dots\beta_m)$. If β is not a conjunctive xregex, then $\beta_1\beta_2\dots\beta_m$ is not sequential and, since $\mathcal{L}_{\text{ref}}(\alpha_1\alpha_2\dots\alpha_m) = \mathcal{L}_{\text{ref}}(\beta_1\beta_2\dots\beta_m)$, this directly implies that $\alpha_1\alpha_2\dots\alpha_m$ is not sequential, which is a contradiction.

It remains to show that $\mathcal{L}(\bar{\alpha}) = \mathcal{L}(\beta)$. Since, for every $i \in [m]$, $\mathcal{L}_{\text{ref}}(\alpha_i) = \mathcal{L}_{\text{ref}}(\beta_i)$, we can conclude that the $\nabla_{\beta_i} = \nabla_{\alpha_i}$. Moreover, this also implies that, for every $i \in [m]$, $\mathcal{L}_{\text{ref}}(\nabla_{\beta_i}\#\beta_i) = \mathcal{L}_{\text{ref}}(\nabla_{\alpha_i}\#\alpha_i)$ and therefore $\mathcal{L}_{\text{ref}}(\langle \alpha_i \rangle_{\text{int}}) = \mathcal{L}_{\text{ref}}(\langle \beta_i \rangle_{\text{int}})$. This directly implies that any (w_1, w_2, \dots, w_m) is a conjunctive match for $\bar{\alpha}$ if and only if it is a conjunctive match for β . \square

4 Conjunctive Xregex Path Queries

Finally, we can define conjunctive xregex path queries:

Definition 5 (CXRPQ). A conjunctive xregex path query (CXRPQ for short) is an \mathfrak{R} -CPQ $q = \bar{z} \leftarrow G_q$, where \mathfrak{R} is the class of xregex, $G_q = (V_q, E_q)$ with $E_q = \{(x_i, \alpha_i, y_i) \mid i \in [m]\}$ and $\bar{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_m)$ is a conjunctive xregex (which is also called the conjunctive xregex of q).

The semantics are defined by combining the semantics of conjunctive path queries and conjunctive xregex in the obvious way. More precisely, let $q = \bar{z} \leftarrow G_q$ be a CXRPQ, where $G_q = (V_q, E_q)$ with $E_q = \{(x_i, \alpha_i, y_i) \mid i \in [m]\}$, and let

\mathcal{D} be a graph database. Then $h : V_q \rightarrow V_{\mathcal{D}}$ is a matching morphism for q and \mathcal{D} if there is a conjunctive match (w_1, w_2, \dots, w_m) of $\bar{\alpha}$, such that, for every $i \in [m]$, \mathcal{D} contains a path from $h(x_i)$ to $h(y_i)$ labelled with w_i . Note that this definition of a matching morphism h for CXRPQ also implies definitions for $q_h(\mathcal{D})$, $q(\mathcal{D})$, $\mathcal{D} \models q$, $\llbracket q \rrbracket$ and $\llbracket \text{CXRPQ} \rrbracket$.

Since, for a $q \in \text{CXRPQ}$ and a fixed graph database \mathcal{D} , the set $q(\mathcal{D})$ is completely determined by the corresponding graph pattern and the set of conjunctive matches of the corresponding conjunctive xregex, the following can directly be concluded from the definitions.

Proposition 2. *Let $q = \bar{z} \leftarrow G_q$ be a CXRPQ with conjunctive xregex $\bar{\alpha} \in m\text{-CXRE}$. Let $\bar{\beta} \in m\text{-CXRE}$ and let $q' = \bar{z} \leftarrow G_{q'}$ be a CXRPQ where $G_{q'}$ is obtained from G_q by replacing each edge label $\bar{\alpha}[i]$ by $\bar{\beta}[i]$, for every $i \in [m]$. If $\mathcal{L}(\bar{\alpha}) = \mathcal{L}(\bar{\beta})$, then $q \equiv q'$.*

Next, we observe that evaluating CXRPQ is surprisingly difficult. Let $\Delta = \{\mathbf{a}, \mathbf{b}, \#\}$ and let $\alpha_{\text{ni}} = \#\mathbf{z}\{(\mathbf{a} \vee \mathbf{b})^*\}(\#\#\mathbf{z})^*\#\#\# \in \text{XRE}_{\Delta, \mathcal{X}_s}$.

Theorem 1. *Deciding whether a given graph-database over $\Sigma \supset \Delta$ contains a path labelled with some $w \in \mathcal{L}(\alpha_{\text{ni}})$ is PSpace-hard.*

Proof. We will prove the result by a reduction from the PSpace-complete NFA-intersection problem over binary alphabet $\{\mathbf{a}, \mathbf{b}\}$, which is defined as follows: Given NFA M_1, \dots, M_k over alphabet $\{\mathbf{a}, \mathbf{b}\}$, decide whether or not $\bigcap_{i=1}^k \mathcal{L}(M_i) \neq \emptyset$.

Without loss of generality, we assume that, for every $i \in [k]$, M_i has state set Q_i , transition function δ_i , initial state $q_{0,i}$, only one accepting state $q_{f,i}$, and we also assume that $\bigcap_{i=1}^k Q_i = \emptyset$. We transform the NFA M_1, \dots, M_k into a graph-database $\mathcal{D} = (V_{\mathcal{D}}, E_{\mathcal{D}})$ over alphabet $\Delta = \{\mathbf{a}, \mathbf{b}, \#\}$ as follows. For the sake of convenience, we assume that arcs can also be labelled by the words $\#\#$ and $\#\#\#$ (technically, these would be paths of length 2 and 3, respectively, instead of single arcs). We set $V_{\mathcal{D}} = (\bigcup_{i=1}^k Q_i) \cup \{s, t\}$, where $\{s, t\} \cap (\bigcup_{i=1}^k Q_i) = \emptyset$, and we set $E_{\mathcal{D}} = (\bigcup_{i=1}^k \delta_i) \cup \{(q_{f,i}, \#\#, q_{0,i+1}) \mid 1 \leq i \leq k-1\} \cup \{(s, \#, q_{0,1}), (q_{f,k}, \#\#\#, t)\}$.

We first note that in \mathcal{D} there is a path labelled with a word from $\mathcal{L}(\alpha_{\text{ni}})$ if and only if there is such a path from node s to node t . Let us now assume that there is a path in \mathcal{D} from s to t labelled with a word $w \in \mathcal{L}(\alpha_{\text{ni}})$. Since $w \in \mathcal{L}(\alpha_{\text{ni}})$, we can conclude that $w = \#w'(\#\#w')^\ell\#\#\#$ for some $w' \in \{\mathbf{a}, \mathbf{b}\}^*$ and $\ell \geq 0$. By the structure of \mathcal{D} , this directly implies that $\ell = k-1$ and, for every $i \in [k]$, there is a path labelled with w' from $q_{0,i}$ to $q_{f,i}$. Consequently, $w' \in \bigcap_{i=1}^k \mathcal{L}(M_i)$. On the other hand, if there is some word $w' \in \bigcup_{i=1}^k \mathcal{L}(M_i)$, then, for every $i \in [k]$, there is a path labelled with w' from $q_{0,i}$ to $q_{f,i}$, which directly implies that there is a path from s to t labelled with $\#w'(\#\#w')^{k-1}\#\#\# \in \mathcal{L}(\alpha_{\text{ni}})$. \square

Theorem 1 constitutes a rather strong negative result: CXRPQ-BOOL-EVAL is PSpace-hard in data-complexity, even for a single-edge query with a conjunctive xregex $\alpha \in \text{XRE}_{\Sigma, \mathcal{X}_s}$, where $|\Sigma| = 3$ and $|\mathcal{X}_s| = 1$. However, as explained in the introduction, α_{ni} gives a hint how to restrict CXRPQ to obtain more tractable fragments. Such fragments are investigated in the following Sections 5 and 6.

5 Variable-Star Free CXRPQ

An xregex α is *variable-star free* (*vstar-free*, for short), if no variable reference is subject to the $+$ -operator.¹ See Example 4 for an illustration. A conjunctive xregex $\bar{\alpha}$ of dimension m is vstar-free if, for every $i \in [m]$, $\bar{\alpha}[i]$ is vstar-free, and a $q \in \text{CXRPQ}$ is vstar-free, if its conjunctive xregex is vstar-free. Finally, let $\text{CXRPQ}^{\text{vsf}}$ be the class of vstar-free CXRPQ.

The main result of this section is as follows.

Theorem 2. $\text{CXRPQ}^{\text{vsf}}\text{-BOOL-EVAL}$ is in ExpSpace with respect to combined-complexity and in NL with respect to data-complexity.

Before discussing at greater length how to prove the upper bounds of Theorem 2, we first observe the following lower bound.

Theorem 3. $\text{CXRPQ}^{\text{vsf}}\text{-BOOL-EVAL}$

- is PSPACE -hard in combined-complexity, even for single-edge queries with xregex $\alpha \in \text{XRE}_{\Sigma, \mathcal{X}_s}$, where $|\Sigma| = 3$ and $|\mathcal{X}_s| = 1$, and
- is NL -hard in data-complexity, even for single-edge queries with xregex α , where $\alpha \in \text{RE}_{\Sigma}$ with $|\Sigma| = 2$.

Proof. For the PSPACE -hardness of $\text{CXRPQ}\text{-BOOL-EVAL}$ in combined-complexity, we can proceed similarly to the proof of Theorem 1. An instance M_1, \dots, M_k of the NFA-intersection problem over alphabet $\{\mathbf{a}, \mathbf{b}\}$ is transformed into a graph-database $\mathcal{D} = (V_{\mathcal{D}}, E_{\mathcal{D}})$ over alphabet $\Sigma = \{\mathbf{a}, \mathbf{b}, \#\}$ in the same way as in the proof of Theorem 1 and into a graph pattern $(x, \alpha_{\text{ni}}^k, y)$ with

$$\alpha_{\text{ni}}^k = \#\mathbf{z}\{(\mathbf{a} \vee \mathbf{b})^*\}(\#\#\mathbf{z})^{k-1}\#\#\#$$

(i. e., α_{ni}^k is obtained from α_{ni} by replacing $(\#\#\mathbf{z})^*$ by k copies of $\#\#\mathbf{z}$). We note that α_{ni}^k is vstar-free, but, due to the explicit $(k-1)$ -times repetition of $(\#\#\mathbf{z})$, it does not have constant size. It follows analogously as in the proof of Theorem 1 that in \mathcal{D} there is a path labelled with some $w \in \mathcal{L}(\alpha_{\text{ni}}^k)$ if and only if $\bigcap_{i=1}^k \mathcal{L}(M_i) \neq \emptyset$.

The NL -hardness of $\text{CXRPQ}\text{-BOOL-EVAL}$ in data-complexity can directly be concluded from the fact that $\text{CRPQ}\text{-BOOL-EVAL}$ is already NL -hard in data-complexity. In order to show the stronger statement of the result, we give a full proof. We consider the Boolean CRPQ q represented by the graph pattern (x, α, z) with $\alpha = \mathbf{ab}^*\mathbf{aa}$. From an arbitrary directed (and unlabelled) graph $G = (V, E)$ and two vertices $s, t \in V$, we can construct a graph-database $\mathcal{D} = (V_{\mathcal{D}}, E_{\mathcal{D}})$, where $V_{\mathcal{D}} = V \cup \{s', t'\}$ and $E_{\mathcal{D}} = \{(u, \mathbf{b}, v) \mid (u, v) \in E\} \cup \{(s', \mathbf{a}, s), (t, \mathbf{a}, t'), (t', \mathbf{a}, t'')\}$. It can be easily verified that in G there is a path from s to t if and only if in \mathcal{D} there is a path from s' to t'' labelled with $\mathbf{ab}^k\mathbf{aa}$ if and only if $\mathcal{D} \models q$ is true. \square

The upper bounds require more work. Before we can give an intuitive idea of our procedure, we have to introduce some more restrictions of xregex. An xregex α is

¹Since we use the Kleene-star r^* only as short hand form for $r^+ \vee \epsilon$, the term “variable-plus free” seems more appropriate. We nevertheless use the term “star free”, since it is much more common in the literature on regular expressions and languages.

- *variable-alternation free (valt-free)* if, for every subexpression $(\beta_1 \vee \beta_2)$ of α , neither β_1 nor β_2 contain any variable definition or variable reference.
- *variable-simple* if it is vstar-free and valt-free.
- *simple* if it is variable-simple and every variable definition $x\{\gamma\}$ is *basic*, i. e., γ is a classical regular expression or a single variable reference.
- in *normal form* if $\alpha = \alpha_1 \vee \alpha_2 \vee \dots \vee \alpha_m$ and, for every $i \in [m]$, α_i is simple.

Let us clarify these definitions with some intuitive explanations and examples. The condition of being vstar-free or valt-free can also be interpreted as follows: an xregex is vstar-free (or valt-free), if every subtree of its syntax tree rooted by a node for a $+$ -operation (for a \vee -operation, respectively) does not contain any nodes for variable definitions or references. If this is true for all $+$ -operation nodes *and* all \vee -operation nodes, then the xregex is variable-simple. Equivalently, α is variable-simple if $\alpha = \beta_1\beta_2 \dots \beta_k$, where each β_i is a classical regular expression, a variable reference or a variable definition $x\{\gamma\}$, where γ is also variable-simple. If additionally each variable definition $x\{\gamma\}$ is such that γ is a classical regular expressions or a single variable reference, then α is simple. Finally, an xregex is in normal form, if it is an alternation of simple xregex.

Example 4. *The xregex $x\{a^*\}(bx(c \vee a))^*b$ is not vstar-free, but valt-free. The xregex $x\{a^*\}y((bx) \vee (ca))^*b^*y$ is vstar-free, but not valt-free. Furthermore, the xregex $ax\{(b \vee c)^*by\{dxa^*\}\}bxa^*z\{d^*\}zy$ is variable-simple, but not simple, and $ax\{(b \vee c)^*da\}bxa^*y\{z\}xy$ is simple.*

We extend these restrictions to conjunctive xregex and CXRPQ in the obvious way, i. e., as also done above for vstar-free xregex.

We can now give a high-level “road map” for the proof of Theorem 2. As an important building block, we first prove that CXRPQ that are simple (in the sense defined above) can be evaluated in nondeterministic space $O(|q| \log(|\mathcal{D}|) + |q| \log(|q|))$ (see Lemma 3 below). Then, in Subsection 5.1, we show that every variable star-free conjunctive xregex can be transformed into an equivalent one in normal form (Lemmas 4, 5 and 6), which also means that $CXRPQ^{vsf}$ can be transformed into equivalent CXRPQ in normal form. A CXRPQ in normal form has a conjunctive xregex whose elements are alternations of simple xregex; thus, they can be evaluated by using Lemma 3. This directly yields a nondeterministic log-space algorithm with respect to data-complexity. For combined-complexity, we have to deal with the problem that the normal form causes a double exponential size blow-up. In order to get the ExpSpace upper bound, we observe that the first exponential blow-up can be handled by nondeterminism.

After this, in Subsection 5.3, we discuss more thoroughly what exactly causes the exponential blow-up in our normal form and how it can be avoided. Our observations can be transformed into the fragment $CXRPQ^{vsf,fl}$ (already mentioned in the introduction) that avoids the exponential blow-up, and therefore has a PSPACE -complete evaluation problem in combined-complexity.

We proceed with Lemma 3 (note that its proof is similar as the corresponding result for CRPQ (see, e. g., [8])).

Lemma 3. *Given a Boolean $q \in \text{CXRPQ}$ that is simple and a graph-database \mathcal{D} , we can nondeterministically decide whether or not $\mathcal{D} \models q$ in space $O(|q| \log(|\mathcal{D}|))$.*

Proof. Let q be defined by a graph pattern G_q with edge set $E_q = \{(x_i, \alpha_i, y_i) \mid i \in [m]\}$ and simple conjunctive xregex $\bar{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_m) \in \text{CXRE}_{\Sigma, \mathcal{X}_s}$. By assumption, for every $i \in [m]$, $\alpha_i = \alpha_{i,1} \alpha_{i,2} \dots \alpha_{i,t_i}$, where, for every $j \in [t_i]$, $\alpha_{i,j}$ is a classical regular expression, a variable reference, or a variable definition $x\{\gamma\}$, where γ is either a classical regular expression or a single variable reference.

We first note that since $\bar{\alpha}$ is simple, every variable $x \in \mathcal{X}_s$ has exactly one definition in every $\bar{v} \in \mathcal{L}_{\text{ref}}(\alpha_1) \times \mathcal{L}_{\text{ref}}(\alpha_2) \times \mathcal{L}_{\text{ref}}(\alpha_m)$. Consequently, if some variable x has a definition $x\{y\}$, then this definition and all references of x can be replaced by references of y without changing the set of conjunctive matches. In the following, we can therefore assume that, for every $i \in [m]$ and $j \in [t_i]$, $\alpha_{i,j}$ is a classical regular expression, a variable reference, or a variable definition $x\{\gamma\}$, where γ is a classical regular expression.

We next modify G_q such that every (x_i, α_i, y_i) is replaced by edges

$$(z_{i,0}, \alpha_{i,1}, z_{i,1}), (z_{i,1}, \alpha_{i,2}, z_{i,2}), \dots, (z_{i,t_i-1}, \alpha_{i,t_i}, z_{i,t_i}),$$

where $z_{i,0} = x_i$ and $z_{i,t_i} = y_i$. We denote this modified CXRPQ by q' and let $\bar{\beta} = (\beta_1, \dots, \beta_{m'})$ be its conjunctive xregex, i. e., $m' = \sum_{i=1}^m t_i$ and $(\beta_1, \dots, \beta_{m'})$ equals

$$(\alpha_{1,1}, \dots, \alpha_{1,t_1}, \alpha_{2,1}, \dots, \alpha_{2,t_2}, \dots, \alpha_{m,1}, \dots, \alpha_{m,t_m}).$$

It can be easily seen that $q \equiv q'$.

Next, for every $i \in [m']$, we define an NFA M_i with state-set Q_i and transition function δ_i as follows:

- If $\beta_i \in \{\gamma, x\{\gamma\}\}$ for some classical regular expression γ , then M_i accepts $\mathcal{L}(\gamma)$.
- If $\beta_i = x$ for some $x \in \mathcal{X}_s$, then M_i accepts Σ^* .

We now show how to decide $\mathcal{D} \models q'$ for a given graph database \mathcal{D} .

For technical reasons, we first add an ε -labelled self-loop to every node in \mathcal{D} and to every state in every M_i . Next, we define a graph $G_{q', \mathcal{D}} = (V_{q', \mathcal{D}}, E_{q', \mathcal{D}})$, where $V_{q', \mathcal{D}} = (V_{\mathcal{D}})^{m'} \times Q_1 \times Q_2 \times \dots \times Q_{m'}$ and $E_{q', \mathcal{D}}$ contains an edge

$$((u_1, \dots, u_{m'}, p_1, \dots, p_{m'}), (v_1, \dots, v_{m'}, r_1, \dots, r_{m'}))$$

if and only if there are $b_1, b_2, \dots, b_{m'} \in \Sigma \cup \{\varepsilon\}$ such that,

- for every $i \in [m']$, $r_i \in \delta_i(p_i, b_i)$ and $(u_i, b_i, v_i) \in E_{\mathcal{D}}$, and
- for every $i, i' \in [m']$, if $\beta_i, \beta_{i'} \in \{x, x\{\gamma\}\}$, for some $x \in \mathcal{X}_s$ and a classical regular expression γ , then $b_i = b_{i'}$.

Let $\bar{u} = (u_1, \dots, u_{m'}, p_1, \dots, p_{m'})$ and $\bar{v} = (v_1, \dots, v_{m'}, r_1, \dots, r_{m'})$ be two vertices from $G_{q', \mathcal{D}}$. We observe that in $G_{q', \mathcal{D}}$ there is a path from \bar{u} to \bar{v} if and only if there is a tuple $(w_1, w_2, \dots, w_{m'}) \in (\Sigma^*)^{m'}$ such that, for every $i \in [m']$, there is a w_i -labelled path from p_i to r_i in M_i , a w_i -labelled path from u_i to v_i in \mathcal{D} , and, for every $i, i' \in [m']$, if $\beta_i, \beta_{i'} \in \{x, x\{\gamma\}\}$, for some $x \in \mathcal{X}_s$ and a classical regular expression γ , then $w_i = w_{i'}$.

We say that a vertex \bar{u} of $G_{q',\mathcal{D}}$ is *initial*, if, for every $i \in [m']$, p_i is the initial state of M_i , and, for every $i, i' \in [m']$, $x_i = x_{i'}$ implies $u_i = u_{i'}$. Analogously, we say that a vertex \bar{v} of $G_{q',\mathcal{D}}$ is *final*, if, for every $i \in [m']$, r_i is the final state of M_i , and, for every $i, i' \in [m']$, $y_i = y_{i'}$ implies $v_i = v_{i'}$. With the observation from above, we can conclude that $\mathcal{D} \models q'$ if and only if in $G_{q',\mathcal{D}}$ there is a path from some initial to some final vertex. For simplicity, we add two additional vertices s and t to $G_{q',\mathcal{D}}$ with an edge (s, \bar{u}) and an edge (\bar{v}, t) for every initial vertex \bar{u} and final vertex \bar{v} . Checking $\mathcal{D} \models q'$ can now be done by checking whether there is a path from s to t in $G_{q',\mathcal{D}}$.

It remains to describe how this can be done in the claimed time and space bound. To this end, we first observe that the initial replacement of variable definition $x\{y\}$ by y can obviously be carried out in space $O(|q|)$, and it does not increase the size of q . Then, transforming $\bar{\alpha}$ into $\bar{\beta}$ can also be done in space $O(|q'|)$, where $|q'| = O(|q|)$. Moreover, we can transform each β_i into an NFA M_i of size $O(|\beta_i|)$ in space $O(|\beta_i|)$. Consequently, we can obtain all M_i in space $O(|q'|)$. In particular, we note that removing ε -transitions may result in NFA of size $O(|\beta_i|^2)$, which justifies our construction of $G_{q',\mathcal{D}}$ from above that can also handle possible ε -transitions of the M_i .

Now a vertex \bar{u} from $G_{q',\mathcal{D}}$ can be represented by $O(m')$ pointers to \mathcal{D} and $O(1)$ pointers to each of the M_i . Moreover, evaluating the edge relation of $G_{q',\mathcal{D}}$, i. e., checking for fixed vertices \bar{u}, \bar{v} of $G_{q',\mathcal{D}}$ whether there is an edge (\bar{u}, \bar{v}) , can be done as follows. Checking whether there are $b_1, b_2, \dots, b_{m'} \in \Sigma \cup \{\varepsilon\}$ such that, for every $i \in [m']$, $r_i \in \delta_i(p_i, b_i)$ and $(u_i, b_i, v_i) \in E_{\mathcal{D}}$ can be done by consulting the pointers that represent \bar{u} and \bar{v} , the transition functions of the M_i and the edge-relation of \mathcal{D} . Moreover, checking, for every $i, i' \in [m']$, whether $\beta_i, \beta_{i'} \in \{x, x\{\gamma\}\}$, for some $x \in \mathcal{X}_s$ and a classical regular expression γ , and $b_i = b_{i'}$, can be done by using two pointers to q' . Consequently, we can nondeterministically check whether there is a path from s to t in $G_{q',\mathcal{D}}$ in space $O(|q'|(\log(|\mathcal{D}|) + \log(|q'|))) = O(|q|(\log(|\mathcal{D}|) + \log(|q|))) = O(|q| \log(|\mathcal{D}|) + |q| \log(|q|)) = O(|q| \log(|\mathcal{D}|))$. \square

5.1 Construction of the Normal Form

The goal of this section is to show the following result (which also implies that $\text{CXRPQ}^{\text{vsf}}$ can be transformed in normal form).

Theorem 4. *Let $\bar{\alpha} \in m\text{-CXRE}_{\Sigma, \mathcal{X}_s}$ be vstar-free. Then there is an equivalent $\bar{\beta} \in m\text{-CXRE}_{\Sigma, \mathcal{X}'_s}$ in normal form with $|\bar{\beta}| = O\left(2^{2^{\text{poly}(|\bar{\alpha}|)} }\right)$.*

Our construction to transform vstar-free $\bar{\alpha} \in \text{CXRE}$ in normal form has three main steps, represented by Lemmas 4, 5 and 6. Before stating these lemmas, we will first explain the three main steps on an intuitive level with an example.

Let $\bar{\gamma} = (\gamma_1, \gamma_2)$ be a variable star free conjunctive xregex with

$$\begin{aligned}\gamma_1 &= x\{a^*y\{b^*\}az\} \vee (x\{b^*\} \cdot (z \vee y\{c^*\})) \\ \gamma_2 &= (a^* \vee x) \cdot z\{y \cdot (a \vee b)\}\end{aligned}$$

Step 1: We will “multiply-out” alternations with definitions or variables, which transforms γ_1 and γ_2 into alternations of variable-simple xregex.

$$\begin{aligned}\gamma_1 &= [x\{a^*y\{b^*\}az\}] \vee [x\{b^*\} \cdot z] \vee [x\{b^*\} \cdot y\{c^*\}] \\ \gamma_2 &= [a^* \cdot z\{y \cdot (a \vee b)\}] \vee [x \cdot z\{y \cdot (a \vee b)\}]\end{aligned}$$

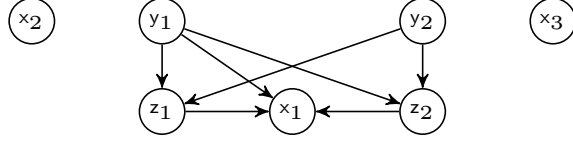


Figure 3: The DAG $G_{\bar{\gamma}}$ at Step 3.

This only results in an equivalent xregex because $\bar{\gamma}$ is vstar-free. This transformation may cause an exponential size blow-up.

Step 2: Next, we relabel the variables in such a way, that every variable has at most one definition in $\bar{\gamma}$. This also requires that variable references are substituted by several variable references (e. g., references of x by $x_1x_2x_3$). The size blow-up is at most quadratic.

$$\begin{aligned}\gamma_1 &= [x_1\{\mathbf{a}^*y_1\{\mathbf{b}^*\}\mathbf{a}z_1z_2\}] \vee [x_2\{\mathbf{b}^*\} \cdot z_1z_2] \vee [x_3\{\mathbf{b}^*\} \cdot y_2\{\mathbf{c}^*\}] \\ \gamma_2 &= [\mathbf{a}^* \cdot z_1\{y_1y_2 \cdot (\mathbf{a} \vee \mathbf{b})\}] \vee [x_1x_2x_3 \cdot z_2\{y_1y_2 \cdot (\mathbf{a} \vee \mathbf{b})\}]\end{aligned}$$

Step 3: We are now able to remove non-basic variable definitions, which remains to be done in order to obtain the normal form. The main idea is to split γ of a variable definition $x\{\gamma\}$ into smaller parts and (if they are not already definitions) store them in new variables. Instead of x , we can then reference the new variables, which allows to remove x . This, however, has to be done in the order given by the DAG induced by the relation $\preceq_{\bar{\alpha}}$ (see Figure 3). More precisely, we first consider the roots of the DAG representing variables x_2, x_3, y_1, y_2 , which have basic variable definitions and therefore can be left unchanged. After deleting them from the DAG, the nodes for z_1, z_2 are roots. The non-basic definition $z_1\{y_1y_2 \cdot (\mathbf{a} \vee \mathbf{b})\}$ is replaced by $u_1\{y_1\}u_2\{y_2\}u_3\{\mathbf{a} \vee \mathbf{b}\}$ and $z_2\{y_1y_2 \cdot (\mathbf{a} \vee \mathbf{b})\}$ is replaced by $u_4\{y_1\}u_5\{y_2\}u_6\{\mathbf{a} \vee \mathbf{b}\}$. All references for z_1 and z_2 are replaced by $\bar{z}_1 = u_1u_2u_3$ and $\bar{z}_2 = u_4u_5u_6$, respectively.

$$\begin{aligned}\gamma_1 &= [x_1\{\mathbf{a}^*y_1\{\mathbf{b}^*\}\mathbf{a}\bar{z}_1\bar{z}_2\}] \vee [x_2\{\mathbf{b}^*\} \cdot \bar{z}_1\bar{z}_2] \vee [x_3\{\mathbf{b}^*\} \cdot y_2\{\mathbf{c}^*\}] \\ \gamma_2 &= [\mathbf{a}^* \cdot u_1\{y_1\}u_2\{y_2\}u_3\{\mathbf{a} \vee \mathbf{b}\}] \vee \\ &\quad [x_1x_2x_3 \cdot u_4\{y_1\}u_5\{y_2\}u_6\{\mathbf{a} \vee \mathbf{b}\}]\end{aligned}$$

Finally, after deleting the nodes for z_1 and z_2 from the DAG, the node for x_1 is the only root left. Therefore, we can now replace

$$x_1\{\mathbf{a}^*y_1\{\mathbf{b}^*\}\mathbf{a}\bar{z}_1\bar{z}_2\} = x_1\{\mathbf{a}^*y_1\{\mathbf{b}^*\}\mathbf{a}u_1u_2u_3u_4u_5u_6\}$$

by the following concatenation of variable definitions:

$$\begin{aligned}u_7\{\mathbf{a}^*\}u_8\{y_1\{\mathbf{b}^*\}\}u_9\{\mathbf{a}\}u_{10}\{u_1\}u_{11}\{u_2\}u_{12}\{u_3\} \\ u_{13}\{u_4\}u_{14}\{u_5\}u_{15}\{u_6\}.\end{aligned}$$

Moreover, all references of x_1 are replaced by $\bar{x}_1 = u_7u_8u_9 \dots u_{15}$. The thus obtained conjunctive xregex is in fact in normal form.

We shall discuss a few particularities. While Steps 1 and 2 are more or less straightforward, Step 3 is more complicated. In particular, it is not easy to see why it is necessary to use new variables u with definition $u\{x\}$ instead of just

using the existing x . For example, if a variable definition $z\{y_1 a^* y_2\}$ would be replaced by $y_1 u\{a^*\} y_2$ and references of z by $y_1 u y_2$, then $y_1 u y_2$ must refer to ε in the case that z is not instantiated, which is not the case if y_1 is instantiated elsewhere. So we use $u_1\{y_1\}u_2\{a^*\}u_3\{y_2\}$ and replace references of z by $u_1 u_2 u_3$, which means that if z is not instantiated, then also none of the u_1 , u_2 and u_3 are instantiated in the modified xregex. Why it is necessary to apply these modifications according to the order given by $\preceq_{\bar{\gamma}}$ will be discussed in detail in the proof of the corresponding lemma further below. The possible exponential size blow-up of this step is also discussed in Subsection 5.3.

We now give the lemmas for the three steps of the construction (it will be helpful to keep in mind that the constructions of these lemmas are illustrated by the example from above). The proofs of the first two lemmas are more or less straightforward proofs of correctness for the constructions of Step 1 and 2 illustrated above. The proof of the third lemma is technically more involved.

Lemma 4. (Step 1) *Let $\bar{\alpha} \in m\text{-CXRE}_{\Sigma, \mathcal{X}_s}$ be vstar-free. Then there is an equivalent $\bar{\beta} \in m\text{-CXRE}_{\Sigma, \mathcal{X}_s}$ with $|\bar{\beta}| = O(2^{|\bar{\alpha}|})$, such that each $\bar{\beta}[i]$ is an alternation of variable-simple xregex.*

Proof. Let $\bar{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_m)$. For every $i \in [m]$, we transform α_i into $\beta_i = \beta_{i,1} \vee \beta_{i,2} \vee \dots \vee \beta_{i,t_i}$, such that, for every $j \in [t_i]$, $\beta_{i,j}$ is variable-simple.

If α_i is not already in the form described above (otherwise we set $\beta_i = \alpha_i$ and are done), i.e., an alternation of variable-simple xregex, then it has the form $\alpha_i = \pi_1 \vee \pi_2 \vee \dots \vee \pi_q$ (note that $q = 1$ is possible), where at least one π_j is not valt-free (since α_i is vstar-free, each π_j is vstar-free as well). This means that π_j has some subexpression $\gamma = (\gamma_1 \vee \gamma_2)$ and γ contains a variable definition or a variable reference. Let $\pi_{j,1}$ be obtained from π_j by replacing γ by γ_1 and let $\pi_{j,2}$ be obtained from π_j by replacing γ by γ_2 . Finally, let $\alpha'_i = \pi_1 \vee \dots \vee \pi_{j,1} \vee \pi_{j,2} \vee \dots \vee \pi_q$ and let $\bar{\alpha}'$ be obtained from $\bar{\alpha}$ by replacing α_i by α'_i .

We first note that $\bar{\alpha}'$ is sequential, variable-acyclic and vstar-free, which means that $\bar{\alpha}'$ is a vstar-free conjunctive xregex. Moreover, since α_i is vstar-free, γ is not under a $+$ -operator, which implies that $\mathcal{L}_{\text{ref}}(\pi_j) = \mathcal{L}_{\text{ref}}(\pi_{j,1} \vee \pi_{j,2})$ (which would not be the case if γ was under a $+$ -operator), and therefore also

$$\mathcal{L}_{\text{ref}}(\pi_1 \vee \pi_2 \vee \dots \vee \pi_q) = \mathcal{L}_{\text{ref}}(\pi_1 \vee \dots \vee \pi_{j,1} \vee \pi_{j,2} \vee \dots \vee \pi_q),$$

Consequently, we have $\mathcal{L}_{\text{ref}}(\alpha_i) = \mathcal{L}_{\text{ref}}(\alpha'_i)$, which, by Lemma 2, means that $\mathcal{L}(\bar{\alpha}) = \mathcal{L}(\bar{\alpha}')$.

By repeating this step, we can therefore transform $\bar{\alpha}$ into an equivalent $\bar{\beta} \in m\text{-CXRE}_{\Sigma, \mathcal{X}_s}$ such that each $\bar{\beta}[i]$ is an alternation of variable-simple xregex. Finally, we note that in each step of this construction, the size of the xregex can double in the worst case. Thus, $|\bar{\beta}| = O(2^{|\bar{\alpha}|})$. \square

Lemma 5. (Step 2) *Let $\bar{\alpha} \in m\text{-CXRE}_{\Sigma, \mathcal{X}_s}$ such that, for every $i \in [m]$, $\bar{\alpha}[i]$ is an alternation of variable-simple xregex. Then there is an equivalent $\bar{\beta} \in m\text{-CXRE}_{\Sigma, \mathcal{X}'_s}$ such that, for every $i \in [m]$, $\bar{\beta}[i]$ is an alternation of variable-simple xregex, and every $x \in \mathcal{X}'_s$ has at most one definition in $\bar{\beta}$. Moreover, $|\bar{\beta}| = O(|\bar{\alpha}|^2)$.*

Proof. Let $\bar{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_m)$, and, for every $i \in [m]$, let

$$\alpha_i = \alpha_{i,1} \vee \alpha_{i,2} \vee \dots \vee \alpha_{i,t_i},$$

where, for every $j \in [t_i]$, $\alpha_{i,j}$ is variable-simple. We transform $\bar{\alpha}$ into an equivalent $\bar{\beta}$ with the property described in the statement of the lemma as follows.

Let $x \in \mathcal{X}_s$ and let $i \in [m]$ be such that α_i contains a definition of x (since $\bar{\alpha}$ is a conjunctive xregex, there is at most one such $i \in [m]$). Moreover, if, for some $j \in [t_i]$, there are at least two definitions for x in $\alpha_{i,j}$, then $\alpha_{i,j}$ is not sequential or not valt-free, which is a contradiction. Thus, for every $j \in [t_i]$, there is at most one definition for x in $\alpha_{i,j}$, and, for the sake of convenience, we shall assume that, for every $j \in [t_i]$, $\alpha_{i,j}$ contains a definition for x (it will be obvious how to adapt the construction when this is not the case). Now, for every $j \in [t_i]$, we replace in $\alpha_{i,j}$ the variable definition $x\{\gamma\}$ by $x^{(j)}\{\gamma\}$, where $x^{(j)}$ is a new variable. Let $\bar{\alpha}' = (\alpha'_1, \alpha'_2, \dots, \alpha'_m)$ be the thus obtained conjunctive xregex. Obviously, every variable $x^{(j)}$ with $j \in [t_i]$ has exactly one definition in $\bar{\alpha}'$, but $\bar{\alpha}'$ is not necessarily equivalent to $\bar{\alpha}$, since we did not change the references for variable x . Before we do this, we note that, for any ref-word $v \in \mathcal{L}_{\text{ref}}(\alpha'_i)$, there is exactly one $j \in [t_i]$ such that v contains a definition of $x^{(j)}$ (or, in other words, for exactly one $j \in [t_i]$ the definition for $x^{(j)}$ is instantiated), while for all other $j' \in [t_i] \setminus \{j\}$, there is no definition of $x^{(j')}$, which means that the image for $x^{(j')}$ with respect to v 's variable mapping is necessarily empty. This is due to the fact that α_i is sequential, variable alternation-free and we assume that, for every $j \in [t_i]$, $\alpha_{i,j}$ has a definition of x . Hence, we can simply substitute *all* variable references for x in $\bar{\alpha}'$ (not only those in α_i) by $\bar{x} = \prod_{j=1}^{t_i} x^{(j)}$. We denote the thus modified variant of $\bar{\alpha}$ by $\bar{\alpha}'' = (\alpha''_1, \alpha''_2, \dots, \alpha''_m)$.

It is straightforward to see that $\bar{\alpha}''$ is sequential and variable-acyclic; thus, it is a conjunctive xregex, and also that, for every $i \in [m]$, α''_i is an alternation of variable-simple xregex. With the observations from above, it is also easy to see that $\bar{\alpha}$ and $\bar{\alpha}''$ are equivalent. More precisely, every tuple $\bar{v} \in \mathcal{L}_{\text{ref}}(\alpha_1) \times \mathcal{L}_{\text{ref}}(\alpha_2) \times \dots \times \mathcal{L}_{\text{ref}}(\alpha_m)$ corresponds to a tuple $\bar{v}'' \in \mathcal{L}_{\text{ref}}(\alpha''_1) \times \mathcal{L}_{\text{ref}}(\alpha''_2) \times \dots \times \mathcal{L}_{\text{ref}}(\alpha''_m)$, obtained from \bar{v} by replacing each variable reference x by \bar{x} , and the variable definition $x \triangleright u \triangleleft x$ by $x^{(j)} \triangleright u \triangleleft x^{(j)}$, where $j \in [t_i]$ is such that the variable definition of x in \bar{v} was generated by $\alpha_{i,j}$.

By repeating this modification step with respect to every variable from \mathcal{X}_s , we can obtain a conjunctive xregex $\bar{\beta} = (\beta_1, \beta_2, \dots, \beta_m)$ with the desired property that is equivalent to $\bar{\alpha}$. Moreover, this modification replaces each variable reference by a sequence of $O(|\bar{\alpha}|)$ variable references; thus, $|\bar{\beta}| = O(|\bar{\alpha}|^2)$. \square

Lemma 6. (Step 3) *Let $\bar{\alpha} \in m\text{-CXRE}_{\Sigma, \mathcal{X}_s}$ such that, for every $i \in [m]$, $\bar{\alpha}[i]$ is an alternation of variable-simple xregex, and every $x \in \mathcal{X}_s$ has at most one definition in $\bar{\alpha}$. Then there is an equivalent $\bar{\beta} \in m\text{-CXRE}_{\Sigma, \mathcal{X}'_s}$ in normal form with $|\bar{\beta}| = O(|\bar{\alpha}|^{|\mathcal{X}_s|+1})$.*

Proof. Let $\bar{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_m)$, and, for every $i \in [m]$, let

$$\alpha_i = \alpha_{i,1} \vee \alpha_{i,2} \vee \dots \vee \alpha_{i,t_i},$$

where, for every $j \in [t_i]$, $\alpha_{i,j}$ is variable-simple, and, for every $x \in \mathcal{X}_s$, there is at most one definition for x in $\bar{\alpha}$. First, we define as general modification step, which can be applied to any variable definition in $\bar{\alpha}$. This modification step will then be used in order to transform $\bar{\alpha}$ into normal form.

Main modification step: Let $z\{\gamma\}$ be a variable definition of $\bar{\alpha}$. Since γ is variable-simple, $\gamma = \gamma_1 \gamma_2 \dots \gamma_p$ such that each γ_ℓ with $\ell \in [p]$ is either a classical regular expression, a variable definition, or a variable reference.

For every $\ell \in [p]$, we define a γ'_ℓ as follows. If γ_ℓ is a variable definition $y_\ell\{\dots\}$, then we set $\gamma'_\ell = \gamma_\ell$. If γ_ℓ is a classical regular expression or a variable reference, then we set $\gamma'_\ell = y_\ell\{\gamma_\ell\}$ for a *new* variable $y_\ell \notin \mathcal{X}_s$. Note that $\gamma'_1\gamma'_2\dots\gamma'_p = y_1\{\dots\}y_2\{\dots\}\dots y_p\{\dots\}$ is a concatenation of variable definitions. Next, we replace the variable definition $\mathbf{z}\{\gamma\}$ in $\bar{\alpha}$ by $\gamma'_1\gamma'_2\dots\gamma'_p$. Then, we replace all variable references of \mathbf{z} in $\bar{\alpha}$ by $y_1y_2\dots y_p$.

Let the thus modified version of $\bar{\alpha}$ be denoted by $\bar{\beta} = (\beta_1, \beta_2, \dots, \beta_m)$. It can be easily seen that $\bar{\beta}$ is sequential and variable acyclic and therefore a conjunctive xregex. Moreover, for every $i \in [m]$, β_i is still an alternation of variable-simple xregex.

Correctness of the main modification step: Next, we show that $\mathcal{L}(\bar{\alpha}) = \mathcal{L}(\bar{\beta})$. The difficulty in doing so is due to the fact that we cannot assume, for every $i \in [m]$, that $\mathcal{L}_{\text{ref}}(\alpha_i) = \mathcal{L}_{\text{ref}}(\beta_i)$; i. e., we cannot conveniently apply Lemma 2, as it has been done in the proof of Lemma 4.

We need a few more notational preliminaries. We assume that $r \in [m]$ and $\hat{j} \in [t_r]$ is such that the modified variable definition $\mathbf{z}\{\gamma\}$ is in $\alpha_{r,\hat{j}}$. Moreover, let $\beta_r = \beta_{r,1} \vee \beta_{r,2} \vee \dots \vee \beta_{r,t_r}$, where, for every $j \in [t_r]$, $\beta_{r,j}$ is obtained from $\alpha_{r,j}$ by the construction from above (i. e., $\beta_{r,\hat{j}}$ is obtained from $\alpha_{r,\hat{j}}$ by replacing the definition $\mathbf{z}\{\gamma\}$ with $\gamma'_1\gamma'_2\dots\gamma'_p$ and all references \mathbf{z} with $y_1y_2\dots y_p$, and all other $\beta_{r,j}$, with $j \in [t_r] \setminus \{\hat{j}\}$, are obtained from $\alpha_{r,j}$ by only replacing all references \mathbf{z} with $y_1y_2\dots y_p$). Furthermore, all other $\beta_{r'}$ with $r' \in [m] \setminus \{r\}$ are obtained from $\alpha_{r'}$ by only replacing each variable reference for \mathbf{z} by $y_1y_2\dots y_p$.

Let $\bar{w} = (w_1, w_2, \dots, w_m)$ be a conjunctive match for $\bar{\alpha}$, with variable mapping $\psi_{\bar{\alpha}}$, i. e., for every $i \in [m]$, there is a ref-word $v_{i,\alpha} = v'_{i,\alpha} \# v''_{i,\alpha} \in \mathcal{L}_{\text{ref}}(\langle \alpha_i \rangle_{\text{int}})$ with variable mapping $\psi_{\bar{\alpha}}$ and with $\text{deref}(v_{i,\alpha}) = \text{deref}(v'_{i,\alpha}) \# w_i$. In order to show that \bar{w} is a conjunctive match for $\bar{\beta}$, we have to show that, for every $i \in [m]$, there are ref-words $v_{i,\beta} = v'_{i,\beta} \# v''_{i,\beta} \in \mathcal{L}_{\text{ref}}(\langle \beta_i \rangle_{\text{int}})$ with $\text{deref}(v_{i,\beta}) = \text{deref}(v'_{i,\beta}) \# w_i$, and $\text{vmap}_{v_{1,\beta}} = \text{vmap}_{v_{2,\beta}} = \dots = \text{vmap}_{v_{m,\beta}}$.

There are two cases, which we shall first discuss on an intuitive level. The simple case is when the ref-word $v''_{r,\alpha} \in \mathcal{L}_{\text{ref}}(\alpha_r)$ (i. e., the ref-word that yields w_r) can be constructed by some $\alpha_{r,j}$ with $j \in [t_r] \setminus \{\hat{j}\}$; more precisely, we have $v''_{r,\alpha} \in \mathcal{L}_{\text{ref}}(\alpha_{r,j})$ for some $j \in [t_r] \setminus \{\hat{j}\}$. In this case, the modified variable definition $\mathbf{z}\{\gamma\}$ is not instantiated. This means that suitable ref-words $v_{i,\beta} = v'_{i,\beta} \# v''_{i,\beta}$ that witness that \bar{w} is a conjunctive match for $\bar{\beta}$ can be easily obtained from the witnesses $v'_{i,\alpha} \# v''_{i,\alpha}$. We only have to take care of the variables y_ℓ , for which references may occur in $v''_{i,\alpha}$. However, these variables, regardless of whether they have been newly introduced by the modification step or whether definitions for them have already been present before, have only definitions in $\alpha_{r,\hat{j}}$ (note that for the variables that are not newly introduced, this is only true since we assume that every variable has at most one definition in $\bar{\alpha}$), and therefore their images must be empty.

The more complicated case is when, for every $j \in [t_r] \setminus \{\hat{j}\}$, $v''_{r,\alpha}$ is not in $\mathcal{L}_{\text{ref}}(\alpha_{r,j})$, which means that the only way for α_r to produce $v''_{r,\alpha}$ is to use the subexpression $\alpha_{r,\hat{j}}$, which is the one that contains the variable definition $\mathbf{z}\{\gamma\} = \mathbf{z}\{\gamma_1\gamma_2\dots\gamma_p\}$ that has been modified by the construction. Moreover, since $\alpha_{r,\hat{j}}$ is variable-simple, we also know that it must be instantiated by $v''_{r,\alpha}$, i. e., $v''_{r,\alpha} \in \mathcal{L}_{\text{ref}}(\alpha_{r,\hat{j}})$ and $v''_{r,\alpha}$ does contain a factor $\overset{z}{\triangleright} g_1g_2\dots g_p \triangleleft^z$, where, for every $\ell \in [p]$, $g_\ell \in \mathcal{L}_{\text{ref}}(\gamma_\ell)$. In this case, we have to obtain the witnesses

$v_{i,\beta} = v'_{i,\beta} \# v''_{i,\beta}$ from the witnesses $v_{i,\alpha} = v'_{i,\alpha} \# v''_{i,\alpha}$ in such a way that the g_ℓ parts are handled by the variables y_ℓ . We now discuss these two cases more formally.

(C1) $[v''_{r,\alpha} \in \mathcal{L}_{\text{ref}}(\alpha_{r,j}) \text{ for some } j \in [t_r] \setminus \{\widehat{j}\}]$

In this case, for every $i \in [m]$, we let $v''_{i,\beta}$ be the ref-word obtained from $v''_{i,\alpha}$ by replacing each variable reference z by $y_1 y_2 \dots y_p$. From the fact that, for every $i \in [m]$, $v''_{i,\alpha} \in \mathcal{L}_{\text{ref}}(\alpha_{r,i})$, and from the assumption made in this case, it follows that also $v''_{i,\beta} \in \mathcal{L}_{\text{ref}}(\beta_{r,i})$.

By construction, we have $\nabla_{\beta_r} = \nabla_{\alpha_r}$. Consequently, we can set $v'_{r,\beta} = v'_{r,\alpha}$ and observe that $v_{r,\beta} = v'_{r,\beta} \# v''_{r,\beta} \in \mathcal{L}_{\text{ref}}(\langle \beta_r \rangle_{\text{int}})$. It remains to show that $\text{deref}(v_{r,\beta})$ equals $\text{deref}(v'_{r,\beta}) \# w_r$. To this end, we first note that, for every variable x different from z or any y_ℓ , we have $\text{vmap}_{v_{r,\beta}}(x) = \text{vmap}_{v'_{r,\beta}}(x)$. Moreover, the variables y_ℓ have no definition in $v''_{r,\beta}$, which means that $\text{vmap}_{v_{r,\beta}}(y_\ell) = \varepsilon$, and since $v''_{r,\alpha} \in \mathcal{L}_{\text{ref}}(\alpha_{r,j})$ for some $j \in [t_r] \setminus \{\widehat{j}\}$, we also have $\text{vmap}_{v_{r,\alpha}}(z) = \varepsilon$. Finally, since $v''_{i,\beta}$ is obtained from $v''_{i,\alpha}$ by replacing each variable-reference z by $y_1 y_2 \dots y_p$ and since $\text{deref}(v_{r,\alpha}) = \text{deref}(v'_{r,\alpha}) \# w_r$, it follows that $\text{deref}(v_{r,\beta}) = \text{deref}(v'_{r,\beta}) \# w_r$.

It remains to define $v_{i,\beta}$ for every $i \in [m] \setminus \{r\}$. To this end, for every $i \in [m] \setminus \{r\}$, we define $v''_{i,\beta}$ analogously to $v''_{r,\beta}$ done above, i.e., we let $v''_{i,\beta}$ be obtained from $v''_{i,\alpha}$ by replacing each variable reference z by $y_1 y_2 \dots y_p$. As above, we note that, for every $i \in [m]$, $v''_{i,\beta} \in \mathcal{L}_{\text{ref}}(\beta_{r,i})$.

For every $i \in [m] \setminus \{r\}$, due to the new variables y_ℓ with $\ell \in [p]$, it is not the case that $\nabla_{\beta_i} = \nabla_{\alpha_i}$. Therefore, for every $i \in [m] \setminus \{r\}$, let $v'_{i,\beta}$ be such that $\text{vmap}_{v'_{i,\beta}}$ is equal to $\text{vmap}_{v'_{i,\alpha}}$ extended by $\text{vmap}_{v'_{i,\beta}}(y_1) = \text{vmap}_{v'_{i,\beta}}(y_2) = \dots = \text{vmap}_{v'_{i,\beta}}(y_p) = \varepsilon$. We observe that $v_{i,\beta} = v'_{i,\beta} \# v''_{i,\beta} \in \mathcal{L}_{\text{ref}}(\langle \beta_i \rangle_{\text{int}})$ and $\text{deref}(v_{i,\beta}) = \text{deref}(v'_{i,\beta}) \# w_i$, since clearly $\text{vmap}_{v_{i,\beta}}(z) = \varepsilon$.

Finally, we note that all $\text{vmap}_{v_{i,\beta}}$ with $i \in [m]$ are identical to $\text{vmap}_{v_{i,\alpha}}$ for all variables $x \notin \{y_1, y_2, \dots, y_p\}$, and that $\text{vmap}_{v_{i,\beta}}(y_1) = \text{vmap}_{v_{i,\beta}}(y_2) = \dots = \text{vmap}_{v_{i,\beta}}(y_p) = \varepsilon$.

(C2) $[v''_{r,\alpha} \in \mathcal{L}_{\text{ref}}(\alpha_{r,\widehat{j}}) \setminus (\bigcup_{i \in [t_r] \setminus \{\widehat{j}\}} \mathcal{L}_{\text{ref}}(\alpha_{r,i}))]$

Since $\alpha_{r,\widehat{j}}$ is valt-free, the assumption made in this case implies that $v''_{r,\alpha}$ must contain a definition $z \triangleright g_1 g_2 \dots g_p \triangleleft^z$, where, for every $\ell \in [p]$, $g_\ell \in \mathcal{L}_{\text{ref}}(\gamma_\ell)$. We let $v''_{r,\beta}$ be obtained from $v''_{r,\alpha}$ by replacing $z \triangleright g_1 g_2 \dots g_p \triangleleft^z$ by $y_1 \triangleright g_1 \triangleleft^{y_2} y_2 \triangleright g_2 \triangleleft^{y_2} \dots y_p \triangleright g_p \triangleleft^{y_p}$ and all occurrences of z by $y_1 y_2 \dots y_p$. By construction, $v''_{r,\beta} \in \mathcal{L}_{\text{ref}}(\beta_{r,\widehat{j}})$ is satisfied. As in Case 1, we have $\nabla_{\beta_r} = \nabla_{\alpha_r}$ and can therefore set $v'_{r,\beta} = v'_{r,\alpha}$ to get that $v_{r,\beta} = v'_{r,\beta} \# v''_{r,\beta} \in \mathcal{L}_{\text{ref}}(\langle \beta_r \rangle_{\text{int}})$. In order to see that $\text{deref}(v_{r,\beta}) = \text{deref}(v'_{r,\beta}) \# w_r$, we again note that for every variable x different from z or any y_ℓ , we have that $\text{vmap}_{v_{r,\beta}}(x) = \text{vmap}_{v'_{r,\beta}}(x)$. Moreover, for the variables y_ℓ , we have that $\text{vmap}_{v_{r,\beta}}(y_\ell) = u_\ell$, such that $u_1 u_2 \dots u_p = \text{vmap}_{v_{r,\alpha}}(z)$. This follows from the fact that $v''_{r,\beta}$ is obtained from $v''_{r,\alpha}$ by replacing $z \triangleright g_1 g_2 \dots g_p \triangleleft^z$ by $y_1 \triangleright g_1 \triangleleft^{y_2} y_2 \triangleright g_2 \triangleleft^{y_2} \dots y_p \triangleright g_p \triangleleft^{y_p}$. This directly implies that $\text{deref}(v_{r,\beta}) = \text{deref}(v'_{r,\beta}) \# w_r$.

For every $i \in [m] \setminus \{r\}$, we let $v''_{i,\beta}$ be obtained from $v''_{i,\alpha}$ by replacing each variable-reference z by $y_1 y_2 \dots y_p$. By construction, it follows that

$v''_{i,\beta} \in \mathcal{L}_{\text{ref}}(\beta_i)$. For every $i \in [m] \setminus \{r\}$, β_i does not contain a definition for y_ℓ , which means that ∇_{β_i} contains the definition $y_\ell\{\Sigma^*\}$. Thus, for every $i \in [m] \setminus \{r\}$, we can define $v'_{i,\beta}$ such that $\text{vmap}_{v'_{i,\beta}}$ is equal to $\text{vmap}_{v'_{i,\alpha}}$ up to the following exceptions: for every $\ell \in [p]$, $\text{vmap}_{v'_{i,\beta}}(y_\ell) = \text{vmap}_{v_{r,\beta}}(y_\ell)$ (recall that $v_{r,\beta}$ contains a definition for each y_ℓ). By construction, it can be easily seen that in fact $v_{i,\beta} = v'_{i,\beta} \# v''_{i,\beta} \in \mathcal{L}_{\text{ref}}(\langle \beta_r \rangle_{\text{int}})$ and $\text{deref}(v_{i,\beta}) = \text{deref}(v'_{i,\beta}) \# w_i$ holds for every $i \in [m] \setminus \{r\}$. Furthermore, we observe that all $v_{i,\beta}$ with $i \in [m]$ are defined in such a way that they have the same variable mapping.

This shows that every conjunctive match for $\bar{\alpha}$ is also a conjunctive match for $\bar{\beta}$. Moreover, on close inspection it can be noted that the way how we obtained appropriate ref-words $v_{i,\beta}$ from the ref-words $v_{i,\alpha}$ can be reversed in order to prove that every conjunctive match for $\bar{\beta}$ is also a conjunctive match for $\bar{\alpha}$. More precisely, we replace $y_1 y_2 \dots y_p$ by z and ${}^{y_1} \triangleright g_1 \triangleleft^{y_2} {}^{y_2} \triangleright g_2 \triangleleft^{y_3} \dots {}^{y_p} \triangleright g_p \triangleleft^{y_p}$ by ${}^z \triangleright g_1 g_2 \dots g_p \triangleleft^z$ instead, and instead of changing the ref-words $v'_{i,\alpha}$ such that the variable mappings are extended for variables y_ℓ , we change the ref-words $v'_{i,\beta}$ such that the variable mappings are restricted accordingly (in particular, in Case 2, $v'_{i,\alpha}$ must be chosen such that the image of z equals the concatenation of the images of y_ℓ as determined by $v'_{r,\alpha}$).

Consequently, $\mathcal{L}(\bar{\alpha}) = \mathcal{L}(\bar{\beta})$, which concludes the proof of the correctness of the main modification step.

Transforming $\bar{\alpha}$ in normal form: It remains to show how $\bar{\alpha}$ can be transformed into normal form by applications of the main modification step.

The idea is to apply the main modification step to each non-basic variable definition of $\bar{\alpha}$, in order to transform them into basic ones. However, this does not necessarily work if we do not apply the modification steps in a specific order. In order to illustrate the potential problem, let $z\{\gamma\}$ be a variable definition and let there be some other variable definition $x\{\gamma'\}$, such that γ' contains a reference of z . Then the main modification step will replace $z\{\gamma\}$ by a concatenation of variable definitions, but also the reference of z in γ' by a concatenation $y_1 y_2 \dots y_p$ of variable references. In the case that γ' is a single variable reference and $p \geq 2$, this will turn a basic variable definition, namely $x\{z\}$, into a non-basic one, namely $x\{y_1 y_2 \dots y_p\}$.

As mentioned above, we avoid this by applying the main modification steps to the non-basic variable definitions in a particular order. To this end, we recall the binary relation $\preceq_{\bar{\alpha}}$ over \mathcal{X}_s with $x \preceq_{\bar{\alpha}} y$ if x has a reference or a definition in the definition of y . In particular, we consider the directed graph $G_{\bar{\alpha}} = (\mathcal{X}_s, \{(x, y) \mid x \preceq_{\bar{\alpha}} y\})$.

Since $\preceq_{\bar{\alpha}}$ is acyclic, we know that $G_{\bar{\alpha}}$ is a directed acyclic graph (DAG). Moreover, the roots of $G_{\bar{\alpha}}$ (i.e., nodes without incoming arcs) correspond to the minimal elements of $\preceq_{\bar{\alpha}}$, which are exactly the variables x whose variable definitions do not contain any references or definitions. We now apply the main modification step to the non-basic variable definition governed by the structure of $G_{\bar{\alpha}}$.

We repeatedly choose some root x of (the current version of) $G_{\bar{\alpha}}$, perform a modification on $\bar{\alpha}$ (which might be the identity) and then we delete this root x . Since $G_{\bar{\alpha}}$ is a DAG, this step either deletes the last node of $G_{\bar{\alpha}}$, which terminates the procedure, or after its application there is at least one other root in $G_{\bar{\alpha}}$.

Consequently, this procedure terminates after $|\mathcal{X}_s|$ steps. A single step works as follows.

If the definition $x\{\gamma\}$ of the root x is basic, i. e., γ is a classical regular expression or a single variable reference, then we will just remove the root x from $G_{\bar{\alpha}}$ without modifying $\bar{\alpha}$. If, on the other hand, $x\{\gamma\}$ is not basic, then we apply the main modification step to $x\{\gamma\}$ (which modifies $\bar{\alpha}$), and then we remove the root x from $G_{\bar{\alpha}}$. Note that in this procedure, $x\{\gamma\}$ always refers to the definition of x in the current version of $\bar{\alpha}$, which may have been changed by applications of the main modification steps, i. e., variable references in the original definition of x might have been replaced by concatenations of new variable references.

Let $\beta = (\beta_1, \beta_2, \dots, \beta_m)$ be the conjunctive xregex obtained by this procedure. Due to the correctness of the main modification step, we know that $\mathcal{L}(\bar{\alpha}) = \mathcal{L}(\beta)$, and that, for every $i \in [m]$, β_i is still an alternation of variable-simple xregex. In order to conclude that β is in normal form, it only remains to show that β does not contain non-basic variable definitions. To this end, we first state some observations about the procedure described above:

1. References of new variables that are introduced by the procedure will never be replaced anymore. References of variables from \mathcal{X}_s that have already been considered by the procedure, but are not removed by an application of the main modification step, will also never be replaced anymore
2. If the definition of a variable $x \in \mathcal{X}_s$ is initially non-basic, then this variable will necessarily be removed by the procedure. This is due to the fact that the procedure considers each variable from \mathcal{X}_s , and that other applications of the main modification step with respect to some variables different from x cannot transform the definition of x into a basic one.
3. The definition of a variable $x\{\gamma\}$, where γ is a classical regular expression, is never changed by the procedure. This holds for the case that $x \in \mathcal{X}_s$ and $x\{\gamma\}$ is its initial definition in $\bar{\alpha}$ before the procedure, and for the case that x is a new variable and $x\{\gamma\}$ is created by the application of the main modification step.
4. Whenever a variable definition $x\{\gamma\}$ is changed in such a way that references in γ are replaced by concatenations of references, then $x \in \mathcal{X}_s$ and x has not yet been considered by the procedure. This can be seen as follows. If $x \notin \mathcal{X}_s$, i. e., x is a new variable, then there was some application of the main modification step with respect to some z that has created variable x with an initial definition $x\{\gamma'\}$. By definition of the main modification step, this means that γ' is either a classical regular expression, or a single variable reference. In the first case, due to Point 3, $x\{\gamma'\}$ cannot be changed anymore, which contradicts our assumption. Thus, $\gamma' = y$. If y is a new variable, then y is not replaced, due to Point 1. If $y \in \mathcal{X}_s$, then this means that y has a reference in the definition of the variable z , which means that it must already have been considered when the modification step is carried out with respect to z . Hence, due to Point 1, it will not be replaced anymore. Therefore, we can assume that $x \in \mathcal{X}_s$.

If x has already been considered by the procedure, then, since it has not been removed, it initially had a basic definition $x\{\gamma'\}$. If γ' is a regular

expression, then, due to Point 3, it cannot be changed, which is a contradiction. Therefore, $\gamma' = y$, and in the same way as above, we can conclude that y cannot be a new variable. This means that y has already been considered by the procedure and has not been removed, so, due to Point 1, this reference will not be changed anymore. This is a contradiction.

Now let us assume that $\bar{\beta}$ contains a non-basic variable definition $x\{\gamma\}$, and let us first consider the case that $x \in \mathcal{X}_s$. Due to Point 2 from above, the original definition of x in $\bar{\alpha}$ must be basic, and due to Point 3, it must have the form $x\{z\}$. If, in the process of the procedure, the definition of x is changed in such a way that references are replaced by concatenations of references (i. e., it is changed into a non-basic definition), then, due to Point 4, it will necessarily be considered by the procedure and the main modification step is applied to it, which removes it. This is a contradiction.

Next, let us assume that $x \notin \mathcal{X}_s$, so x is introduced in some application of the main modification step. This means that initially its definition is basic and has the form $x\{z\}$. By assumption, its definition is not basic after termination of the procedure, thus, it must be changes such that z is replaced by a concatenation of references. Due to Point 3, this means that $x \in \mathcal{X}_s$, which is a contradiction.

This shows that after the procedure terminates, there are no non-basic variable definitions and therefore $\bar{\beta}$ is in normal form. We next estimate the size of $\bar{\beta}$, and the time and space required to construct it.

In the procedure that transforms $\bar{\alpha}$ into $\bar{\beta}$, there are $k \leq |\mathcal{X}_s|$ applications of the main modification step. For every $i \in [k] \cup \{0\}$, let $\bar{\alpha}^{(i)}$ be the version of $\bar{\alpha}$ after the i^{th} application of the main modification step in the procedure that transforms $\bar{\alpha}$ into $\bar{\beta}$. In particular, $\bar{\alpha}^{(0)} = \bar{\alpha}$ and $\bar{\alpha}^{(k)} = \bar{\beta}$. For every $x \in \mathcal{X}_s$, let k_x be the number of references of x in $\bar{\alpha}$. We claim that, for every $i \in [k] \cup \{0\}$, $|\bar{\alpha}^{(i)}| = O(|\bar{\alpha}|^{i+1})$. For $i = 0$, this obviously holds. Now let $i \in [k-1] \cup \{0\}$ and assume that $|\bar{\alpha}^{(i)}| = O(|\bar{\alpha}|^{i+1})$. Moreover, let the $(i+1)^{\text{th}}$ application of the main modification step apply to variable x . This means that $\bar{\alpha}^{(i+1)}$ is obtained in the $(i+1)^{\text{th}}$ application of the main modification step by replacing k_x symbols in $\bar{\alpha}^{(i)}$ by at most $|\bar{\alpha}^{(i)}|$ symbols. Hence, $|\bar{\alpha}^{(i+1)}| = O(k_x |\bar{\alpha}^{(i)}|) = O(|\bar{\alpha}| |\bar{\alpha}^{(i)}|) = O(|\bar{\alpha}| |\bar{\alpha}|^{i+1}) = O(|\bar{\alpha}|^{i+2})$. Consequently, $|\bar{\beta}| = O(|\bar{\alpha}|^{k+1}) = O(|\bar{\alpha}|^{|\mathcal{X}_s|+1})$.

It can be easily verified that $\bar{\beta}$ can also be computed in space $O(|\bar{\beta}|)$ and in time $O^*(|\bar{\beta}|)$

□

We are now sufficiently prepared to give a proof of the upper bound claimed in Theorem 2 in the next subsection.

5.2 Proof of Theorem 2

Theorem 2 is a consequence from the following lemma.

Lemma 7. *Given a Boolean $q \in \text{CXR PQ}^{\text{vsf}}$ and a graph database \mathcal{D} , we can nondeterministically check whether $\mathcal{D} \models q$ in space $O(2^{\text{poly}(|q|)} \log(|\mathcal{D}|))$.*

Proof. Let $q \in \text{CXR PQ}^{\text{vsf}}$ be Boolean and represented by the graph pattern G_q with $E_{\mathcal{D}} = \{(x_i, \alpha_i, y_i) \mid i \in [m]\}$, and let \mathcal{D} be a graph-database. We define a nondeterministic procedure to check whether or not $\mathcal{D} \models q$. First, for every $i \in [m]$, we transform α_i as follows. As long as α_i is not valt-free, we choose

some subexpression $(\gamma_1 \vee \gamma_2)$ where γ_1 or γ_2 contains a variable definition or a variable reference, and we nondeterministically replace it by γ_1 or γ_2 . In this way, we obtain a variable-simple conjunctive xregex $\bar{\alpha}' = (\alpha'_1, \alpha'_2, \dots, \alpha'_m)$. Moreover, this can clearly be done in space $O(|q|)$, and, for every $\bar{w} \in (\Sigma^*)^m$, if $\bar{w} \in \mathcal{L}(\bar{\alpha})$, then it is possible to perform the nondeterministic guesses in such a way that also $\bar{w} \in \mathcal{L}(\bar{\alpha}')$. In order to see this, it is sufficient to observe that this nondeterministic construction is similar to how vstar-free conjunctive xregex are made variable-simple in the proof of Lemma 4; the difference is that we cannot afford to explicitly store all possibilities of resolving alternations with variable references or definitions, so they need to be nondeterministically guessed. In particular, this implies that if $\mathcal{D} \models q$, then it is possible to perform the nondeterministic choices such that $\mathcal{D} \models q'$, where q' is the CXRPQ^{vsf} obtained from q by replacing each α_i by α'_i (note that this uses Proposition 2).

Next, we transform $\bar{\alpha}'$ into $\bar{\alpha}''$ according to the constructions used in the proof of Lemma 5 and we note that $|\bar{\alpha}''| = O(|\bar{\alpha}'|^2)$. Then, we transform $\bar{\alpha}''$ into $\bar{\beta} = (\beta_1, \beta_2, \dots, \beta_m)$ according to the construction used in the proof of Lemmas 6 and we note that $|\bar{\beta}| = O(|\bar{\alpha}''|^{|\bar{\alpha}''|+1}) = O((|\bar{\alpha}'|^2)^{|\bar{\alpha}'|^2+1}) = O(2^{\text{poly}(|\bar{\alpha}'|)})$. By q'' we denote the query obtained from q' by replacing each α'_i by β_i . We note that $|q''| = O(2^{\text{poly}(|q|)})$ and, according to Lemmas 5 and 6, $q' \equiv q''$. Thus, if $\mathcal{D} \models q$, then it is also possible to perform the initial nondeterministic guesses in such a way that also $\mathcal{D} \models q''$.

Next, we use Lemma 3 to nondeterministically decide whether $\mathcal{D} \models q''$ in space

$$O(|q''| \log(|\mathcal{D}|) + |q''| \log(|q''|)) = O(2^{\text{poly}(|q|)} \log(|\mathcal{D}|)).$$

Thus, this whole procedure decides nondeterministically whether $\mathcal{D} \models q$ in space $O(2^{\text{poly}(|q|)} \log(|\mathcal{D}|))$. \square

5.3 PSpace Combined-Complexity

Step 3 of the normal form construction (Lemma 6) works for $\bar{\alpha}$, where, for every $i \in [m]$, $\bar{\alpha}[i]$ is an alternation of variable-simple xregex and every $x \in \mathcal{X}_s$ has at most one definition. However, in the proof of Lemma 7, when we apply Step 3, we can make the much stronger assumption that $\bar{\alpha}$ is even variable-simple and every $x \in \mathcal{X}_s$ has at most one definition in $\bar{\alpha}$. Unfortunately, as illustrated by the following example, this does not make any difference with respect to the possible exponential size blow-up caused by Step 3.

$$\alpha = x_1 \{ \mathbf{a} \} x_2 \{ x_1 x_1 \} x_3 \{ x_2 x_2 \} x_4 \{ x_3 x_3 \} \dots x_n \{ x_{n-1} x_{n-1} \}$$

is a conjunctive xregex that is variable-simple and every variable has at most one definition. However, the procedure of Lemma 6 will apply the main modification step with respect to all variables x_2, x_3, \dots, x_n in this order (note that $G_{\bar{\alpha}}$ is a path), which replaces each references of x_2 by 2 variables, each references of x_4 by 4 variables, each references of x_5 by 8 variables, and so on. For example, the first application of the modification step will produce

$$x_1 \{ \mathbf{a} \} u_1 \{ x_1 \} u_2 \{ x_1 \} x_3 \{ (u_1 u_2)^2 \} x_4 \{ x_3 x_3 \} \dots x_n \{ x_{n-1} x_{n-1} \}$$

The crucial point seems to be non-basic definitions for variables with references in other non-basic definitions. If we restrict vstar-free conjunctive xregex accordingly, then the exponential size blow-up does not occur in Step 3 of the normal form construction.

A variable $x \in \mathcal{X}_s$ is *flat* (in some $\bar{\alpha} \in \text{CXRE}_{\Sigma, \mathcal{X}_s}$) if its definition is basic, or it has no reference in any other definition. For example, let $\alpha_1 = \text{ub}^*x\{y\{\mathbf{a}^*\}(\mathbf{a} \vee \mathbf{b})^*zy\}$ and $\alpha_2 = \text{u}\{\text{cbz}\{\mathbf{a}^*(\mathbf{b} \vee \mathbf{c}\mathbf{a})\}\}\mathbf{a}x$, then in $\bar{\alpha} = (\alpha_1, \alpha_2)$ every variable is flat. Finally, let $\text{CXRPQ}^{\text{vsf}, \text{fl}}$ be the class of vstar-free CXRPQ with only flat variables.

We next show that if a conjunctive xregex satisfies the preconditions of Lemma 6 *and* only has flat variables, then Step 3 of the normal form construction does not cause an exponential size blow-up.

Lemma 8. *Let $\bar{\alpha} \in m\text{-CXRE}_{\Sigma, \mathcal{X}_s}$ such that, for every $i \in [m]$, $\bar{\alpha}[i]$ is an alternation of variable-simple xregex, every $x \in \mathcal{X}_s$ has at most one definition in $\bar{\alpha}$, and all variables are flat. Then there is an equivalent $\bar{\beta} \in m\text{-CXRE}_{\Sigma, \mathcal{X}'_s}$ in normal form with $|\bar{\beta}| = O(|\bar{\alpha}|^2)$.*

Proof. Since the conjunctive xregex $\bar{\alpha}$ satisfies the conditions of the statement of Lemma 6, we can transform $\bar{\alpha}$ to $\bar{\beta}$ in exactly the same way as in the proof of Lemma 6. In order to conclude the proof, we only have to show that if all variables of $\bar{\alpha}$ are flat, then $|\bar{\beta}| = O(|\bar{\alpha}|^2)$.

Let $x \in \mathcal{X}_s$. If the definition of x is basic, then it is not changed by the procedure that changes $\bar{\alpha}$ to $\bar{\beta}$. If the definition of x is not basic, then, by assumption, it has no reference in any other variable definition. Thus, when the main modification step is applied to x , then the size of all other variable definitions does not increase. Consequently, in the procedure that changes $\bar{\alpha}$ to $\bar{\beta}$, every reference of x is not changed if the definition of x is basic, and it is replaced by at most $|\bar{\alpha}|$ other variable references, if it is not basic. This implies that, in the worst case, every symbol of $\bar{\alpha}$ can only be replaced by at most $|\bar{\alpha}|$ symbols, which means that $|\bar{\beta}| = O(|\bar{\alpha}|^2)$. \square

With this lemma, and the observation that Steps 1 and 2 of the normal form construction preserve flat variables, we can now show that the space upper bound for Boolean evaluation of vstar-free conjunctive xregex path queries is polynomial.

Lemma 9. *Given a Boolean $q \in \text{CXRPQ}^{\text{vsf}}$ with only flat variables and a graph database \mathcal{D} , we can nondeterministically check whether $\mathcal{D} \models q$ in space $O(\text{poly}(|q|) \log(|\mathcal{D}|))$.*

Proof. We can proceed analogously to the proof of Lemma 7. We only have to note that the initial nondeterministic transformation as well as the procedure of Lemma 5 preserves the property of having only flat variables. Consequently, the application of the procedure Lemma 6 yields a CXRPQ in normal form that is of size polynomial in the initial query. \square

Finally, let us summarise all these observations more concisely as follows.

Theorem 5. *$\text{CXRPQ}^{\text{vsf}, \text{fl}}\text{-BOOL-EVAL}$ is PSpace-complete in combined complexity.*

6 CXRPQ with Bounded Image Size

The fragment to be defined in this section directly follows from the definition of the subsets $\mathcal{L}^{\leq k}(\alpha) \subseteq \mathcal{L}(\alpha)$ that contain the words that match xregex α with

a witness v such that $|\text{vmap}_v(x)| \leq k$ for every $x \in \mathcal{X}_s$. This can be easily done as follows.

Let $\alpha \in \text{XRE}_{\Sigma, \mathcal{X}_s}$ with $\mathcal{X}_s = \{x_1, x_2, \dots, x_n\}$ and $\bar{v} \in (\Sigma^*)^n$. We define $\mathcal{L}_{\text{ref}}^{\bar{v}}(\alpha) = \{u \in \mathcal{L}_{\text{ref}}(\alpha) \mid \text{vmap}_{u, \mathcal{X}_s} = \bar{v}\}$ and $\mathcal{L}^{\bar{v}}(\alpha) = \text{deref}(\mathcal{L}_{\text{ref}}^{\bar{v}}(\alpha))$. For every $k \geq 1$, let $\mathcal{L}_{\text{ref}}^{\leq k}(\alpha) = \bigcup_{\bar{v} \in (\Sigma^{\leq k})^n} \mathcal{L}_{\text{ref}}^{\bar{v}}(\alpha)$. Finally, we define $\mathcal{L}^{\leq k}(\alpha) = \text{deref}(\mathcal{L}_{\text{ref}}^{\leq k}(\alpha))$.

The notions $\mathcal{L}_{\text{ref}}^{\bar{v}}$ and $\mathcal{L}_{\text{ref}}^{\leq k}$ also extend to conjunctive xregex in the obvious way, and therefore to CXRPQ as follows. For a $q \in \text{CXRPQ}$ with conjunctive xregex $\bar{\alpha} \in \text{CXRE}_{\Sigma, \mathcal{X}_s}$, a $\bar{v} \in (\Sigma^*)^n$, and a graph database \mathcal{D} , by $q^{\bar{v}}(\mathcal{D})$ we denote the subset of $q(\mathcal{D})$ that contains those $q_h(\mathcal{D}) \in q(\mathcal{D})$, where h is a matching morphism with respect to some matching words $\bar{w} = (w_1, w_2, \dots, w_m) \in \mathcal{L}^{\bar{v}}(\bar{\alpha})$; $q^{\leq k}(\mathcal{D})$ is defined analogously by restricting the matching words to be from $\mathcal{L}^{\leq k}(\bar{\alpha})$. In the Boolean case, we also set $\mathcal{D} \models^{\bar{v}} q$ and $\mathcal{D} \models^{\leq k} q$ to denote that $q^{\bar{v}}(\mathcal{D})$ or $q^{\leq k}(\mathcal{D})$, respectively, contains the empty tuple.

The above defined restrictions do not restrict the class CXRPQ (as it is the case for $\text{CXRPQ}^{\text{vsf}}$ considered in Section 5), but rather how the results of queries from CXRPQ are defined. However, it shall be convenient to allow a slight abuse of notation and define, for every $k \in \mathbb{N}$, the class $\text{CXRPQ}^{\leq k}$, which is the same as CXRPQ, but for every $q \in \text{CXRPQ}^{\leq k}$ it is understood that $q(\mathcal{D})$ actually means $q^{\leq k}(\mathcal{D})$ (and $\mathcal{D} \models q$ actually means $\mathcal{D} \models^{\leq k} q$). In this way, for every $k \in \mathbb{N}$, also the problems $\text{CXRPQ}^{\leq k}\text{-BOOL-EVAL}$ and $\text{CXRPQ}^{\leq k}\text{-CHECK}$ are defined.

For the classes $\text{CXRPQ}^{\leq k}$, we can show the following upper complexity bounds (which shall be proven in Subsection 6.1).

Theorem 6. *For every $k \in \mathbb{N}$, $\text{CXRPQ}^{\leq k}\text{-BOOL-EVAL}$ is in NP with respect to combined-complexity and in NL with respect to data-complexity.*

Since matching lower bounds directly carry over from CRPQs, the evaluation complexity of $\text{CXRPQ}^{\leq k}$ and CRPQ seems to be the same. However, we can state a much stronger hardness result, which points out an important difference between $\text{CXRPQ}^{\leq k}$ and CRPQ in terms of evaluation complexity: for $\text{CXRPQ}^{\leq 1}$, we have NP-hardness in combined-complexity even for single-edge graph patterns (with even simple xregex). This is not the case for CRPQ, which can be evaluated in polynomial-time if the structure of the underlying graph pattern is acyclic (see [10, 8, 6]).

Theorem 7. *$\text{CXRPQ}^{\leq k}\text{-BOOL-EVAL}$ is*

- *NP-hard in combined-complexity, even for $k = 1$ and single-edge queries with simple xregex $\alpha \in \text{XRE}_{\Sigma, \mathcal{X}_s}$ and $|\Sigma| = 3$,*
- *NL-hard in data-complexity, even for $k = 0$ and for single-edge queries with xregex α , where $\alpha \in \text{RE}_{\Sigma}$ with $|\Sigma| = 2$.*

Proof. Since $\text{CRPQ} \subseteq \text{CXRPQ}^{\leq 1}$, the NL-hardness for data-complexity follows in the same way as in the proof for Theorem 3.

In order to prove the NP-hardness for combined-complexity, we devise a reduction from the problem HITTING SET, which is defined as follows. Given subsets A_1, A_2, \dots, A_m of some universe U and $k \in \mathbb{N}$, decide whether there is a *hitting set* of size at most k , i. e., a set $B \subseteq U$ with $|B| \leq k$ and $B \cap A_i \neq \emptyset$ for every $i \in [m]$.

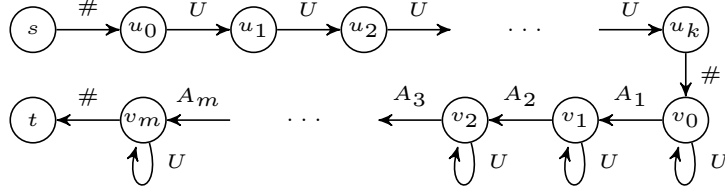


Figure 4: Sketch of the HITTING SET Reduction from Theorem 7. An arc labelled with U or A_i stands for arcs labelled by $\langle z \rangle$ for every $z \in U$ or $z \in A_i$, respectively.

Now let $A_1, A_2, \dots, A_m \subseteq U = \{z_1, z_2, \dots, z_n\}$ and $k \in \mathbb{N}$ be an instance of HITTING SET. Let $\Sigma = \{\mathbf{a}, \mathbf{b}, \#\}$ and, for every $z_i \in U$, we define $\langle z_i \rangle = \mathbf{b}\mathbf{a}^i\mathbf{b}$.

Next, we define a graph database $\mathcal{D} = (V_{\mathcal{D}}, E_{\mathcal{D}})$ over Σ as follows. For the sake of convenience, we also allow words as edge labels in graph databases, which stand for paths in the obvious way, i. e., for some $w \in \Sigma^*$, an arc (u, w, v) in a graph-database represents a path from u to v labelled with w (using some distinct intermediate nodes). We set

$$V_{\mathcal{D}} = \{s, u_0, u_1, \dots, u_k, v_0, v_1, \dots, v_m, t\}$$

and the set $E_{\mathcal{D}}$ is defined as follows.

- There are arcs $(s, \#, u_0)$, $(u_k, \#, v_0)$ and $(v_m, \#, t)$,
- For every $i \in [k]$ and $z \in U$ there is an arc $(u_{i-1}, \langle z \rangle, u_i)$,
- for every $i \in [m]$ and $z \in A_i$ there is an arc $(v_{i-1}, \langle z \rangle, v_i)$,
- for every $i \in [m] \cup \{0\}$ and every $z \in U$ there is an arc $(v_i, \langle z \rangle, v_i)$.

See Figure 4 for an illustration of \mathcal{D} . Next, we define the xregex

$$\alpha = \# \prod_{i=1}^{(n+2)k} x_i \{\mathbf{a} \vee \mathbf{b} \vee \varepsilon\} \# \left(\prod_{i=1}^{(n+2)k} x_i \right)^m \# ,$$

and the Boolean $q \in \text{CXRPQ}$ is defined by the single edge graph pattern $\{\{x, y\}, (x, \alpha, y)\}$. It is crucial to note that, for every $k \geq 1$, $\mathcal{L}^{\leq k}(\alpha) = \mathcal{L}(\alpha)$, which is due to the fact that for every possible ref-word in $\mathcal{L}_{\text{ref}}(\alpha)$, each variable image is either \mathbf{a} , \mathbf{b} or ε . Consequently, it does not matter for what k we interpret q to be a $\text{CXRPQ}^{\leq k}$.

In order to prove the correctness, we first make some observations:

1. Due to the occurrences of $\#$, there is a path in \mathcal{D} from a node y to a node z labelled with a word from $\mathcal{L}(\alpha)$ if and only if there is such a path in \mathcal{D} from s to t .
2. The language $\mathcal{L}(\alpha)$ contains exactly the words $w = \#w_1\#w_2\#$, where $w_1 \in \{\mathbf{a}, \mathbf{b}\}^*$ with $|w_1| \leq (n+2)k$, and $w_2 = (w_1)^m$.

3. Every path in \mathcal{D} from s to t is labelled by $w = \#w_1\#w_2\#$, where $w_1 = \langle z_{j_1} \rangle \langle z_{j_2} \rangle \dots \langle z_{j_k} \rangle$, for some $\{j_1, j_2, \dots, j_k\} \subseteq [n]$, and

$$w_2 = u_1 \langle z_{r_1} \rangle u_2 \langle z_{r_2} \rangle u_3 \dots u_m \langle z_{r_m} \rangle u_{m+1},$$

such that $\{r_1, r_2, \dots, r_m\} \subseteq [n]$ and, for every $i \in [m]$, $z_{r_i} \in A_i$; in particular, this means that $\{z_{r_1}, z_{r_2}, \dots, z_{r_m}\}$ is a hitting set (with respect to the considered problem-instance).

Next, we assume that in \mathcal{D} there is a path p labelled with a word $w \in \mathcal{L}(\alpha)$. With Point 1, we conclude that p is a path from s to t . Moreover, with Point 3, we have that $w = \#w_1\#w_2\#$, where $w_1 = \langle z_{j_1} \rangle \langle z_{j_2} \rangle \dots \langle z_{j_k} \rangle$ and $w_2 = u_1 \langle z_{r_1} \rangle u_2 \langle z_{r_2} \rangle u_3 \dots u_m \langle z_{r_m} \rangle u_{m+1}$, such that $\{z_{r_1}, z_{r_2}, \dots, z_{r_m}\}$ is a hitting set. Finally, Point 2 means that $w_2 = (w_1)^m$, which directly implies that $\{z_{r_1}, z_{r_2}, \dots, z_{r_m}\} \subseteq \{z_{j_1}, z_{j_2}, \dots, z_{j_k}\}$ and therefore $|\{z_{r_1}, z_{r_2}, \dots, z_{r_m}\}| \leq k$.

On the other hand, if $\{z_{j_1}, z_{j_2}, \dots, z_{j_k}\}$ is a hitting set of size k , then we can construct a word $w \in \mathcal{L}(\alpha)$ and a path from s to t that is labelled with w as follows. We set $w = \#w_1\#w_2\#$ with $w_1 = \langle z_{j_1} \rangle \langle z_{j_2} \rangle \dots \langle z_{j_k} \rangle$ and $w_2 = (w_1)^m$. We note that, according to Point 2, $w \in \mathcal{L}(\alpha)$ and we have to show that there is a path from s to t labelled with w . There is obviously a path from s to v_0 labelled with $\#w_1\#$. The set $\{z_{j_1}, z_{j_2}, \dots, z_{j_k}\}$ is a hitting set and each of the m occurrences of factor w_1 in w_2 contains a factor $\langle z_{j_i} \rangle$ for every $i \in [k]$. Thus, w_2 can be factorised into $w_2 = u_1 \langle z_{r_1} \rangle u_2 \langle z_{r_2} \rangle u_3 \dots u_m \langle z_{r_m} \rangle u_{m+1}$, such that, for every $i \in [m]$, $\langle z_{r_i} \rangle \in A_i$. This means that there is a path from v_0 to t labelled with $w_2\#$ and therefore a path from s to t labelled with w . \square

6.1 Proof of Theorem 6

The main building block of our algorithm is that if we have fixed some variable mapping $\bar{v} = (v_1, v_2, \dots, v_m)$, then the subset $\mathcal{L}^{\bar{v}}(\bar{\alpha})$ of $\mathcal{L}(\bar{\alpha})$ can be represented by a conjunctive xregex without variable definitions, i. e., a tuple of classical regular expressions (see Lemma 10). However, the corresponding procedure is not as simple as “replace each $x_i\{\dots\}$ and x_i by v_i ”. We shall now illustrate this with an example and some intuitive explanations. Let $\alpha = (\alpha_1, \alpha_2) \in \text{CXRE}_{\Sigma, \mathcal{X}_s}$ be defined by

$$\begin{aligned} \alpha_1 &= x_3 \{x_1 \{ca^*c\}x_2^*\} \vee [(x_1 \{cb^*\} \vee x_1 \{x_4c^*\})(b \vee x_2^*)x_3 \{x_1x_2x_1^*\}], \\ \alpha_2 &= (x_1 \vee x_2)^*x_4 \{(b \vee c)^*x_2^*\}x_2 \{(a \vee b)^*a\}, \end{aligned}$$

and let $\bar{v} = (v_1, \dots, v_4) = (ca, a, caaca, ca)$.

For computing $\beta = (\beta_1, \beta_2) \in \text{CXRE}_{\Sigma, \emptyset}$ with $\mathcal{L}(\bar{\beta}) = \mathcal{L}^{\bar{v}}(\bar{\alpha})$, replacing $x_i\{\gamma\}$ by v_i can only be correct, if γ can produce v_i (in a match with variable mapping \bar{v}). So we should treat the variable definitions and references of γ as their intended images and then check whether the thus obtained classical regular expression can produce v_i . For example, we transform $x_3 \{x_1 \{ca^*c\}x_2^*\}$ in α_1 to $x_3 \{v_1v_2^*\} = x_3 \{caa^*\}$ and check $v_3 = ca \in \mathcal{L}(caa^*)$. However, this assumes that $x_1 \{ca^*c\}$ can really produce v_1 , which is not the case, so we should rather remove $x_3 \{x_1 \{ca^*c\}x_2^*\}$ altogether, since it can never be involved in a conjunctive match from $\mathcal{L}^{\bar{v}}(\bar{\alpha})$. Hence, we also need to cut whole alternation branches in $\bar{\alpha}$, and, moreover, we have to make sure that if $x_i \neq \varepsilon$, then at least one definition of x_i

is necessarily instantiated, while this is not required if $x_i = \varepsilon$. Let us illustrate the correct transformations with this example.

For every variable definition $x_i\{\gamma\}$, where γ is a classical regular expression, we check whether $v_i \in \mathcal{L}(\gamma)$, and we mark the definition accordingly with 1 or 0, respectively. Then, we cut all branches that necessarily produce a definition marked with 0. This transform α_1 as follows

$$\begin{aligned}\alpha_1 &\rightsquigarrow x_3\{x_1\{0\}x_2^*\} \vee [(x_1\{0\} \vee x_1\{1\})(b \vee x_2^*)x_3\{1\}] \\ &\rightsquigarrow x_1\{1\}(b \vee x_2^*)x_3\{1\}.\end{aligned}$$

If there were variable definitions left, we would repeat this step, but for checking whether v_i can be generated by γ , we would treat variable definitions marked with 1 as their intended images. With respect to α_2 , we get

$$(x_1 \vee x_2)^*x_4\{(b \vee c)^*x_2^*\}x_2\{(a \vee b)^*a\} \rightsquigarrow (x_1 \vee x_2)^*x_4\{1\}x_2\{1\}.$$

Note that for $x_4\{(b \vee c)^*x_2^*\}$, we check whether $(b \vee c)^*(v_2)^* = (b \vee c)^*a^*$ can produce $v_4 = ca$, which is the case.

Finally, we replace all definitions and references by the intended images to obtain $(\beta_1, \beta_2) = (ca(b \vee a^*)caaca, ((ca) \vee a)^*caa)$.

In the general case, the situation can be slightly more complicated, since we also have to make sure that every ref-word necessarily instantiates a definition for x_i if $v_i \neq \emptyset$, while for $v_i = \emptyset$ ref-words without definition for x_i should still be possible. It can also happen that this procedure reduces an xregex to \emptyset , which means that $\mathcal{L}^{\bar{v}}(\bar{\alpha}) = \emptyset$.

These exemplary observations can be turned into a general procedure, which yields the following results.

Lemma 10. *For every $\bar{\alpha} \in m\text{-CXRE}_{\Sigma, \mathcal{X}_s}$ with $\mathcal{X}_s = \{x_1, x_2, \dots, x_n\}$ and every $\bar{v} \in (\Sigma^*)^n$, there is a $\bar{\beta} \in m\text{-CXRE}_{\Sigma, \emptyset}$ such that $\mathcal{L}(\bar{\beta}) = \mathcal{L}^{\bar{v}}(\bar{\alpha})$. Furthermore, $|\bar{\beta}| = O(|\bar{\alpha}|k)$, where $k = \max\{|\bar{v}[i]| \mid i \in [n]\}$, and $\bar{\beta}$ can be computed in time polynomial in $|\bar{\alpha}|$ and $|\bar{v}|$.*

Proof. Let $\bar{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_m)$ and let $\bar{v} = (v_1, v_2, \dots, v_n) \in \Sigma^*$. We give an algorithm that transforms $\bar{\alpha}$ into $\bar{\beta}$ with the desired property.

The general idea of the procedure is as follows. For every considered variable definition $x_i\{\gamma\}$ in some α_j , we want to determine whether γ can produce v_i under the assumption that all variable references and variable definitions in γ can produce their corresponding images (according to \bar{v}). If this is not the case, then we delete the alternation branch that instantiates the variable definition $x_i\{\gamma\}$ to make sure that it cannot be instantiated. When this procedure terminates, we can replace all remaining variable definitions and variable references by their corresponding images in order to obtain the components β_i of $\bar{\beta}$, which are then classical regular expressions. We now describe this procedure in more detail.

Step 1: We assume that the variable definitions in $\bar{\alpha}$ can be either *marked* or *unmarked* and that they are initially all *unmarked*. The algorithm considers each variable definition separately in such an order that every considered variable definition does only contain other variable definitions (if any) that are already marked. We repeat the following step until all variable definitions are *marked* (recall that initially all variable definitions are *unmarked*). Let $x_i\{\gamma\}$

be some *unmarked* variable definition in some α_j such that γ does not contain any *unmarked* variable definition. Let γ' be the classical regular expression obtained from γ by replacing each reference and definition for a variable $x_{i'}$ by $v_{i'}$. If $v_i \in \mathcal{L}(\gamma')$, then we *mark* $x_i\{\gamma\}$. If $v_i \notin \mathcal{L}(\gamma')$, then we have to modify α_j in such a way that $x_i\{\gamma\}$ is never instantiated. This is achieved as follows. We start at the node of the syntax-tree of α_j that represents $x_i\{\gamma\}$, we move up in the syntax tree and simply delete every node that we encounter (including the node that represents $x_i\{\gamma\}$ where we started, which also means that the whole subtree rooted by this node is deleted), and we stop as soon as we encounter an alternation node, which is then replaced by its other child (i. e., the sibling of the node from which we entered the alternation node). In particular, if no alternation node is encountered, then we can conclude that the definition $x_i\{\gamma\}$ under consideration will necessarily be instantiated by every ref-word of α_j , which immediately implies that there is no conjunctive match of $\bar{\alpha}$ with variable mapping (v_1, v_2, \dots, v_n) . In this case, the procedure will replace α_j by \emptyset (the unique regular expressions with $\mathcal{L}(\emptyset) = \emptyset$). Let $\bar{\alpha}' = (\alpha'_1, \alpha'_2, \dots, \alpha'_m)$ be the conjunctive xregex obtained when the procedure of this step terminates.

Step 2: Next, we have to modify $\bar{\alpha}'$ with respect to every $i \in [n]$ with $v_i \neq \varepsilon$ as follows. If, for some $j \in [m]$, α'_j contains a definition of x_i (note that this is possible for at most one α'_j), then we have to modify it such that it necessarily instantiates a definition for x_i , which is done as follows. For every node of the syntax tree for α'_j that corresponds to a definition for x_i , we mark it and then we move from this node up to the root and mark every visited node along the way. Next, for every marked alternation-node, we remove its unmarked child nodes (note that every such marked alternation-node has either one or two marked child nodes).

In particular, we note that this modification is only necessary for $i \in [n]$ with $v_i \neq \emptyset$, since ref-words that have no definition for x_i correspond to variable mappings with image ε for variable x_i , which should not be excluded if $v_i = \varepsilon$. Let $\bar{\alpha}'' = (\alpha''_1, \alpha''_2, \dots, \alpha''_m)$ be the conjunctive xregex obtained when the procedure of this step terminates.

After these two modification steps, it is possible that, for some $i \in [m]$ with $v_i \neq \emptyset$, there is no definition of x_i in $\bar{\alpha}''$. If this is the case, we replace each α''_j by \emptyset .

Finally, for every $i \in [n]$, we replace each definition and each occurrence of x_i by v_i in order to obtain a $\bar{\beta} \in m\text{-CXRE}_{\Sigma, \emptyset}$. It can be verified with moderate effort that $\mathcal{L}(\bar{\beta}) = \mathcal{L}^{\bar{v}}(\bar{\alpha})$. Moreover, this procedure can obviously be carried out in time polynomial in $|\bar{\alpha}|$ and $|\bar{v}|$, and also $|\bar{\beta}| = O(|\bar{\alpha}|k)$, where $k = \max\{|\bar{v}[i]| \mid i \in [n]\}$. \square

The statement of Lemma 10 directly carries over from conjunctive xregex to CXRPQ as follows.

Lemma 11. *For every $q \in \text{CXRPQ}$ with conjunctive xregex $\bar{\alpha} \in m\text{-CXRE}_{\Sigma, \mathcal{X}_s}$ with $\mathcal{X}_s = \{x_1, x_2, \dots, x_n\}$ and every $\bar{v} \in (\Sigma^*)^n$, there is a $q' \in \text{CRPQ}$, such that, for every database \mathcal{D} , we have that $q^{\bar{v}}(\mathcal{D}) = q'(\mathcal{D})$. Furthermore, $|q'| = O(|q|k)$, where $k = \max\{|\bar{v}[i]| \mid i \in [n]\}$, and q' can be computed in time polynomial in $|q|$ and $|\bar{v}|$.*

Proof. According to Lemma 10, we can compute in time polynomial in $|\bar{\alpha}|$ and $|\bar{v}|$ (and therefore in time polynomial in $|q|$ and $|\bar{v}|$) a $\bar{\beta} \in m\text{-CXRE}_{\Sigma, \emptyset}$ such that $\mathcal{L}(\bar{\beta}) = \mathcal{L}^{\bar{v}}(\bar{\alpha})$. Thus, q' can be obtained from q by replacing each edge label α_i by β_i . In particular, we have $|\bar{\beta}| = O(|\bar{\alpha}|k)$, where $k = \max\{|\bar{v}[i]| \mid i \in [n]\}$, and therefore also $|q'| = O(|q|k)$. \square

We are now ready to give a formal proof for Theorem 6.

Proof. (of Theorem 6) Let $q \in \text{CXPQ}^{\leq k}$ be Boolean with a conjunctive xregex $\bar{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_m)$ over Σ and $\mathcal{X}_s = \{x_1, x_2, \dots, x_n\}$, and let \mathcal{D} be a graph-database. We note that $\mathcal{D} \models q$ if and only if there is a $\bar{v} = (v_1, v_2, \dots, v_n) \in (\Sigma^{\leq k})^n$ such that $\mathcal{D} \models^{\bar{v}} q$. Consequently, the following is a nondeterministic algorithm that checks whether $\mathcal{D} \models q$.

1. Nondeterministically guess some $\bar{v} = (v_1, v_2, \dots, v_n) \in (\Sigma^{\leq k})^n$.
2. Compute $q' \in \text{CRPQ}$ such that, for every database \mathcal{D} , we have that $q^{\bar{v}}(\mathcal{D}) = q'(\mathcal{D})$.
3. Check whether $\mathcal{D} \models q'$.

Point 2 can be done according to Lemma 11, and Point 3 can be done according to Lemma 1. It remains to show that this nondeterministic algorithm requires polynomial time in combined-complexity and logarithmic space in data-complexity.

We first note that Point 1 can be done in time $O(nk)$, which is polynomial in combined complexity. Moreover, the required space for this does only depend on q and k , which means that it is constant in data-complexity.

According to Lemma 11, q' can be computed in time that is polynomial in $|q|$ and $|\bar{v}|$ in combined-complexity, which, since k is a constant, is polynomial in $|q|$. Again, the required space for this does only depend on q and k , which means that it is constant in data-complexity.

According to Lemma 1, we can nondeterministically check $\mathcal{D} \models q'$ in time polynomial in $|q'|$ and $|\mathcal{D}|$ in combined complexity. Moreover, according to Lemma 11, $|q'| = O(|q|k) = O(|q|)$, which means that we can nondeterministically check $\mathcal{D} \models q'$ in time polynomial in $|q|$ and $|\mathcal{D}|$, so polynomial in data-complexity. Finally, we also note that Lemma 1 implies that $\mathcal{D} \models q'$ can be checked in nondeterministic space that is logarithmic in \mathcal{D} with respect to data-complexity. \square

6.2 Logarithmically Bounded Image Size

Analogously to $q^k(\mathcal{D})$ for every $k \geq 1$, $q \in \text{CXPQ}$ and graph database \mathcal{D} , we can also define $q^{\log}(\mathcal{D}) = \bigcup_{k=0}^{\log(|\mathcal{D}|)} q^k(\mathcal{D})$. Just like we derived the classes $\text{CXPQ}^{\leq k}$, this gives rise to the class CXPQ^{\log} . Note that a Boolean $q \in \text{CXPQ}^{\log}$ matches a graph database \mathcal{D} , if and only if there is a matching morphism with image size bounded by $\log(|\mathcal{D}|)$.

Analogously to the proof of Theorem 6, we can show the following upper bound for CXPQ^{\log} .

Corollary 1. *$\text{CXPQ}^{\log}\text{-BOOL-EVAL}$ can nondeterministically be solved in polynomial time in combined-complexity and in space $O(\log^2(|\mathcal{D}|))$ in data-complexity.*

Proof. We can apply the same nondeterministic algorithm from the proof of Theorem 6, with the only difference that we initially guess $\bar{v} = (v_1, v_2, \dots, v_n) \in (\Sigma^{\leq \log(|\mathcal{D}|)})^n$. More precisely, we use the following nondeterministic algorithm:

1. Nondeterministically guess some $\bar{v} = (v_1, v_2, \dots, v_n) \in (\Sigma^{\leq \log(|\mathcal{D}|)})^n$.
2. Compute $q' \in \text{CRPQ}$ such that, for every database \mathcal{D} , we have that $q^{\bar{v}}(\mathcal{D}) = q'(\mathcal{D})$.
3. Check whether $\mathcal{D} \models q'$.

Point 1 can be done in time and space $O(n \log(|\mathcal{D}|))$. According to Lemma 11, q' can be computed in time that is polynomial in $|q|$ and $|\bar{v}|$, which, since $|\bar{v}| = n \log(|\mathcal{D}|)$, is polynomial in $|q|$ and $|\mathcal{D}|$. Moreover, according to Lemma 11, $|q'| = O(|q| \log(|\mathcal{D}|))$.

We can use Lemma 1 in order to conclude that we can nondeterministically check $\mathcal{D} \models q'$ in polynomial time. With respect to space complexity, we note that since $q' \in \text{CRPQ}$, we also have that q' is a simple CXRPQ. Therefore, we can conclude with Lemma 3 that we can check $\mathcal{D} \models q'$ in space

$$\begin{aligned} & O(|q'| \log(|\mathcal{D}|) + |q'| \log(|q'|)) \\ &= O(|q| \log(|\mathcal{D}|) \log(|\mathcal{D}|) + |q| \log(|\mathcal{D}|) \log(|q| \log(|\mathcal{D}|))) \\ &= O(|q| \log(|q|) \log^2(|\mathcal{D}|)). \end{aligned}$$

□

7 Expressive Power

In this section, we compare the expressive power of CXRPQs and their fragments with other established classes of graph queries.

The *extended conjunctive regular path queries* (ECRPQs), already mentioned in the introduction, have been introduced in [8]. We shall define them now in more detail (a definition in full details can be found in [8]). Extended conjunctive regular path queries (ECRPQs) have the form $q = \bar{z} \leftarrow G_q, \bigwedge_{j \in [t]} R_j(\bar{\omega}_j)$, where $\bar{z} \leftarrow G_q$ is a CRPQ and, for every $j \in [t]$, $\bar{\omega}_j$ is an s_j -tuple over E_q and R_j is a regular expression that describes a regular relations over Σ^* of arity s_j . The semantics of ECRPQ can be derived from the semantics of CRPQ as follows. We interpret q as a CRPQ, but we add to the concept of a matching morphism the requirement that there must be a tuple $(w_{e_1}, w_{e_2}, \dots, w_{e_m})$ of matching words such that, for each $\bar{\omega}_j = (e_{p_1}, e_{p_2}, \dots, e_{p_{s_j}})$, we have $(w_{e_{p_1}}, w_{e_{p_2}}, \dots, w_{e_{p_{s_j}}}) \in \mathcal{L}(R_j)$.

By ECRPQ with *equality relations* (ECRPQ^{er} for short), we denote the class of ECRPQ for which each R_j is the *equality relation* (for some arity s_j), i. e., the relation $\{(u_1, u_2, \dots, u_{s_j}) \in (\Sigma^*)^{s_j} \mid u_1 = u_2 = \dots = u_{s_j}\}$. In order to ease our notations, we sometimes represent the equality relations as a partition $\{E_{q,1}, E_{q,2}, \dots, E_{q,t}\}$ of E_q (i. e., for every $j \in [t]$, the edges of $E_{q,j}$ are subject to an equality relation), or, for simple queries, we also state which edges are required to be equal without formally stating the equality relations.

We recall that for some conjunctive path query q , the mapping $\mathcal{D} \mapsto q(\mathcal{D})$ from the set of graph-databases to the set of relations over Σ that is defined

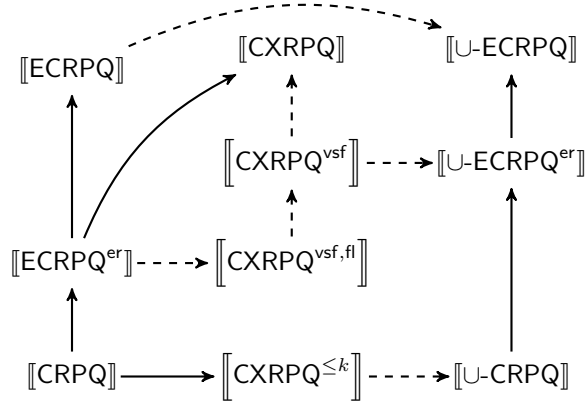


Figure 5: Illustration of the relations between the considered classes of conjunctive path queries. A dashed or solid arrow from A to B denotes $A \subseteq B$ or $A \subsetneq B$, respectively.

by q is denoted by $\llbracket q \rrbracket$, and for any class A of conjunctive path queries, $\llbracket A \rrbracket = \{\llbracket q \rrbracket \mid q \in A\}$.

For a class Q of conjunctive path queries, a *union of Q s* (or \cup - Q , for short) is a query of the form $q = q_1 \vee q_2 \vee \dots \vee q_k$, where, for every $i \in [k]$, $q_i \in Q$. For a graph-database \mathcal{D} , we define $q(\mathcal{D}) = \bigcup_{i \in [k]} q_i(\mathcal{D})$.

The results of this section can be summarised as follows.

Theorem 8. *The inclusion-diagram of Figure 5 is correct.*

Next, we first discuss these results in Subsection 7.1, and then, in the following subsections, we formally prove all the inclusions of Figure 5.

7.1 Discussion of Results

The diagram of Figure 5 shows three vertical layers of query classes of increasing expressive power. More precisely, on the left side we have the *classical* conjunctive path query classes CRPQ , ECRPQ^{er} and ECRPQ . Then, our new fragments of CXPQ follow. Finally, we have the classes of unions of the classical conjunctive path query classes. All these layers are naturally ordered by the subset relation, and these subset relations follow all directly by definition.

The more interesting inclusion relations are the vertical ones. The inclusion $\llbracket \text{ECRPQ}^{\text{er}} \rrbracket \subseteq \llbracket \text{CXPQ}^{\text{vsf, fl}} \rrbracket$ is as expected, but nevertheless points out that also quite restricted classes of CXPQ still cover CRPQ that are extended by equality relations. However, this does not seem to be the case for the classes $\llbracket \text{CXPQ}^{\leq k} \rrbracket$, which nevertheless contains the class $\llbracket \text{CRPQ} \rrbracket$.

Less obvious are the inclusions of the class $\llbracket \text{CXPQ}^{\text{vsf}} \rrbracket$ in $\llbracket \text{U-ERPQ}^{\text{er}} \rrbracket$, and of the class $\llbracket \text{CXPQ}^{\leq k} \rrbracket$ in $\llbracket \text{U-CRPQ} \rrbracket$. Both of them are more or less a result from the upper bounds for these CXPQ -fragments shown in Sections 5 and 6. In a sense, this means that the queries of our CXPQ -fragments can be “decomposed” into unions of the more classical CRPQ s and ECRPQ^{er} s. Thus, semantically, the CXPQ^{vsf} and $\text{CXPQ}^{\leq k}$ can still be described in the formalisms of

CRPQ and ECRPQ^{er}. However, there is a significant syntactical difference: our conversions of CXRPQ^{vsf} or CXRPQ^{≤k} into \cup -CRPQ or \cup -ECRPQ^{er}, respectively, require exponential (or even double exponential) size blow-ups.

It is to be expected, that CXRPQ^{vsf} are strictly more powerful than ECRPQ^{er}. However, we can only show that CXRPQ are strictly more powerful than ECRPQ^{er}. On the other hand, that also CXRPQ^{≤k}s are strictly more powerful than CRPQs seems less obvious, since the string variables of CXRPQ^{≤k}s, which syntactically set them apart from CRPQs, can only range over finite sets of possible images. So it is more surprising that, in fact, even CXRPQ^{≤1} contains queries that are not expressible as CRPQ.

7.2 $\llbracket \text{CRPQ} \rrbracket \subsetneq \llbracket \text{ECRPQ}^{\text{er}} \rrbracket \subsetneq \llbracket \text{ECRPQ} \rrbracket$ and $\llbracket \cup\text{-CRPQ} \rrbracket \subsetneq \llbracket \cup\text{-ECRPQ}^{\text{er}} \rrbracket \subsetneq \llbracket \cup\text{-ECRPQ} \rrbracket$

Theorem 9. $\llbracket \text{CRPQ} \rrbracket \subsetneq \llbracket \text{ECRPQ}^{\text{er}} \rrbracket \subsetneq \llbracket \text{ECRPQ} \rrbracket$.

Proof. The inclusions follow immediately, since CRPQ can be interpreted as ECRPQ^{er} without any equality relations, and ECRPQ^{er} \subseteq ECRPQ. We shall next prove that they are proper.

Let $\Sigma = \{\mathbf{a}, \mathbf{b}, \mathbf{c}, \mathbf{d}\}$ and let $q_{\mathbf{a}^n \mathbf{b}^n}$ be the Boolean ECRPQ over Σ defined by the graph-pattern $G_{q_{\mathbf{a}^n \mathbf{b}^n}} = (V_{q_{\mathbf{a}^n \mathbf{b}^n}}, E_{q_{\mathbf{a}^n \mathbf{b}^n}})$ with

$$\begin{aligned} V_{q_{\mathbf{a}^n \mathbf{b}^n}} &= \{x, y_1, y_2, z, x', y'_1, y'_2, z'\}, \\ E_{q_{\mathbf{a}^n \mathbf{b}^n}} &= \{(x, \mathbf{c}, y_1), (y_1, \mathbf{a}^*, y_2), (y_2, \mathbf{c}, z), \\ &\quad (x', \mathbf{d}, y'_1), (y'_1, \mathbf{b}^*, y'_2), (y'_2, \mathbf{d}, z')\}, \end{aligned}$$

and with only one equal-length relation (i. e., the relation $\{(u_1, u_2) \in (\Sigma^*)^2 \mid |u_1| = |u_2|\}$) that applies to the edges (y_1, \mathbf{a}^*, y_2) and $(y'_1, \mathbf{b}^*, y'_2)$. See Figure 6 for an illustration. We note that $\llbracket q_{\mathbf{a}^n \mathbf{b}^n} \rrbracket$ is the set of graph-databases \mathcal{D} that contain (not necessarily distinct) vertices u, v, u', v' and a path from u to v labelled with $\mathbf{ca}^n \mathbf{c}$ and a path from u' to v' labelled with $\mathbf{db}^n \mathbf{d}$, respectively, for some $n \geq 0$.

Claim 1: $\llbracket q_{\mathbf{a}^n \mathbf{b}^n} \rrbracket \notin \llbracket \text{ECRPQ}^{\text{er}} \rrbracket$.

Proof of Claim 1: For the sake of convenience, we relabel $q_{\mathbf{a}^n \mathbf{b}^n}$ to q in the proof of the claim. We assume that there is a Boolean $q' \in \text{ECRPQ}^{\text{er}}$, such that $\llbracket q' \rrbracket = \llbracket q \rrbracket$. Moreover, let q' be defined by a graph pattern $G_{q'} = (V_{q'}, E_{q'})$ with $E_{q'} = \{(\hat{x}_i, \alpha_i, \hat{y}_i) \mid i \in [m]\}$ and some equality relations.

For every $n \in \mathbb{N}$, let $\mathcal{D}_{n,n}$ be the graph-database given by two node-disjoint paths $(r_0, r_1, \dots, r_{n+2})$ and $(s_0, s_1, \dots, s_{n+2})$ labelled with $\mathbf{ca}^n \mathbf{c}$ and $\mathbf{db}^n \mathbf{d}$, respectively. Obviously, for every $n \in \mathbb{N}$, $\mathcal{D}_{n,n} \in \llbracket q \rrbracket = \llbracket q' \rrbracket$, which means that there is at least one matching morphism h for q' and $\mathcal{D}_{n,n}$. In the following, for every $n \in \mathbb{N}$, let h_n be some fixed matching morphism for q' and $\mathcal{D}_{n,n}$. By the structure of $\mathcal{D}_{n,n}$, we also know that, for every $i \in [m]$, the arc $(\hat{x}_i, \alpha_i, \hat{y}_i)$ is matched to some sub-path of either $(r_0, r_1, \dots, r_{n+2})$ or $(s_0, s_1, \dots, s_{n+2})$; more precisely, there are $\ell_i, \ell'_i \in [n+2] \cup \{0\}$ with $0 \leq \ell_i \leq \ell'_i \leq n+2$ such that either $h_n(\hat{x}_i) = r_{\ell_i}$ and $h_n(\hat{y}_i) = r_{\ell'_i}$, or $h_n(\hat{x}_i) = s_{\ell_i}$ and $h_n(\hat{y}_i) = s_{\ell'_i}$. Now let $\{C_n, D_n\}$ be a partition of $[m]$ such that $i \in C_n$ if $(\hat{x}_i, \alpha_i, \hat{y}_i)$ is matched to some sub-path of $(r_0, r_1, \dots, r_{n+2})$ and $i \in D_n$ if $(\hat{x}_i, \alpha_i, \hat{y}_i)$ is matched to some sub-path of $(s_0, s_1, \dots, s_{n+2})$. In particular, we note that $(r_0, r_1, \dots, r_{n+2})$ only contains labels \mathbf{a} and \mathbf{c} , while $(s_0, s_1, \dots, s_{n+2})$ only contains labels \mathbf{b} and \mathbf{d} . This

means that for each single equality relation of q' with arity p that applies to a set $\{e_{j_1}, e_{j_2}, \dots, e_{j_p}\}$ of arcs, there are three possibilities: (1) $\{j_1, j_2, \dots, j_p\} \subseteq C_n$, (2) $\{j_1, j_2, \dots, j_p\} \subseteq D_n$, or (3), for every $i \in [p]$, $h_n(\hat{x}_{j_i}) = h_n(\hat{y}_{j_i})$ (i. e., they cover paths of length 0 that are labelled with ε).

Since there is only a finite number of partitions of $[m]$ into two sets, there must be some $n_1, n_2 \in \mathbb{N}$ with $n_1 \neq n_2$ such that $C_{n_1} = C_{n_2}$ and $D_{n_1} = D_{n_2}$. We now define a morphism $h : V_{q'} \rightarrow \{r_0, r_1, \dots, r_{n_1+2}\} \cup \{s_0, s_1, \dots, s_{n_2+2}\}$ by setting, for every $i \in C_{n_1} = C_{n_2}$, $h(\hat{x}_i) = h_{n_1}(\hat{x}_i)$ and $h(\hat{y}_i) = h_{n_1}(\hat{y}_i)$, and, for every $i \in D_{n_1} = D_{n_2}$, $h(\hat{x}_i) = h_{n_2}(\hat{x}_i)$ and $h(\hat{y}_i) = h_{n_2}(\hat{y}_i)$. We note that h is a matching morphism for q' and \mathcal{D}_{n_1, n_2} . In particular, due to our observation from above, each equality relation is satisfied. Since $n_1 \neq n_2$, we have that $\mathcal{D}_{n_1, n_2} \notin \llbracket q' \rrbracket$, which is a contradiction. \square (*Claim 1*)

Let $q_{\mathbf{a}^n \mathbf{a}^n}$ be the ECRPQ^{er} over Σ that is defined in almost the same way as $q_{\mathbf{a}^n \mathbf{b}^n}$, with the only differences that we label the arc from y'_1 to y'_2 of the graph-pattern by \mathbf{a}^* instead of \mathbf{b}^* and that the binary equal-length relation on (y_1, \mathbf{a}^*, y_2) and $(y'_1, \mathbf{b}^*, y'_2)$ becomes a binary equality relation on (y_1, \mathbf{a}^*, y_2) and $(y'_1, \mathbf{a}^*, y'_2)$.

More formally, let $q_{\mathbf{a}^n \mathbf{a}^n} \in \text{ECRPQ}^{\text{er}}$ over $\Sigma = \{\mathbf{a}, \mathbf{b}, \mathbf{c}, \mathbf{d}\}$ be defined by the graph-pattern $G_{q_{\mathbf{a}^n \mathbf{a}^n}} = (V_{q_{\mathbf{a}^n \mathbf{a}^n}}, E_{q_{\mathbf{a}^n \mathbf{a}^n}})$ with

$$\begin{aligned} V_{q_{\mathbf{a}^n \mathbf{a}^n}} &= \{x, y_1, y_2, z, x', y'_1, y'_2, z'\}, \\ E_{q_{\mathbf{a}^n \mathbf{a}^n}} &= \{(x, \mathbf{c}, y_1), (y_1, \mathbf{a}^*, y_2), (y_2, \mathbf{c}, z), \\ &\quad (x', \mathbf{d}, y'_1), (y'_1, \mathbf{a}^*, y'_2), (y'_2, \mathbf{d}, z')\}, \end{aligned}$$

and only one binary equality relation that applies to the edges (y_1, \mathbf{a}^*, y_2) and $(y'_1, \mathbf{a}^*, y'_2)$. We note that $\llbracket q_{\mathbf{a}^n \mathbf{a}^n} \rrbracket$ is the set of graph-databases \mathcal{D} that contain (not necessarily distinct) vertices u, v, u', v' and a path from u to v labelled with $\mathbf{ca}^n \mathbf{c}$ and a path from u' to v' labelled with $\mathbf{da}^n \mathbf{d}$, respectively, for some $n \geq 0$.

Claim 2: $\llbracket q_{\mathbf{a}^n \mathbf{a}^n} \rrbracket \notin \llbracket \text{CRPQ} \rrbracket$.

Proof of Claim 2: We assume that there is a $q' \in \text{CRPQ}$, such that $\llbracket q' \rrbracket = \llbracket q \rrbracket$. Moreover, let q' be defined by a graph pattern $G_{q'} = (V_{q'}, E_{q'})$ with $E_{q'} = \{(\hat{x}_i, \alpha_i, \hat{y}_i) \mid i \in [m]\}$.

We can now obtain a contradiction analogously as in to the proof of Claim 1. In fact, the argument is almost the same, but we argue with graph-databases $\mathcal{D}_{n, n}$ given by two node-disjoint paths $(r_0, r_1, \dots, r_{n+2})$ and $(s_0, s_1, \dots, s_{n+2})$ labelled with $\mathbf{ca}^n \mathbf{c}$ and $\mathbf{da}^n \mathbf{d}$, respectively. In general, the argument is simpler, because we do not have to take care of possible equality relations.

\square (*Claim 2*)

This concludes the proof. \square

Theorem 10. $\llbracket \cup\text{-CRPQ} \rrbracket \subsetneq \llbracket \cup\text{-ECRPQ}^{\text{er}} \rrbracket \subsetneq \llbracket \cup\text{-ECRPQ} \rrbracket$.

Proof. The Inclusions follow by definition. We next show that the inclusion $\llbracket \cup\text{-ECRPQ}^{\text{er}} \rrbracket \subseteq \llbracket \cup\text{-ECRPQ} \rrbracket$ is proper. To this end, we first recall the proof of Claim 1 in the proof of Theorem 9, which showed that for the query $q = q_{\mathbf{a}^n \mathbf{b}^n} \in \text{ECRPQ} \subseteq \cup\text{-ECRPQ}$ (see also Figure 6), we have that $\llbracket q \rrbracket \notin \text{ECRPQ}^{\text{er}}$. We have demonstrated that if there is a query $q' \in \text{ECRPQ}^{\text{er}}$ with $\mathcal{D}_{n, n} \models q'$ for every $n \in \mathbb{N}$, then there are $n_1, n_2 \in \mathbb{N}$ with $n_1 \neq n_2$, such that q can be matched to both \mathcal{D}_{n_1, n_1} and \mathcal{D}_{n_2, n_2} in such a way that the partition of the edges of q'

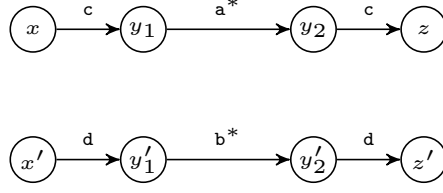


Figure 6: Illustration of the graph-pattern for $q_{a^n b^n}$.

according to whether they are matched to the $ca^n c$ path or the $db^n d$ path are exactly the same. This has lead to the contradiction that $\mathcal{D}_{n_1, n_2} \models q'$.

If we now instead consider a $q' \in \cup\text{-ECRPQ}^{\text{er}}$, then q' contains only a finite number of graph patterns $G_q^{(1)}, G_q^{(2)}, \dots, G_q^{(\ell)}$, and, for every $n \in \mathbb{N}$, there is at least one $j \in [\ell]$, such that $G_q^{(j)}$ can be matched to $\mathcal{D}_{n, n}$. This implies again, that there are $n_1, n_2 \in \mathbb{N}$ with $n_1 \neq n_2$, such that, for some $j \in [\ell]$, $G_q^{(j)}$ can be matched to both \mathcal{D}_{n_1, n_1} and \mathcal{D}_{n_2, n_2} in such a way that the partition of the edges of $G_q^{(j)}$ according to whether they are matched to the $ca^n c$ path or the $db^n d$ are exactly the same. Again, this leads to the contradiction that $G_q^{(j)}$ can be matched to \mathcal{D}_{n_1, n_2} and therefore $\mathcal{D}_{n_1, n_2} \models q'$.

In order to show that the inclusion $[\cup\text{-ECRPQ}^{\text{er}}] \subseteq [\cup\text{-ECRPQ}]$ is proper, we argue analogously, but with $q = q_{a^n a^n} \in \text{ECRPQ}^{\text{er}} \subseteq \cup\text{-ECRPQ}^{\text{er}}$ (i. e., we extend the argument of the proof of Claim 2 of the proof of Theorem 9 to the case of unions of queries just as it is done above with respect to Claim 1 of the proof of Theorem 9). \square

7.3 $[\text{ECRPQ}^{\text{er}}] \subseteq [\text{CXPQ}^{\text{vsf, fl}}] \subseteq [\text{CXPQ}^{\text{vsf}}] \subseteq [\text{CXPQ}]$

All inclusion of

$$[\text{ECRPQ}^{\text{er}}] \subseteq [\text{CXPQ}^{\text{vsf, fl}}] \subseteq [\text{CXPQ}^{\text{vsf}}] \subseteq [\text{CXPQ}]$$

follow by definition, except the first one, which is due to the following Lemma 12. Note that this inclusion chain also implies the inclusion $[\text{ECRPQ}^{\text{er}}] \subseteq [\text{CXPQ}]$ depicted in Figure 5; its strictness is shown later in Subsection 7.6.

Lemma 12. $[\text{ECRPQ}^{\text{er}}] \subseteq [\text{CXPQ}^{\text{vsf, fl}}]$.

Proof. Let $q \in \text{ECRPQ}^{\text{er}}$ be of the form $q = \bar{z} \leftarrow G_q$ with $E_q = \{e_i = (x_i, \alpha_i, y_i) \mid i \in [m]\}$ and let $\{E_{q,1}, E_{q,2}, \dots, E_{q,t}\}$ with $E_{q,j} = \{e_{p_1}, e_{p_2}, \dots, e_{p_{s_j}}\}$ be the partition of E_q that represents the equality constraints. For every $j \in [t]$, we successively modify q as follows. We replace $(x_{p_1}, \alpha_{p_1}, y_{p_1})$ by $(x_{p_1}, \beta, y_{p_1})$, where β is a regular expression for $\bigcap_{i=1}^{s_j} \mathcal{L}(\alpha_{p_i})$, and, for every i with $2 \leq i \leq p_{s_j}$, we replace $(x_{p_i}, \alpha_{p_i}, y_{p_i})$ by $(x_{p_i}, \Sigma^*, y_{p_i})$. We denote the ECRPQ^{er} constructed in this way by q' and we note that q' is equivalent to q . Moreover, q' is represented by a graph $(V_{q'}, E_{q'})$ and a partition $\{E_{q',1}, E_{q',2}, \dots, E_{q',t}\}$ such that, for every $j \in [t]$, $E_{q',j} = \{(x_{p_1}, \beta_{p_1}, y_{p_1}), (x_{p_2}, \Sigma^*, y_{p_2}), \dots, (x_{p_{s_j}}, \Sigma^*, y_{p_{s_j}})\}$ for some regular expression β_{j_1} . We can now translate q' into a $q'' \in \text{CXPQ}^{\text{vsf, fl}}$ by replacing, for every $j \in [t]$, edge $(x_{j_1}, \beta_{j_1}, y_{j_1})$ by $(x_{j_1}, z_j \{\beta_{j_1}\}, y_{j_1})$ and every

edge $(x_{j_\ell}, \Sigma^*, y_{j_\ell})$, $2 \leq \ell \leq s_j$, by $(x_{j_\ell}, z_j, y_{j_\ell})$. It can be easily verified that q'' is equivalent to q . □

7.4 $\llbracket \text{ECRPQ} \rrbracket \subseteq \llbracket \text{U-ECRPQ} \rrbracket$ and $\llbracket \text{CXPQ}^{\text{vsf}} \rrbracket \subseteq \llbracket \text{U-ECRPQ}^{\text{er}} \rrbracket$

The inclusion $\llbracket \text{ECRPQ} \rrbracket \subseteq \llbracket \text{U-ECRPQ} \rrbracket$ follows trivially by definition, whereas $\llbracket \text{CXPQ}^{\text{vsf}} \rrbracket \subseteq \llbracket \text{U-ECRPQ}^{\text{er}} \rrbracket$ is shown by the following Lemma 13.

Lemma 13. $\llbracket \text{CXPQ}^{\text{vsf}} \rrbracket \subseteq \llbracket \text{U-ECRPQ}^{\text{er}} \rrbracket$.

Proof. Let $q = \bar{z} \leftarrow G_q$ with xregex $\bar{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_m)$. By Lemmas 4, 5 and 6, we can assume that $\bar{\alpha}$ is in normal form, i. e., it is such that, for every $i \in [m]$, $\alpha_i = \alpha_{i,1} \vee \alpha_{i,2} \vee \dots \vee \alpha_{i,t_i}$ and every $\alpha_{i,j}$ for $j \in [t_i]$ is simple.

We can now transform q into $q_1, q_2, \dots, q_r \in \text{CXPQ}^{\text{vsf}}$, such that, for every graph-database \mathcal{D} , $q(\mathcal{D}) = \bigcup_{i=1}^r q_i(\mathcal{D})$, and, for every $i \in [r]$, q_i is simple. More precisely, the q_j with $j \in [r]$ are obtained by considering, for every $i \in [m]$, all possibilities of replacing α_i by exactly one of the $\alpha_{i,1}, \alpha_{i,2}, \dots, \alpha_{i,t_i}$, which obviously results in CXPQ^{vsf} that are simple (since q is in normal form).

Next, for every $i \in [r]$, we can transform q_i into an equivalent $q_i'' \in \text{ECRPQ}^{\text{er}}$ as follows. Let $\bar{\alpha}^{(i)} = (\alpha_1^{(i)}, \alpha_2^{(i)}, \dots, \alpha_m^{(i)})$ be the conjunctive xregex of q_i . We first observe that, analogously to the proof of Lemma 3, since $\bar{\alpha}^{(i)}$ is simple, we can replace definitions $x\{y\}$ and all occurrences of x by references of y without changing the set of conjunctive matches. Let the thus modified version of $\bar{\alpha}^{(i)}$ be denoted by $\bar{\alpha}'^{(i)} = (\alpha_1'^{(i)}, \alpha_2'^{(i)}, \dots, \alpha_m'^{(i)})$. For every $j \in [m]$, $\alpha_j'^{(i)} = \pi_1 \pi_2 \dots \pi_t$, where each π_ℓ is a classical regular expression γ , a variable definition $x\{\gamma\}$, where γ is a classical regular expression, or a variable reference x . Thus, for every $j \in [m]$, we can break up the edge labelled with $\alpha_j'^{(i)}$ into a path of size t with the edge labels π_ℓ . If we do this for all $\alpha_j'^{(i)}$, then we have turned q_i into a $q_i' \in \text{CXPQ}^{\text{vsf}}$, such that every edge is labelled by a classical regular expression, a variable definition over a classical regular expression, or a variable reference. For every variable x , we now do the following. We replace the edge label $x\{\gamma\}$ by γ (since q_i' is simple, there can be at most one such edge label) and all edge labels x by Σ^* and add an equality constraint that applies to exactly the edges modified by this step. It can be easily seen that the thus obtained $q_i'' \in \text{ECRPQ}^{\text{er}}$ is equivalent to q_i (note that the tuple of output nodes \bar{z} remains unchanged). □

7.5 $\llbracket \text{CXPQ}^{\leq k} \rrbracket \subseteq \llbracket \text{U-CRPQ} \rrbracket$ and $\llbracket \text{CRPQ} \rrbracket \subseteq \llbracket \text{CXPQ}^{\leq k} \rrbracket$

That, for every $k \geq 1$, $\llbracket \text{CXPQ}^{\leq k} \rrbracket \subseteq \llbracket \text{U-CRPQ} \rrbracket$, is due to the following Lemma 14. Moreover, that, for every $k \geq 1$, $\llbracket \text{CRPQ} \rrbracket \subseteq \llbracket \text{CXPQ}^{\leq k} \rrbracket$, follows trivially by definition and the strictness of the inclusions is shown in Subsection 7.6.

Lemma 14. For every $k \geq 1$, $\llbracket \text{CXPQ}^{\leq k} \rrbracket \subseteq \llbracket \text{U-CRPQ} \rrbracket$.

Proof. Let $k \geq 1$ and let $q \in \text{CXPQ}^{\leq k}$ be defined by a graph pattern $G_q = (V_q, E_q)$ with $E_q = \{(x_i, \alpha_i, y_i) \mid i \in [m]\}$, where $\bar{\alpha} \in \text{CXRE}_{\Sigma, \mathcal{X}_s}$ with $\mathcal{X}_s = \{x_1, x_2, \dots, x_n\}$. Now, for every $\bar{v} \in (\Sigma^{\leq k})^n$, let $q[\bar{v}]$ be a CRPQ with the



Figure 7: The graph patterns for q_1 and q_2 .

property that, for every graph database \mathcal{D} , $q[\bar{v}](\mathcal{D}) = q^{\bar{v}}(\mathcal{D})$. Such $q[\bar{v}]$ exist due to Lemma 11. This directly implies that, for every graph database \mathcal{D} , $q(\mathcal{D}) = \bigcup_{\bar{v} \in (\Sigma^{\leq k})^n} q[\bar{v}](\mathcal{D})$. Moreover, we can conclude that $\llbracket q \rrbracket = \llbracket q' \rrbracket$, where $q' \in \cup\text{-CRPQ}$ is defined by $q' = \bigvee_{\bar{v} \in (\Sigma^{\leq k})^n} q[\bar{v}]$. \square

7.6 $\llbracket \text{CXPQ}^{\leq k} \rrbracket \neq \llbracket \text{CRPQ} \rrbracket$ and $\llbracket \text{ECRPQ}^{\text{er}} \rrbracket \neq \llbracket \text{CXPQ} \rrbracket$

Interestingly, we can also show that the expressive power of CXPQ properly exceeds that of ECRPQ^{er}, and that the expressive power of $\text{CXPQ}^{\leq k}$, for every $k \geq 1$, properly exceeds that of CRPQ.

Lemma 15. *For every $k \geq 1$, there is a Boolean $q \in \text{CXPQ}^{\leq k}$ with $\llbracket q \rrbracket \notin \llbracket \text{CRPQ} \rrbracket$.*

Proof. Let the Boolean $q_1 \in \text{CXPQ}$ over $\Sigma = \{\mathbf{a}, \mathbf{b}, \mathbf{c}, \mathbf{d}\}$ be defined by the graph pattern $G_{q_1} = (V_{q_1}, E_{q_1})$ with $V_{q_1} = \{u_1, u_2, u_3, u_4\}$ and

$$E_{q_1} = \{(u_1, \alpha_1, u_2), (u_3, \alpha_2, u_2), (u_3, \alpha_3, u_4)\},$$

where $\alpha_1 = x\{\mathbf{a} \vee \mathbf{b}\}$, $\alpha_2 = \mathbf{d}$ and $\alpha_3 = x \vee \mathbf{c}$ (see Figure 7). We note that, for every $k \geq 1$ and every graph database \mathcal{D} , $q_1^k(\mathcal{D}) = q_1^1(\mathcal{D}) = q_1(\mathcal{D})$. Thus, in order to prove the statement of the lemma for every $k \geq 1$, it is sufficient to show that $\llbracket q_1 \rrbracket \notin \llbracket \text{CRPQ} \rrbracket$, where q_1 is interpreted as a CXPQ without any restrictions.

For every $\sigma_1, \sigma_2 \in \Sigma$, let $\mathcal{D}_{\sigma_1, \sigma_2} = (V_{\sigma_1, \sigma_2}, E_{\sigma_1, \sigma_2})$ be a graph-database with $V_{\sigma_1, \sigma_2} = \{v_1, v_2, v_3, v_4\}$ and

$$E_{\sigma_1, \sigma_2} = \{(v_1, \sigma_1, v_2), (v_3, \mathbf{d}, v_2), (v_3, \sigma_2, v_4)\}.$$

We note that $\mathcal{D}_{\sigma_1, \sigma_2} \models q_1$ for every $\sigma_1 \in \{\mathbf{a}, \mathbf{b}\}$ and $\sigma_2 \in \{\mathbf{a}, \mathbf{b}, \mathbf{c}\}$ with $\sigma_1 = \sigma_2$ or $\sigma_2 = \mathbf{c}$. We assume that there is a $q' \in \text{CRPQ}$ with $q' \equiv q_1$, defined by the graph pattern $G_{q'} = (V_{q'}, E_{q'})$ with $E_{q'} = \{(x_i, \beta_i, y_i) \mid i \in [m]\}$.

If, for some $i \in [m]$, there is no $w_i \in \mathcal{L}(\beta_i)$ with $|w_i|_{\mathbf{a}} = 0$, then $\mathcal{D}_{\mathbf{b}, \mathbf{b}} \not\models q'$, which, since $\mathcal{D}_{\mathbf{b}, \mathbf{b}} \models q_1$, is a contradiction. Therefore, we can assume that, for every $i \in [m]$, there is some $w_i \in \mathcal{L}(\beta_i)$ with $|w_i|_{\mathbf{a}} = 0$ (note that $w_i = \varepsilon$ is possible). Next, we consider a matching morphism h for q' and $\mathcal{D}_{\mathbf{a}, \mathbf{a}}$, which, since $\mathcal{D}_{\mathbf{a}, \mathbf{a}} \models q_1$, must exist. Let $A \subseteq [m]$ be exactly the set of $i \in [m]$ with $h(x_i) = v_1$ and $h(y_i) = v_2$. If $A = \emptyset$, then we can conclude that the edge (v_1, \mathbf{a}, v_2) is not part of any of the paths between some $h(x_i)$ and $h(y_i)$ that are necessary for h being a matching morphism (this is due to the fact that in $\mathcal{D}_{\mathbf{a}, \mathbf{a}}$ there are no paths of length strictly greater than 2). Consequently, we can remove (v_1, \mathbf{a}, v_2) from $\mathcal{D}_{\mathbf{a}, \mathbf{a}}$ in order to obtain a graph database \mathcal{D}' , such that h would still be a matching morphism for q' and \mathcal{D}' . This, however, is a contradiction, since $\mathcal{D}' \not\models q_1$. Thus, $A \neq \emptyset$. Now let $\mathcal{D}'_{\mathbf{a}, \mathbf{a}}$ be the graph database obtained from $\mathcal{D}_{\mathbf{a}, \mathbf{a}}$, by deleting the edge (v_1, \mathbf{a}, v_2) and, for every $i \in A$, adding a

path labelled with w_i from v_1 to v_2 . We note that h is still a matching morphism for q' and $\mathcal{D}'_{\mathbf{a},\mathbf{a}}$, since, for every $i \in A$, there is a path from v_1 to v_2 labelled with $w_i \in \mathcal{L}(\beta_i)$, and, furthermore, as already observed above, the deleted edge (v_1, \mathbf{a}, v_2) was exclusively covered by edges (x_i, β_i, y_i) with $i \in A$. However, it can be easily seen that $\mathcal{D}'_{\mathbf{a},\mathbf{a}} \not\models q_1$, which is a contradiction. \square

Lemma 16. *There is a Boolean $q \in \text{CXPQ}$ such that $\llbracket q \rrbracket \notin \llbracket \text{ECRPQ}^{\text{er}} \rrbracket$.*

Proof. Let q_2 be defined by a graph pattern with just a single edge (u_1, β, u_2) with $\beta = \#y\{x\{\mathbf{a}^+\mathbf{b}\}x^*\}cy\#$ (see Figure 7). We note that $\mathcal{D} \models q_2$ if and only if \mathcal{D} contains a path labelled with $\#(\mathbf{a}^{n_1}\mathbf{b})^{n_2}\mathbf{c}(\mathbf{a}^{n_1}\mathbf{b})^{n_2}\#$ for some $n_1, n_2 \geq 1$. Let us assume that there is some $q' \in \text{ECRPQ}^{\text{er}}$ with $q_2 \equiv q'$ and q' is defined by a graph pattern $G_{q'} = (V_{q'}, E_{q'})$ with $E_{q'} = \{(x_i, \alpha_i, y_i) \mid i \in [m]\}$ and some equality relations. Moreover, for every $i \in [m]$, let p_i be the pumping lemma constant of $\mathcal{L}(\alpha_i)$ and let $p = \max\{p_i \mid i \in [m]\}$. We consider the graph database $\mathcal{D} = (V_{\mathcal{D}}, E_{\mathcal{D}})$ with $V_{\mathcal{D}} = \{v_0, v_1, \dots, v_t\}$, where $t = 2(p^2m + pm) + 3$ and (v_0, v_1, \dots, v_t) is a path labelled with $\#(\mathbf{a}^p\mathbf{b})^{pm}\mathbf{c}(\mathbf{a}^p\mathbf{b})^{pm}\#$.

Since $\mathcal{D} \models q'$, there is a matching morphism h for q' and \mathcal{D} with some matching words (w_1, w_2, \dots, w_m) , such that, for every $i \in [m]$, $h(x_i) = v_{j_i}$ and $h(y_i) = v_{j'_i}$ with $0 \leq j_i \leq j'_i \leq t$. We now partition $[m]$ into $S = \{i \mid j'_i - j_i < 2p + 1\}$ and $L = [m] \setminus S$, i. e., (x_i, α_i, y_i) is matched to a *long* sub-path of (v_0, v_1, \dots, v_t) of length at least $2p + 1$ if $i \in L$ and to a *short* sub-path of (v_0, v_1, \dots, v_t) of length strictly less than $2p + 1$ otherwise.

If there is an arc (v_r, σ, v_{r+1}) of the path (v_0, v_1, \dots, v_t) that is not covered by some (x_i, α_i, y_i) (i. e., for every $i \in [m]$, it is not the case that $j_i \leq r < j'_i$), then we can contract nodes v_r and v_{r+1} and h is still a matching morphism for q' and the thus modified graph database \mathcal{D}' , which is not in $\llbracket q_2 \rrbracket$ anymore. This can be seen by observing that removing a single symbol from a word $w \in \mathcal{L}(\beta)$ yields a word that is not in $\mathcal{L}(\beta)$ anymore. Therefore, we can assume that every arc (v_r, σ, v_{r+1}) of \mathcal{D} is covered by some (x_i, α_i, y_i) , i. e., $j_i \leq r < j'_i$. Since, for every $i \in S$, at most $2p$ edges can be covered by (x_i, α_i, y_i) , we also know that $L \neq \emptyset$ (since otherwise not all arcs are covered).

We now modify \mathcal{D} as follows. For every $i \in [m]$, we add a *shortcut* from v_{j_i} to $v_{j'_i}$, which is a new path of the same length and with the same label as the path $(v_{j_i}, v_{j_i+1}, \dots, v_{j'_i})$. In particular, we note that for every $i \in [m]$, the labels of the shortcuts are identical to the matching words of h .

We now pump some of the shortcuts depending on the equality relation of q' as follows. Assume that $A \subseteq [m]$ represents an equality relation of q' , i. e., exactly the edges $\{(x_i, \alpha_i, y_i) \mid i \in A\}$ are subject to the equality relation. This also means that either all (x_i, α_i, y_i) with $i \in A$ cover a short path, i. e., $A \subseteq S$, or all (x_i, α_i, y_i) with $i \in A$ cover a long path, i. e., $A \subseteq L$. Moreover, as observed above, $L \neq \emptyset$, so there is at least one such equality relation A (note that we assume that the equality relations are represented by a partition of the edge-set, i. e., every edge is subject to exactly one equality relation, possibly a unary one).

Recall that (w_1, w_2, \dots, w_m) are the matching words for h . Thus, for every $i \in A$, w_i is the label of the sub-path $(v_{j_i}, v_{j_i+1}, \dots, v_{j'_i})$. Since h is a matching morphism and since we have the equality relation represented by A , we know that, for some u , we have $u = w_i$ for every $i \in A$; moreover, the corresponding shortcuts for edges (x_i, α_i, y_i) with $i \in A$ are also all labelled with u . Since

$A \subseteq L$, we have $|u| \geq 2p + 1$, which means that $u = u' \mathbf{a}^p u''$. Hence, for every $i \in A$, there is a d_i such that $u' \mathbf{a}^{p+\delta d_i} u'' \in \mathcal{L}(\alpha_i)$ for every $\delta \geq 0$. This means that, for every $i \in A$, $w' = u' \mathbf{a}^{p+d} u'' \in \mathcal{L}(\alpha_i)$, where $d = \prod_{i \in A} d_i$. Consequently, we can pump all shortcuts for the edges (x_i, α, y_i) with $i \in A$ by the factor \mathbf{a}^d , i. e., we replace them by paths of length $|u| + d$ labelled by w' . We repeat this pumping-step with respect to all equality relations that refer to edges that cover long paths. After this modification, we have the property that in the tuple (w_1, w_2, \dots, w_m) of matching words for h , we can arbitrarily replace some w_i by the label of the corresponding shortcut for (x_i, α_i, y_i) (regardless of whether it has been pumped or not) and still h is a matching morphism with respect to this modified tuple of matching words.

We now choose an arbitrary edge $(v_\ell, \sigma, v_{\ell+1})$ of the original path (v_0, \dots, v_t) , which is only covered by edges (x_i, α_i, y_i) with $i \in L$, which means that their shortcuts have been pumped. Such an edge must exist, since otherwise all edges are covered by edges from S , which is not possible. Then, we delete this edge and we denote the obtained graph database by \mathcal{D}' . After this modification, due to the shortcuts, the matching morphism h must still be a valid matching morphism for q' and \mathcal{D} , i. e., $\mathcal{D} \models q'$. We now conclude the proof by showing that $\mathcal{D}' \not\models q_2$, which clearly is a contradiction.

For $\mathcal{D}' \models q_2$, there must be a path in \mathcal{D}' that is labelled by a word of the form $\#(\mathbf{a}^{n_1} \mathbf{b})^{n_2} \mathbf{c}(\mathbf{a}^{n_1} \mathbf{b})^{n_2} \#$ for some $n_1, n_2 \geq 1$. We note that this is only possible for paths from v_0 to v_t . Now let us consider an arbitrary path from v_0 and v_t in \mathcal{D}' . Since we deleted the edge $(v_\ell, \sigma, v_{\ell+1})$, this path must use at least one shortcut that corresponds to an edge (x_i, α_i, y_i) with $j_i \leq \ell < j'_i$. However, by our choice of $(v_\ell, \sigma, v_{\ell+1})$, all such shortcuts have been pumped, which means that the path is labelled with a word \hat{w} that can be obtained from $\#(\mathbf{a}^p \mathbf{b})^{pm} \mathbf{c}(\mathbf{a}^p \mathbf{b})^{pm} \#$ by pumping some unary factors over \mathbf{a} . Furthermore, it is not possible that all maximal unary factors over \mathbf{a} , i. e., factors of the form $\# \mathbf{a}^p \mathbf{b}$, $\mathbf{b} \mathbf{a}^p \mathbf{b}$ or $\mathbf{c} \mathbf{a}^p \mathbf{b}$, have been pumped, since the considered path can take at most m shortcuts. \square

8 Conclusions and Open Problems

The fact that evaluation for CXRPQ is at least PSpace-hard even in data-complexity is reason enough to look at fragments of CXRPQ instead. Nevertheless, an upper bound for evaluating unrestricted CXRPQs would be interesting from a theoretical point of view; similarly, a lower bound for CXRPQ^{vsf}-evaluation in data-complexity would be interesting. Several of our results implicitly pose conciseness questions. Each CXRPQ ^{$\leq k$} can be represented as the union of $O(|\Sigma| + 1)^{nk}$ many CRPQs, and each CXRPQ^{vsf} can be represented as the union of exponentially many ECRPQ^{ers} of exponential size. Are these exponential blow-ups necessary? Theorem 6 gives some evidence that for CXRPQ ^{$\leq k$} this is the case.

All our algorithms for BOOL-EVAL can be extended to the problem CHECK. With respect to also extracting paths instead of only nodes from the graph-database, we can use the general techniques of [8] to some extent. More precisely, for a $q \in \text{CXRPQ}^{\text{vsf}}$ (or $q \in \text{CXRPQ}^{\leq k}$), a graph database \mathcal{D} and a tuple $\bar{t} \in (V_{\mathcal{D}})^\ell$, our evaluation algorithms can be adapted to produce an automaton that represents all tuples of paths corresponding to the matching morphisms of q and \mathcal{D} , but they are rather large. A thorough analysis of CXRPQ-fragments with

respect to how they can be used as queries that also extract path-variables is left for future work.

Interestingly enough, restrictions that lower the complexity for matching xregex to strings do not seem to help at all if we use them for querying graphs, and vice versa. In the string case, the NP-complete matching problem for xregex trivially becomes polynomial-time solvable, if the number of variables is bounded by a constant (see [40]), while this does not help for graphs (see Theorem 1). Variable-star freeness (Section 5) or bounding the image size (Section 6) has a positive impact with respect to graphs. However, the matching problem for xregex remains NP-hard, even if we require variable-star freeness *and* that variables can only range over words of length at most 1 (see [22, Theorem 3]).

8.1 Acknowledgments

The author thanks Nicole Schweikardt for helpful discussions. The author is supported by DFG-grant SCHM 3485/1-1.

References

- [1] Serge Abiteboul, Peter Buneman, and Dan Suciu. *Data on the Web: From Relations to Semistructured Data and XML*. Morgan Kaufmann, 1999.
- [2] Serge Abiteboul, Dallan Quass, Jason McHugh, Jennifer Widom, and Janet L. Wiener. The lorel query language for semistructured data. *Int. J. on Digital Libraries*, 1(1):68–88, 1997.
- [3] Alfred V. Aho. Algorithms for finding patterns in strings. In *Handbook of Theoretical Computer Science, Volume A: Algorithms and Complexity (A)*, pages 255–300. 1990.
- [4] Renzo Angles, Marcelo Arenas, Pablo Barceló, Aidan Hogan, Juan L. Reutter, and Domagoj Vrgoc. Foundations of modern query languages for graph databases. *ACM Comput. Surv.*, 50(5):68:1–68:40, 2017.
- [5] Renzo Angles, Juan L. Reutter, and Hannes Voigt. Graph query languages. In *Encyclopedia of Big Data Technologies*. 2019.
- [6] Pablo Barceló. Querying graph databases. In *Proceedings of the 32nd ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems, PODS 2013, New York, NY, USA - June 22 - 27, 2013*, pages 175–188, 2013.
- [7] Pablo Barceló, Diego Figueira, and Miguel Romero. Boundedness of conjunctive regular path queries. In *46th International Colloquium on Automata, Languages, and Programming, ICALP 2019, July 9-12, 2019, Patras, Greece*, pages 104:1–104:15, 2019.
- [8] Pablo Barceló, Leonid Libkin, Anthony Widjaja Lin, and Peter T. Wood. Expressive languages for path queries over graph-structured data. *ACM Transactions on Database Systems (TODS)*, 37(4):31:1–31:46, 2012.

- [9] Pablo Barceló, Leonid Libkin, and Juan L. Reutter. Querying regular graph patterns. *Journal of the ACM*, 61(1):8:1–8:54, 2014.
- [10] Pablo Barceló, Miguel Romero, and Moshe Y. Vardi. Semantic acyclicity on graph databases. *SIAM J. Comput.*, 45(4):1339–1376, 2016.
- [11] Angela Bonifati, Wim Martens, and Thomas Timm. An analytical study of large SPARQL query logs. *PVLDB*, 11(2):149–161, 2017.
- [12] Angela Bonifati, Wim Martens, and Thomas Timm. Navigating the maze of wikidata query logs. In *The World Wide Web Conference, WWW 2019, San Francisco, CA, USA, May 13-17, 2019*, pages 127–138, 2019.
- [13] Diego Calvanese, Giuseppe De Giacomo, Maurizio Lenzerini, and Moshe Y. Vardi. Containment of conjunctive regular path queries with inverse. In *KR 2000, Principles of Knowledge Representation and Reasoning Proceedings of the Seventh International Conference, Breckenridge, Colorado, USA, April 11-15, 2000.*, pages 176–185, 2000.
- [14] Diego Calvanese, Giuseppe De Giacomo, Maurizio Lenzerini, and Moshe Y. Vardi. Reasoning on regular path queries. *SIGMOD Record*, 32(4):83–92, 2003.
- [15] Cezar Câmpeanu, Kai Salomaa, and Sheng Yu. A formal study of practical regular expressions. *Int. J. Found. Comput. Sci.*, 14(6).
- [16] Mariano P. Consens and Alberto O. Mendelzon. Graphlog: a visual formalism for real life recursion. In *Proceedings of the Ninth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, April 2-4, 1990, Nashville, Tennessee, USA*, pages 404–416, 1990.
- [17] Isabel F. Cruz, Alberto O. Mendelzon, and Peter T. Wood. A graphical query language supporting recursion. In *Proceedings of the Association for Computing Machinery Special Interest Group on Management of Data 1987 Annual Conference, San Francisco, California, May 27-29, 1987*, pages 323–330, 1987.
- [18] Wojciech Czerwinski, Wim Martens, Matthias Niewerth, and Pawel Parys. Minimization of tree patterns. *J. ACM*, 65(4):26:1–26:46, 2018.
- [19] Alin Deutsch and Val Tannen. Optimization properties for classes of conjunctive regular path queries. In *Database Programming Languages, 8th International Workshop, DBPL 2001, Frascati, Italy, September 8-10, 2001, Revised Papers*, pages 21–39, 2001.
- [20] Johannes Doleschal, Benny Kimelfeld, Wim Martens, Yoav Nahshon, and Frank Neven. Split-correctness in information extraction. In *Proceedings of the 38th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, PODS 2019, Amsterdam, The Netherlands, June 30 - July 5, 2019*, pages 149–163, 2019.
- [21] Henning Fernau, Florin Manea, Robert Mercas, and Markus L. Schmid. Pattern matching with variables: Fast algorithms and new hardness results.

- In *32nd International Symposium on Theoretical Aspects of Computer Science, STACS 2015, March 4-7, 2015, Garching, Germany*, pages 302–315, 2015.
- [22] Henning Fernau and Markus L. Schmid. Pattern matching with variables: A multivariate complexity analysis. *Information and Computation (I&C)*, 242:287–305, 2015.
- [23] Henning Fernau, Markus L. Schmid, and Yngve Villanger. On the parameterised complexity of string morphism problems. *Theory of Computing Systems (ToCS)*, 59(1):24–51, 2016.
- [24] Daniela Florescu, Alon Y. Levy, and Dan Suciu. Query containment for conjunctive queries with regular expressions. In *Proceedings of the Seventeenth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, June 1-3, 1998, Seattle, Washington, USA*, pages 139–148, 1998.
- [25] Dominik D. Freydenberger. Extended regular expressions: Succinctness and decidability. *Theory of Computing Systems (ToCS)*, 53(2):159–193, 2013.
- [26] Dominik D. Freydenberger. A logic for document spanners. *Theory Comput. Syst.*, 63(7):1679–1754, 2019.
- [27] Dominik D. Freydenberger, Benny Kimelfeld, and Liat Peterfreund. Joining extractions of regular expressions. In *Proceedings of the 37th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems, Houston, TX, USA, June 10-15, 2018*, pages 137–149, 2018.
- [28] Dominik D. Freydenberger and Markus L. Schmid. Deterministic regular expressions with back-references. *J. Comput. Syst. Sci.*, 105:1–39, 2019.
- [29] Jeffrey E. F. Friedl. *Mastering regular expressions - understand your data and be more productive: for Perl, PHP, Java, .NET, Ruby, and more (3. ed.)*. O’Reilly, 2006.
- [30] Steve Harris and Andy Seaborne. Sparql 1.1 query language. W3C recommendation. Technical Report <https://www.w3.org/TR/sparql11-query/>, W3C, March 2013.
- [31] The IEEE and The Open Group. IEEE std 1003.1-2008, 2016 edition, chapter 9. <http://pubs.opengroup.org/onlinepubs/9699919799/>, 2016.
- [32] Egor V. Kostylev, Juan L. Reutter, and Domagoj Vrgoc. Containment of queries for graphs with data. *J. Comput. Syst. Sci.*, 92:65–91, 2018.
- [33] Leonid Libkin, Wim Martens, and Domagoj Vrgoc. Querying graphs with data. *Journal of the ACM*, 63(2):14:1–14:53, 2016.
- [34] Katja Losemann and Wim Martens. The complexity of regular expressions and property paths in SPARQL. *ACM Transactions on Database Systems (TODS)*, 38(4):24:1–24:39, 2013.

- [35] Wim Martens, Matthias Niewerth, and Tina Trautner. A trichotomy for regular trail queries. In *37th Symposium on Theoretical Aspects of Computer Science, STACS 2020, March 10-13, 2020, Montpellier, Germany, 2020*. Accepted for publication.
- [36] Wim Martens and Tina Trautner. Evaluation and enumeration problems for regular path queries. In *21st International Conference on Database Theory, ICDT 2018, March 26-29, 2018, Vienna, Austria*, pages 19:1–19:21, 2018.
- [37] Wim Martens and Tina Trautner. Bridging theory and practice with query log analysis. *SIGMOD Record*, 48(1):6–13, 2019.
- [38] Alberto O. Mendelzon and Peter T. Wood. Finding regular simple paths in graph databases. *SIAM Journal on Computing (SICOMP)*, 24(6):1235–1258, 1995.
- [39] Markus L. Schmid. Inside the class of REGEX languages. *International Journal of Foundations of Computer Science (IJFCS)*, 24(7):1117–1134, 2013.
- [40] Markus L. Schmid. Characterising REGEX languages by regular languages equipped with factor-referencing. *Information and Computation (I&C)*, 249:1–17, 2016.
- [41] The Neo4j Team. The neo4j cypher manual v3.5. 2016.
- [42] Apache TinkerPopTM. Gremlin language. 2013.
- [43] Peter T. Wood. Query languages for graph databases. *SIGMOD Record*, 41(1):50–60, 2012.