

Unique Characterisability and Learnability of Temporal Instance Queries

M. Fortin¹, B. Konev¹, V. Ryzhikov², Y. Savateev^{2,3}, F. Wolter¹, M. Zakharyashev²

¹Department of Computer Science, University of Liverpool, UK

²Department of Computer Science, Birkbeck, University of London, UK

³HSE University, Moscow, Russia

{marie.fortin,boris.konev,wolter@liverpool.ac.uk}, {vlad,yury,michael}@dcs.bbk.ac.uk

Abstract

We aim to determine which temporal instance queries can be uniquely characterised by a (polynomial-size) set of positive and negative temporal data examples. We start by considering queries formulated in fragments of propositional linear temporal logic *LTL* that correspond to conjunctive queries (CQs) or extensions thereof induced by the until operator. Not all of these queries admit polynomial characterisations, but by imposing a further restriction to path-shaped queries we identify natural classes that do. We then investigate how far the obtained characterisations can be lifted to temporal knowledge graphs queried by 2D languages combining *LTL* with concepts in description logics \mathcal{EL} or \mathcal{ELI} (i.e., tree-shaped CQs). While temporal operators in the scope of description logic constructors can destroy polynomial characterisability, we obtain general transfer results for the case when description logic constructors are within the scope of temporal operators. We finally apply our characterisations to establish (polynomial) learnability of temporal instance queries using membership queries in the active learning framework.

1 Introduction

Constructing queries or, more generally, logical concepts describing individuals of interest, can be difficult. Supporting a user addressing this problem has been a major research topic in databases, logic, and knowledge representation. For instance, in reverse engineering of database queries and concept descriptions (Martins 2019; Lehmann and Hitzler 2010; Jung et al. 2020), one aims to identify a query using a set of positively and negatively labelled examples of answers and non-answers, respectively; and in active learning approaches, one aims to identify a query by asking an oracle (e.g., domain specialist) whether an example is an answer or a non-answer to the query (Angluin, Frazier, and Pitt 1992; Funk, Jung, and Lutz 2021; ten Cate and Dalmau 2021).

Recently, the *unique characterisation* of a query by a finite (ideally, polynomial-size) set of positive and negative example answers has been identified as a fundamental link between queries and data (ten Cate and Dalmau 2021). Namely, we say that a query q fits a pair $E = (E^+, E^-)$ of sets E^+ and E^- of pointed databases (\mathcal{D}, a) if $\mathcal{D} \models q(a)$ for $(\mathcal{D}, a) \in E^+$, and $\mathcal{D} \not\models q(a)$ for $(\mathcal{D}, a) \in E^-$. Then E uniquely characterises q within a class \mathcal{Q} of queries if q is the only (up to equivalence) query in \mathcal{Q} that fits E .

Unique (polynomial) characterisations can be used to illustrate, explain, and construct queries. They are also a ‘non-procedural’ necessary condition for (polynomial) learnability using membership queries in Angluin’s (1987b) framework of active learning, where membership queries to the oracle take the form ‘does $\mathcal{D} \models q(a)$ hold?’. It is shown by ten Cate and Dalmau (2021) that, for classes of conjunctive queries (CQs), it is often a sufficient condition as well.

In many applications, queries are required to capture the temporal evolution of individuals, making their formulation even harder. The aim of this paper is to start an investigation of the (polynomial) characterisability of temporal instance queries. We first consider one-dimensional data instances of the form $(\delta_0, \dots, \delta_n)$, where δ_i is the set of atomic propositions that are true at timestamp i , describing the temporal behaviour of a single individual, and queries formulated in fragments of propositional linear temporal logic *LTL*. Although rather basic as a temporal data model, this restriction allows us to focus on the purely temporal aspect of unique characterisability. We then generalise our results, where possible, to standard two-dimensional temporal data instances, in which the δ_i are replaced by non-temporal data instances and queries are obtained by combining fragments of *LTL* with \mathcal{ELI} -concept queries (or tree-shaped CQs), thereby combining a well established formalism for accessing temporal data (Chomicki and Toman 2018) with the basic concept descriptions for tractable data access from description logic (Baader et al. 2017).

Our initial observation is that already very primitive temporal queries are not uniquely characterisable. For example, consider the query $q = \diamond_r(A \wedge B)$ with the operator \diamond_r ‘now or later’ (interpreted by \leq over linearly ordered timestamps). By the pigeonhole principle, no finite example set E can distinguish q from a query $q' = \diamond_r(A \wedge (\diamond_r B \wedge \diamond_r(A \wedge \dots)))$ with sufficiently many alternating A and B . Similarly, the query $q = \circ A$ with the ‘next time’ operator \circ is not distinguishable by a finite example set from $q' = (\circ \dots \circ A) \cup A$ with the strict ‘until’ operator \cup and sufficiently many \circ on its left-hand side.

Aiming to identify natural and useful classes of temporal queries enjoying (polynomial) characterisability, in this paper we consider the conjunctive fragment of *LTL*. To begin with, we focus on two classes of path CQs: the class

$\mathcal{Q}_p[\circ, \diamond_r]$ of queries of the form

$$q = \rho_0 \wedge \mathbf{o}_1(\rho_1 \wedge \mathbf{o}_2(\rho_2 \wedge \dots \wedge \mathbf{o}_n \rho_n)), \quad (1)$$

where $\mathbf{o}_i \in \{\circ, \diamond_r\}$ and ρ_i is a conjunction of atomic propositions, and also the class $\mathcal{Q}_p^\sigma[\mathbf{U}]$ of \mathbf{U} -queries of the form

$$q = \rho_0 \wedge (\lambda_1 \mathbf{U} (\rho_1 \wedge (\lambda_2 \mathbf{U} (\dots (\lambda_n \mathbf{U} \rho_n) \dots))), \quad (2)$$

where λ_i is a conjunction of atoms or \perp . The superscript σ in $\mathcal{Q}_p^\sigma[\mathbf{U}]$ indicates that queries are formulated in a finite signature σ , a condition required because of the universal quantification implicit in \mathbf{U} . Our first main result is a syntactic criterion of (polynomial) characterisability of $\mathcal{Q}_p[\circ, \diamond_r]$ -queries. In fact, it turns out that the query $\diamond_r(A \wedge B)$ mentioned above epitomises the cause of non-characterisability in $\mathcal{Q}_p[\circ, \diamond_r]$. It follows, in particular, that the restriction $\mathcal{Q}_p[\circ, \diamond]$ of $\mathcal{Q}_p[\circ, \diamond_r]$ to queries with \circ and strict eventuality $\diamond = \circ \diamond_r$ is polynomially characterisable. Our second main result is that all $\mathcal{Q}_p^\sigma[\mathbf{U}]$ -queries with \subseteq -incomparable λ_i and ρ_i , for each i , are polynomially characterisable within $\mathcal{Q}_p^\sigma[\mathbf{U}]$. Although we show that all $\mathcal{Q}_p^\sigma[\mathbf{U}]$ -queries are exponentially characterisable, it remains open whether they are polynomially characterisable in $\mathcal{Q}_p^\sigma[\mathbf{U}]$.

The essential property that distinguishes $\mathcal{Q}_p[\circ, \diamond_r]$ and $\mathcal{Q}_p^\sigma[\mathbf{U}]$ from other queries is that they do not admit temporal branching as, for instance, in $\diamond A \wedge \diamond B$. In fact, we show that even within the class of queries using only \wedge and \diamond and with a bound on the number of branches, not all queries are polynomially characterisable. A first step towards positive results covering non-path queries is made for the case of queries in which all branches are of equal length.

Our next aim is to generalise the obtained results to 2D temporal queries combining *LTL* and the descriptions logic constructor $\exists P$ of \mathcal{ELI} . Our first main result is negative: even queries of the form $\exists P.q_1 \wedge \dots \wedge \exists P.q_n$ with $q_i \in \mathcal{Q}_p[\circ, \diamond]$ are not polynomially characterisable. The situation changes drastically, however, if we consider queries of the form (1) or (2), in which ρ_i and λ_i are \mathcal{ELI} -queries. Indeed, we generalise our polynomial characterisability results for $\mathcal{Q}_p[\circ, \diamond_r]$ and $\mathcal{Q}_p^\sigma[\mathbf{U}]$ to such queries using recent results on the computation of frontiers in the lattice of \mathcal{ELI} -queries (ten Cate and Dalmau 2021) and proving a new result on split partners in the lattice of \mathcal{EL} -queries (where \mathcal{EL} denotes \mathcal{ELI} without inverse roles).

Finally, we discuss applications of our results to learning temporal instance queries using membership queries of the form ‘does $\mathcal{D} \models q$ hold?’. As we always construct example sets effectively, our unique (exponential) characterisability results imply (exponential-time) learnability with membership queries. Obtaining polynomial-time learnability from polynomial characterisations is more challenging and, in fact, not always possible. A main result here is that $\mathcal{Q}_p[\circ, \diamond_r]$ with \mathcal{ELI} -queries is polynomial-time learnable with membership queries.

Omitted proofs can be found in the appendix to the full version of the paper.

2 Related Work

Our contribution is closely related to work on active learnability of formal languages and on learning temporal logic

formulas interpreted over finite and infinite traces. It is also related to learning database queries and other formal expressions for accessing data. In the former area, the seminal paper by Angluin (1987a) has given rise to a large body of work on active learning of regular languages or variations, for example, (Shahbaz and Groz 2009; Aarts and Vaandrager 2010; Cassel et al. 2016; Howar and Steffen 2018). This work has mainly focused on learning various types of finite state machines or automata using a combination of membership queries with other powerful types of queries such as equivalence queries. The use of two or more types of queries is motivated by the fact that otherwise one cannot efficiently learn a wide variety of important languages, including regular languages. In fact, the main difference between this work and our contribution is that we focus on queries for which the corresponding formal languages form only a small subset of the regular languages and it is this restriction that enables us to focus on characterisability and learnability with membership queries.

Rather surprisingly, there has hardly been any work on active learning of temporal formulas over finite or infinite traces; we refer the reader to (Camacho and McIlraith 2019), also for a discussion of the relationship between learning automata and *LTL*-formulas. In contrast, passive learning of *LTL*-formulas has recently received significant attention; see (Lemieux, Park, and Beschastnikh 2015; Neider and Gavran 2018; Fijalkow and Lagarde 2021) and, in the context of explainable AI, also (Camacho and McIlraith 2019) for an overview.

In the database and KR communities, there has been extensive work on identifying queries and concept descriptions from data examples. For instance, in reverse engineering of queries, the goal is typically to decide whether there is a query that fits (or separates) a set of positive and negative examples. Relevant work under the closed world assumption include (Arenas and Diaz 2016; Barceló and Romero 2017) and under the open world assumption (Gutiérrez-Basulto, Jung, and Sabellek 2018; Funk et al. 2019). Related work on active learning not yet discussed include the identification of \mathcal{EL} -queries (Funk, Jung, and Lutz 2021) and ontologies (Konev, Ozaki, and Wolter 2016; Konev et al. 2017), and of schema-mappings (ten Cate, Dalmau, and Kolaitis 2013; ten Cate et al. 2018). Again this work differs from our contribution as it focusses on learning using membership and equivalence queries rather than only the former. The use of unique characterisations to explain and construct schema mappings has been promoted and investigated in (Kolaitis 2011; Alexe et al. 2011).

Combining *LTL* and description logics for temporal conceptual modeling and data access has a long tradition (Lutz, Wolter, and Zakharyashev 2008; Artale et al. 2017). For querying purposes sometimes description logic concepts have been replaced by general CQs. Our restriction to \mathcal{ELI} -concepts instead of general CQs is motivated by results of (ten Cate and Dalmau 2021) showing that only CQs that are acyclic modulo cycles through the answer variables are polynomially characterisable within the class of CQs. Hence very strong acyclicity conditions are needed to ensure polynomial characterisability. We conjecture that our results can

be extended to this class.

The class of queries in which no $\exists P$ is within the scope of temporal operators was first introduced by (Baader, Borgwardt, and Lippmann 2015; Borgwardt and Thost 2015) in the context of monitoring applications.

3 Preliminaries

By a *signature* we mean any finite set $\sigma \neq \emptyset$ of *atomic concepts* A, B, C, \dots representing observations, measurements, events, etc. A σ -*data instance* is any finite sequence $\mathcal{D} = (\delta_0, \dots, \delta_n)$ with $\delta_i \subseteq \sigma$, saying that $A \in \delta_i$ happened at moment i . The *length* of \mathcal{D} is $\max(\mathcal{D}) = n$ and the *size* of \mathcal{D} is $|\mathcal{D}| = \sum_{i \leq n} |\delta_i|$. We do not distinguish between \mathcal{D} and its *variants* of the form $(\delta_0, \dots, \delta_n, \emptyset, \dots, \emptyset)$.

We access data by means of *queries*, q , constructed from atoms, \perp and \top using \wedge and the temporal operators $\circ, \diamond, \diamond_r$ and \cup . The set of atomic concepts occurring in q is denoted by $\text{sig}(q)$. The set of queries that only use the operators from $\Phi \subseteq \{\circ, \diamond, \diamond_r, \cup\}$ is denoted by $\mathcal{Q}[\Phi]$; $\mathcal{Q}^\sigma[\Phi]$ is the restriction of $\mathcal{Q}[\Phi]$ to a signature σ . The *size* $|q|$ of q is the number of symbols in q , and the *temporal depth* $\text{tdp}(q)$ of q is the maximum number of nested temporal operators in q .

$\mathcal{Q}[\circ, \diamond, \diamond_r]$ -queries can be equivalently defined as *tree-shaped conjunctive queries* (CQs) with the binary predicates *suc*, $<$, \leq over \mathbb{N} , and atomic concepts as unary predicates. Each such CQ is a set $Q(t_0)$ of assertions of the form $A(t)$, $\text{suc}(t, t')$, $t < t'$, and $t \leq t'$, with a distinguished variable t_0 , such that, for every variable t in $Q(t_0)$, there exists exactly one path from t_0 to t along the binary predicates *suc*, $<$, \leq .

The set of $\mathcal{Q}[\circ, \diamond, \diamond_r]$ -queries with *path-shaped CQ* counterparts is denoted by $\mathcal{Q}_p[\circ, \diamond, \diamond_r]$. Such queries q take the form (1), where $\mathbf{o}_i \in \{\circ, \diamond, \diamond_r\}$ and ρ_i is a conjunction of atoms (the empty conjunction is \top). Similarly, $\mathcal{Q}_p[\cup]$ -queries take the form (2).

Given a data instance $\mathcal{D} = (\delta_0, \dots, \delta_n)$, the *truth-relation* $\mathcal{D}, \ell \models q$, for $\ell < \omega$, is defined as follows:

$$\begin{aligned} \mathcal{D}, \ell \models A & \text{ iff } A \in \delta_\ell, & \mathcal{D}, \ell \models \top, & \mathcal{D}, \ell \not\models \perp, \\ \mathcal{D}, \ell \models \mathbf{q}_1 \wedge \mathbf{q}_2 & \text{ iff } \mathcal{D}, \ell \models \mathbf{q}_1 \text{ and } \mathcal{D}, \ell \models \mathbf{q}_2, \\ \mathcal{D}, \ell \models \circ \mathbf{q} & \text{ iff } \mathcal{D}, \ell + 1 \models \mathbf{q}, \\ \mathcal{D}, \ell \models \diamond \mathbf{q} & \text{ iff } \mathcal{D}, m \models \mathbf{q}, \text{ for some } m > \ell, \\ \mathcal{D}, \ell \models \diamond_r \mathbf{q} & \text{ iff } \mathcal{D}, m \models \mathbf{q}, \text{ for some } m \geq \ell, \\ \mathcal{D}, \ell \models \mathbf{q}_1 \cup \mathbf{q}_2 & \text{ iff there is } m > \ell \text{ such that } \mathcal{D}, m \models \mathbf{q}_2 \\ & \text{ and } \mathcal{D}, k \models \mathbf{q}_1, \text{ for all } k \text{ with } \ell < k < m. \end{aligned}$$

Note that $\mathcal{D}, n \models \diamond \top \wedge \circ \top \wedge (\mathbf{q} \cup \top)$ as $(\delta_0, \dots, \delta_n, \emptyset)$ is a variant of \mathcal{D} . We write $q \models q'$ if $\mathcal{D}, \ell \models q$ implies $\mathcal{D}, \ell \models q'$, for any \mathcal{D} and ℓ . If $q \models q'$ and $q' \models q$, we call q and q' *equivalent* and write $q \equiv q'$. Since $\circ \mathbf{q} \equiv \perp \cup \mathbf{q}$, $\diamond \mathbf{q} \equiv \top \cup \mathbf{q}$ and $\diamond \mathbf{q} \equiv \circ \diamond_r \mathbf{q}$, one can assume that $\mathcal{Q}[\circ, \diamond] \subseteq \mathcal{Q}[\cup]$, $\mathcal{Q}[\diamond] \subseteq \mathcal{Q}[\circ, \diamond_r]$ and $\mathcal{Q}[\circ, \diamond_r] = \mathcal{Q}[\circ, \diamond_r, \diamond]$.

4 Unique Characterisability

An *example set* is a pair $E = (E^+, E^-)$ with finite sets E^+ and E^- of data instances that are called *positive* and *negative examples*, respectively. A query q *fits* E if $\mathcal{D}^+, 0 \models q$ and $\mathcal{D}^-, 0 \not\models q$, for all $\mathcal{D}^+ \in E^+$ and $\mathcal{D}^- \in E^-$. E

uniquely characterises $q \in \mathcal{Q}$ within a class \mathcal{Q} of queries if q fits E and $q \equiv q'$ for any $q' \in \mathcal{Q}$ that fits E . If all $q \in \mathcal{Q}$ are characterised by some E within $\mathcal{Q}' \supseteq \mathcal{Q}$, we call \mathcal{Q} *uniquely characterisable* within \mathcal{Q}' . Further, \mathcal{Q} is *polynomially characterisable* within $\mathcal{Q}' \supseteq \mathcal{Q}$ if there is a polynomial f such that every $q \in \mathcal{Q}$ is characterised within \mathcal{Q}' by some E of size $|E| \leq f(|q|)$, where $|E| = \sum_{\mathcal{D} \in E^+ \cup E^-} |\mathcal{D}|$. Let \mathcal{Q}^n be the set of queries in \mathcal{Q} of size at most n . We say that \mathcal{Q} is *polynomially characterisable for bounded query size* if there is a polynomial f such that every $q \in \mathcal{Q}^n$ is characterised by some E of size $\leq f(n)$ within \mathcal{Q}^n .

Observe that (polynomial) characterisability is anti-monotone: if a query q is (polynomially) characterisable within \mathcal{Q} and $\mathcal{Q}' \subseteq \mathcal{Q}$, then q is (polynomially) characterisable within \mathcal{Q}' . In counterexamples to characterisability, we therefore only provide the smallest natural class of queries within which non-characterisability holds. The following examples illustrate (non-)characterisability within the classes $\mathcal{Q}_p[\diamond_r]$ and $\mathcal{Q}_p[\circ, \diamond_r]$.

Example 1. (i) Recall from Section 1 that $\diamond_r(A \wedge B)$ is not uniquely characterisable within $\mathcal{Q}_p[\diamond_r]$. The same argument shows non-characterisability of $\diamond(A \wedge B)$ within $\mathcal{Q}_p[\diamond_r, \diamond]$. On the other hand, the query $\diamond(A \wedge B)$ is characterised within $\mathcal{Q}_p[\diamond, \circ]$ by the example set (E^+, E^-) with positive examples $E^+ = \{(\emptyset, \{A, B\}), (\emptyset, \emptyset, \{A, B\})\}$ and negative $E^- = \{(\emptyset, \{A\}), (\emptyset, \{B\})\}$.

(ii) The conjunction of atoms does not always lead to non-characterisability within classes of queries with \diamond_r . For example, $q = \diamond_r(A \wedge \circ(A \wedge B))$ is characterised within $\mathcal{Q}_p[\circ, \diamond_r]$ by $E = (E^+, E^-)$ with E^+ containing two data instances $\{\{A\}, \{A, B\}\}$ and $\{\emptyset, \{A\}, \{A, B\}\}$ and E^- also two instances:

$$(\emptyset, \emptyset, \{A, B\}), \quad (\emptyset, \{A\}, \{A\}, \{B\}, \{A, B\}).$$

The intuition here is that some instances from E^- have to satisfy the query $\diamond_r(A \wedge \circ(B \wedge \diamond_r(A \wedge B)))$ as well as the query $\diamond_r(A \wedge \circ(A \wedge \diamond_r(A \wedge B)))$.

(iii) While the query $\diamond_r(A \wedge B)$ from (i) is not characterisable, there is a polynomial f such that, for all $n \in \mathbb{N}$, it is characterisable within $\mathcal{Q}_p^n[\circ, \diamond_r]$ by some E_n of size $\leq f(n)$. Namely, we take $E^+ = \{(\{A, B\}), (\emptyset, \{A, B\})\}$ and $E^- = \underbrace{\{(\{A\}, \{B\}), \dots, \{A\}, \{B\})\}}_{n \text{ times}}$.

Observe that one can always separate $q \in \mathcal{Q}[\circ, \diamond_r]$ from any other $q' \in \mathcal{Q}[\circ, \diamond_r]$ with $\text{sig}(q') \supseteq \text{sig}(q) = \sigma$ using the positive example (σ, \dots, σ) with $\text{tdp}(q)+1$ -many copies of σ . One can therefore focus on characterisability within the relevant class of queries over the same signature as the input query. However, this not the case for $\mathcal{Q}[\cup]$:

Example 2. The query $q = \perp \cup A \equiv \circ A$ is not uniquely characterisable within $\mathcal{Q}_p[\cup]$. Indeed, suppose q fits E and σ comprises all atoms occurring in E . Then $\mathcal{D}, 0 \models C \cup A$ iff $\mathcal{D}, 0 \models \circ A$, for all \mathcal{D} in E and $C \notin \sigma$, and so E does not characterise q . On the other hand, for the signature $\sigma = \{A, B\}$, the query q is characterised within $\mathcal{Q}_p^\sigma[\cup]$ by the example set (E^+, E^-) in which $E^+ = \{(\emptyset, \{A\})\}$ and $E^- = \{(\sigma, \{B\}, \{A\})\}$ as $A \cup A \equiv (A \wedge B) \cup A \equiv \circ A$.

As noted in Section 1, $\perp U A$ is not uniquely characterisable within $\mathcal{Q}^{\{A\}}[U]$ because of nested U-operators on the left-hand side of U. This observation prompts us to consider the subclass $\mathcal{Q}_-^\sigma[U]$ of $\mathcal{Q}^\sigma[U]$ -queries q in which any subquery $q' \cup q''$ does not contain occurrences of U in q' . Note that $\mathcal{Q}_-^\sigma[U] \subseteq \mathcal{Q}^\sigma[U]$. We show that $\mathcal{Q}_-^\sigma[U]$ is uniquely characterisable. To simplify notation, we give σ -data instances as *words* over the alphabet 2^σ using the standard notation of regular languages. Instead of $\mathcal{D}, 0 \models q$ we simply write $\mathcal{D} \models q$. By the semantics of U, for any $q \in \mathcal{Q}_-^\sigma[U]$, we have

$$\sigma^d \not\models q \text{ for } d \leq \text{tdp}(q), \quad \sigma^d \models q \text{ for } d > \text{tdp}(q) \quad (3)$$

where σ^d is a word with d -many σ . Note also that there are finitely-many, say $N_d < \omega$, pairwise non-equivalent queries of any depth $d < \omega$ in $\mathcal{Q}_-^\sigma(U)$.

Lemma 3. *If $q, q' \in \mathcal{Q}_-^\sigma[U]$ are of depth d and $q \not\models q'$, then there is \mathcal{D} such that $\max(\mathcal{D}) \leq N_d$, $\mathcal{D} \models q$ and $\mathcal{D} \not\models q'$.*

Proof. Consider \mathcal{D} of minimal length such that $\mathcal{D} \models q$ and $\mathcal{D} \not\models q'$. Let $tp(i)$ comprise all of the subqueries s of q and q' with $\mathcal{D}, i \models s$. By the choice of \mathcal{D} , we have $tp(i) \neq tp(j)$ for any distinct $i, j \in [0, \max(\mathcal{D})]$ (otherwise we could cut the interval $[i, j)$ out of \mathcal{D} to obtain a shorter instance separating q from q'). It follows that $\max(\mathcal{D}) \leq N_d$. \square

Theorem 4. *For any σ , $\mathcal{Q}_-^\sigma[U]$ is uniquely characterisable.*

Proof. Any $q \in \mathcal{Q}_-^\sigma(U)$ is uniquely characterised by E with

$$E^+ = \{\mathcal{D} \models q \mid \max(\mathcal{D}) \leq N_{\text{tdp}(q)}\}, \\ E^- = \{\mathcal{D} \not\models q \mid \max(\mathcal{D}) \leq N_{\text{tdp}(q)}\}.$$

Indeed, let $q' \in \mathcal{Q}_-^\sigma(U)$ fit E . Then $\text{tdp}(q') = \text{tdp}(q)$ by (3), and so $q \equiv q'$ by Lemma 3. \square

5 Characterisability in $\mathcal{Q}_p[\circ, \diamond_r]$

In this section, we prove a criterion of (polynomial) unique characterisability of queries within $\mathcal{Q}_p[\circ, \diamond_r]$. The criterion is applicable to $\mathcal{Q}_p[\circ, \diamond, \diamond_r]$ -queries in a normal form, which is defined and illustrated below.

Example 5. It is readily checked that the $\mathcal{Q}_p[\circ, \diamond_r]$ -query $q = \circ \diamond_r \circ \diamond_r (A \wedge B \wedge C \wedge \diamond_r (B \wedge \diamond_r (B \wedge C)))$ is equivalent to the $\mathcal{Q}_p[\circ, \diamond, \diamond_r]$ -query $q^{nf} = \diamond \diamond (A \wedge B \wedge C)$.

We define the normal form for $\mathcal{Q}_p[\circ, \diamond, \diamond_r]$ -queries represented as a first-order CQ by a list of atoms. For example, the query q^{nf} above is given by the CQ

$$q^{nf}(t_0) = t_0 < t_1, t_1 < t_2, A(t_2), B(t_2), C(t_2)$$

with one free (answer) variable t_0 and existentially quantified t_1 and t_2 . In general, any $q \in \mathcal{Q}_p[\circ, \diamond, \diamond_r]$ is represented as a CQ

$$\rho_0(t_0), R_1(t_0, t_1), \dots, \rho_{n-1}(t_{n-1}), R_n(t_{n-1}, t_n), \rho_n(t_n),$$

where ρ_i is a set of atoms, $\rho_i(t_i) = \{A(t_i) \mid A \in \rho_i\}$ and $R_i \in \{suc, <, \leq\}$. We divide q into *blocks* q_i such that

$$q = q_0 \mathcal{R}_1 q_1 \dots \mathcal{R}_n q_n \quad (4)$$

with $\mathcal{R}_i = R_1^i(t_0, t_1) \dots R_{n_i}^i(t_{n_i-1}, t_{n_i})$, for $R_j^i \in \{<, \leq\}$,

$$q_i = \rho_0^i(s_0^i) suc(s_0^i, s_1^i) \rho_1^i(s_1^i) \dots suc(s_{k_i-1}^i, s_{k_i}^i) \rho_{k_i}^i(s_{k_i}^i)$$

and $s_{k_i}^i = t_0^{i+1}$, $t_{n_i}^i = s_0^i$. If $k_i = 0$, the block q_i is *primitive*. A primitive block $q_i = \rho_0^i(s_0^i)$ with $i > 0$ and $|\rho_0^i| \geq 2$ is called a *lone conjunct* of q .

Example 6. The query $\diamond_r(A \wedge B)$ in Example 1(i), whose CQ representation is $t_0 \leq t_1, \rho_1(t_1)$, for $\rho_1 = \{A, B\}$, has a lone conjunct $\rho_1(t_1)$. In $\diamond_r(A \wedge \circ(A \wedge B))$ from Example 1(ii), represented as $t_0 \leq t_1, A(t_1), suc(t_1, t_2), \rho_1(t_2)$, the conjunct $\rho_1(t_2)$ is not lone.

Now, we say that q given by (4) is in *normal form* if the following conditions hold:

- (n1) $\rho_0^i \neq \emptyset$ if $i > 0$, and $\rho_{k_i}^i \neq \emptyset$ if $i > 0$ or $k_i > 0$ (thus, of all the first/last ρ in a block only ρ_0^0 can be empty);
- (n2) each \mathcal{R}_i is either a single $t_0^i \leq t_1^i$ or a sequence of $<$;
- (n3) $\rho_{k_i}^i \not\supseteq \rho_0^{i+1}$ if q_{i+1} is primitive and R_{i+1} is \leq ;
- (n4) $\rho_{k_i}^i \not\supseteq \rho_0^{i+1}$ if $i > 0$, q_i is primitive and R_{i+1} is \leq .

The queries in Example 6 are in normal form with two blocks each; the query q^{nf} above is in normal form with two blocks $q_0 = \top(t_0)$ and $q_1 = A(t_2) \wedge B(t_2) \wedge C(t_2)$.

Lemma 7. *Every query in $\mathcal{Q}_p[\circ, \diamond_r]$ is equivalent to a query in normal form that can be computed in linear time.*

A query $q \in \mathcal{Q}_p[\circ, \diamond]$ is *safe* if it is equivalent to a query $q' \in \mathcal{Q}_p[\circ, \diamond]$ in normal form not containing lone conjuncts. We are now in the position to formulate the criterion.

Theorem 8. (i) *A query $q \in \mathcal{Q}_p[\circ, \diamond_r]$ is uniquely characterisable within $\mathcal{Q}_p[\circ, \diamond_r]$ iff q is safe.*

(ii) *Those queries that are uniquely characterisable within $\mathcal{Q}_p[\circ, \diamond_r]$ are actually polynomially characterisable within $\mathcal{Q}_p[\circ, \diamond_r]$.*

(iii) *The class $\mathcal{Q}_p[\circ, \diamond_r]$ is polynomially characterisable for bounded query size.*

(iv) *The class $\mathcal{Q}_p[\circ, \diamond]$ is polynomially characterisable.*

Proof. A detailed proof is given in the full version. Here, we define a polynomial-size example set $E = (E^+, E^-)$ characterising a query q in normal form (4), which does not contain lone conjuncts. Let b be the number of \circ and \diamond in q plus 1. For each block q_i in (4), we take two words

$$\bar{q}_i = \rho_0^i \dots \rho_{k_i}^i, \quad \bar{q}_i \bowtie \bar{q}_{i+1} = \rho_0^i \dots (\rho_{k_i}^i \cup \rho_0^{i+1}) \dots \rho_{k_{i+1}}^{i+1}.$$

The set E^+ contains the data instances given by the words

- $\mathcal{D}_b = \bar{q}_0 \emptyset^b \dots \bar{q}_i \emptyset^b \bar{q}_{i+1} \dots \emptyset^b \bar{q}_n$,
- $\mathcal{D}_i = \bar{q}_0 \emptyset^b \dots \bar{q}_i \bowtie \bar{q}_{i+1} \dots \emptyset^b \bar{q}_n$ if \mathcal{R}_{i+1} is \leq ,
- $\mathcal{D}_i = \bar{q}_0 \emptyset^b \dots \bar{q}_i \emptyset^{n_{i+1}} \bar{q}_{i+1} \dots \emptyset^b \bar{q}_n$ otherwise.

Here, \emptyset^b is a sequence of b -many \emptyset and similarly for $\emptyset^{n_{i+1}}$. The set E^- contains all data instances of the form

- $\mathcal{D}_i^- = \bar{q}_0 \emptyset^b \dots \bar{q}_i \emptyset^{n_{i+1}-1} \bar{q}_{i+1} \dots \emptyset^b \bar{q}_n$ if $n_{i+1} > 1$;
- $\mathcal{D}_i^- = \bar{q}_0 \emptyset^b \dots \bar{q}_i \bowtie \bar{q}_{i+1} \dots \emptyset^b \bar{q}_n$ if \mathcal{R}_{i+1} is a single $<$,

and also the data instances obtained from \mathcal{D}_b by

- (a) removing a single atom from some $\rho_j^i \neq \emptyset$ or removing the whole $\rho_j^i = \emptyset$, for $i \neq 0$ and $j \neq 0$, from some \bar{q}_i ;
- (b) replacing $\bar{q}_i = \rho_0^i \dots \rho_l^i \rho_{l+1}^i \dots \rho_{k_i}^i$ ($k_i > 0$) by $\bar{q}'_i \emptyset^b \bar{q}''_i$, where $\bar{q}'_i = \rho_0^i \dots \rho_l^i$, $\bar{q}''_i = \rho_{l+1}^i \dots \rho_{k_i}^i$ and $l \geq 0$;
- (c) replacing some $\rho_l^i \neq \emptyset$, $0 < l < k_i$, by $\rho_l^i \emptyset^b \rho_l^i$;
- (d) replacing $\rho_{k_i}^i$ ($k_i > 0$, $|\rho_{k_i}^i| \geq 2$) with $\rho_{k_i}^i \setminus \{A\} \emptyset^b \rho_{k_i}^i$, for some $A \in \rho_{k_i}^i$, or replacing ρ_0^i ($k_i > 0$, $|\rho_0^i| \geq 2$) with $\rho_0^i \emptyset^b \rho_0^i \setminus \{A\}$, for some $A \in \rho_0^i$;
- (e) replacing $\rho_0^i \neq \emptyset$ with $\rho_0^i \setminus \{A\} \emptyset^b \rho_0^i$, for some $A \in \rho_0^i$, if $k_0 = 0$, and with $\rho_0^i \emptyset^b \rho_0^i$ if $k_0 > 0$.

The size of E is clearly polynomial in $|q|$. It is readily seen that $\mathcal{D} \models q$ for all $\mathcal{D} \in E^+$. To continue the proof sketch, note that $\mathcal{D} \models q$ iff there is a homomorphism h from the set $\text{var}(q)$ of variables in q to $[0, \max(\mathcal{D})]$, i.e., $h(t_0) = 0$, $A(h(t)) \in \mathcal{D}$ if $A(t) \in q$, $h(t') = h(t) + 1$ if $\text{suc}(t, t') \in q$, and $h(t) R h(t')$ if $R(t, t') \in q$ for $R \in \{<, \leq\}$. Using the assumption that q is in normal form, one can show that there is no homomorphism witnessing $\mathcal{D} \models q$, for any $\mathcal{D} \in E^-$.

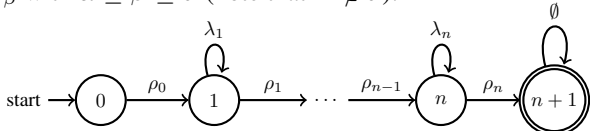
Suppose now that $q' \in \mathcal{Q}_p[\diamond, \diamond_r]$ in normal form is given and $q' \neq q$. If $\mathcal{D}_b \not\models q'$, we are done as $\mathcal{D}_b \in E^+$. Otherwise, let h be a homomorphism witnessing $\mathcal{D}_b \models q'$. Then one can show that either the restriction of h to the blocks of q' is an isomorphism onto the blocks of q or there exists a data instance \mathcal{D} obtained using one of the rules (a)–(e) such that a suitably modified h is a homomorphism from q' to \mathcal{D} . In the latter case, we are done as $\mathcal{D} \in E^-$ and $\mathcal{D} \models q'$. In the former case, q and q' coincide with the exception of the sequences of \diamond and \diamond_r between blocks. Then q can be separated from q' using the examples \mathcal{D}_i and \mathcal{D}_i^- . \square

6 Characterisability in $\mathcal{Q}_p^\sigma[U]$

LTL-queries with U do not correspond to CQs (because of the universal quantification in its semantics), and so require a different approach. We view them as defining regular languages. With each $\mathcal{Q}_p^\sigma[U]$ -query of the form (2) we associate the following regular expression over the alphabet 2^σ :

$$q = \rho_0 \lambda_1^* \rho_1 \lambda_2^* \dots \lambda_n^* \rho_n \lambda_{n+1}^* \quad (5)$$

where $\lambda_{n+1} = \emptyset$ and $\perp^* = \varepsilon$. We regard the words of the language $L(q)$ over 2^σ as data instances. Clearly, $\mathcal{D}' \models q$ iff there is $\mathcal{D} \in L(q)$ such that $\mathcal{D} \in \mathcal{D}'$, i.e., $\mathcal{D} = (\delta_0, \dots, \delta_k)$ and $\mathcal{D}' = (\delta'_0, \dots, \delta'_k)$, for some $k < \omega$, and $\delta_i \subseteq \delta'_i$, for all $i \leq k$. The language L_q of all σ -data instances $\mathcal{D} \models q$ (regarded as words over 2^σ) can be given by the NFA \mathfrak{A}_q below, where each \rightarrow_α , for $\alpha \neq \perp$, stands for all transitions \rightarrow_β with $\alpha \subseteq \beta \subseteq \sigma$ (note that $\perp \notin \sigma$):



Without loss of generality we assume that all our q are *minimal* in the sense that by replacing any $\lambda_i \neq \perp$ with \perp in q we obtain a query that is *not equivalent* to q . For example, in minimal q , $\rho_j \supseteq \dots \supseteq \rho_i \supseteq \lambda_i$ and $\lambda_l = \perp$ for all $l \in (j, i)$ imply $\rho_j \not\subseteq \lambda_j$ as otherwise $\lambda_j U(\rho_j \wedge (\perp U \dots (\lambda_i U \varphi) \dots))$

is equivalent to $\perp U(\rho_j \wedge (\perp U \dots (\lambda_i U \varphi) \dots))$. Using standard automata-theoretic techniques, one can show:

Theorem 9. Any $\mathcal{Q}_p^\sigma[U]$ -queries $q \neq q'$ can be separated by some \mathcal{D} with $\max(\mathcal{D}) \leq O((\min\{\text{tdp}(q), \text{tdp}(q')\})^2)$.

Using Theorem 9 in the proof of Theorem 4 we obtain:

Corollary 10. The class $\mathcal{Q}_p^\sigma[U]$ is exponentially characterisable within $\mathcal{Q}_p^\sigma[U]$.

The following examples illustrate difficulties in finding short unique characterisations of $\mathcal{Q}_p^\sigma[U]$ -queries, namely, that in general, data instances of different shapes and forms are needed to separate $\mathcal{Q}_p^\sigma[U]$ -queries. To unclutter notation we omit $\{\}$ in singletons like $\{A\}$.

Example 11. (a) The shortest data instance separating

$$q = X \emptyset^* A \perp^* B \perp^* A B^* A A^* B \emptyset^*,$$

$$q' = X \emptyset^* A \perp^* B A^* A B^* A \perp^* B \emptyset^*$$

is $\mathcal{D} = X A B A B B A A B$ with $\mathcal{D} \models q$ and $\mathcal{D} \not\models q'$ (e.g., $X A B A B A A B$ satisfies both q and q').

(b) For $l > 0$, let $q_l = (A B^*)^l A A^* B B^*$. Then

$$X A^* q_{l_1} q_{l_2} \dots q_{l_k} X \emptyset^* \neq X \perp^* q_{l_1} q_{l_2} \dots q_{l_k} X \emptyset^*,$$

$$X A^* q_{l_1} q_{l_2} \dots q_{l_k} A \emptyset^* \equiv X \perp^* q_{l_1} q_{l_2} \dots q_{l_k} A \emptyset^*.$$

If $l_1 \leq \dots \leq l_k$, the former inequality is witnessed by the instance $X A^{l_1} B A^{l_2} B \dots A^{l_k} B A^{l_k} B X$. Less generally, $X A^* q_1 q_2 X \emptyset^* \neq X \perp^* q_1 q_2 X \emptyset^*$ can be shown by $X A A B A A A B A A A B X$ or by $X A A B A B A A B A B X$ (spot the difference and see (n₂) below).

Here, we consider the class $\mathcal{P}^\sigma[U]$ of *peerless* queries given by (5), in which, for any i , either $\lambda_i = \perp$ or the sets λ_i and ρ_i are *incomparable* with respect to \subseteq . Our main result is that $\mathcal{P}^\sigma[U]$ is polynomially characterisable within $\mathcal{Q}_p^\sigma[U]$.

We start with a general observation. Consider two queries $q = \rho_0 \lambda_1^* \dots \lambda_n^* \rho_n \emptyset^*$ and $q' = \rho_0 \mu_1^* \dots \mu_n^* \rho_n \emptyset^*$. We say that $\lambda_i \neq \perp$ *subsumes* $\mu_j \neq \perp$ if either $i = j$ and $\mu_j \subseteq \lambda_i$, or $j < i$ and $\mu_j \rho_j \dots \rho_{i-1} \subseteq \rho_j \dots \rho_{i-1} \lambda_i$, or $j > i$ and $\rho_i \dots \rho_{j-1} \mu_j \subseteq \lambda_i \rho_i \dots \rho_{j-1}$. In the last two cases,

$$\mu_j \subseteq \rho_j \subseteq \dots \subseteq \rho_{i-1} \subseteq \lambda_i, \quad \mu_j \subseteq \rho_{j-1} \subseteq \dots \subseteq \rho_i \subseteq \lambda_i,$$

respectively. Note that, for peerless q , the last inclusion is impossible. If λ_i and μ_j subsume each other, in which case $\lambda_i = \mu_j$, we call (λ_i, μ_j) a *matching pair*. Observe also that, for $\mathcal{D}_q^i = \rho_0 \dots \rho_{i-1} \lambda_i \rho_i \dots \rho_n$, if $\mathcal{D}_q^i \models q'$, then λ_i subsumes some μ_j : $\rho_0 \dots \rho_n \emptyset \in \mathcal{D}_q^i$ means that λ_i subsumes $\mu_{n+1} = \emptyset$, and $\rho_0 \dots \mu_j \dots \rho_n \in \mathcal{D}_q^i$ that λ_i subsumes μ_j . The proof of the next lemma is given in the full version:

Lemma 12. For any queries q and q' as above, either (i) each $\lambda_i \neq \perp$ subsumes μ_j occurring in some matching pair (λ_k, μ_j) or (ii) q and q' are separated by a data instance of the form \mathcal{D}_q^i or $\mathcal{D}_{q'}^j$. Also, if q is peerless, λ_i can only subsume μ_j in the matching pair (λ_i, μ_j) with $i \geq j$, in which case $\mu_j = \rho_j = \dots = \rho_{i-1} = \lambda_i$.

Note that the number of data instances of the form \mathcal{D}_q^i for all possible $\mathcal{Q}_p^\sigma[U]$ -queries q' can be exponential in $|\sigma|$. The following example indicates how to overcome this issue.

Example 13. Let $\sigma = \{A, B, C, D, X\}$. To separate the query $X\{C, D\}^*A\emptyset^*$ from any $X\lambda^*A\emptyset^*$ with $A, D \notin \lambda$, we can use $\mathcal{D} = X\sigma \setminus \{A, D\}A$.

Theorem 14. *The class $\mathcal{P}^\sigma[\mathbb{U}]$ is polynomially characterisable within $\mathcal{Q}_p^\sigma[\mathbb{U}]$.*

Proof. We show that any $\mathbf{q} = \rho_0\lambda_1^*\rho_1\lambda_2^*\dots\lambda_n^*\rho_n\emptyset^*$ in $\mathcal{P}^\sigma[\mathbb{U}]$ is characterised by the example set $E = (E^+, E^-)$, where E^+ contains all data instances of the following forms:

- (p₀) $\rho_0 \dots \rho_n$,
- (p₁) $\rho_0 \dots \rho_{i-1}\lambda_i\rho_i \dots \rho_n = \mathcal{D}_q^i$,
- (p₂) $\rho_0 \dots \rho_{i-1}\lambda_i^k\rho_i \dots \rho_{j-1}\lambda_j\rho_j \dots \rho_n = \mathcal{D}_{i,k}^j$, for $i < j$ and $k = 1, 2$;

and E^- has all instances that are *not* in $L(\mathbf{q})$ of the forms:

- (n₀) σ^n and $\sigma^{n-i}\sigma \setminus \{A\}\sigma^i$, for $A \in \rho_i$,
- (n₁) $\rho_0 \dots \rho_{i-1}\sigma \setminus \{A, B\}\rho_i \dots \rho_n$, for $A \in \lambda_i \cup \{\perp\}$ and $B \in \rho_i \cup \{\perp\}$,
- (n₂) for all i and $A \in \lambda_i \cup \{\perp\}$, *some* data instance

$$\mathcal{D}_A^i = \rho_0 \dots \rho_{i-1}(\sigma \setminus \{A\})\rho_i\lambda_{i+1}^{k_{i+1}} \dots \lambda_n^{k_n}\rho_n, \quad (6)$$

if any, such that $\max(\mathcal{D}_A^i) \leq (n+1)^2$ and $\mathcal{D}_A^i \not\models \mathbf{q}^\dagger$ for \mathbf{q}^\dagger obtained from \mathbf{q} by replacing λ_j , for all $j \leq i$, with \perp . Note that $\mathcal{D}_A^i \not\models \mathbf{q}$ for peerless \mathbf{q} .

By definition, \mathbf{q} fits E and $|E|$ is polynomial in $|\mathbf{q}|$. We prove in the full version that E uniquely characterises \mathbf{q} . \square

One reason why this construction does not generalise to the whole $\mathcal{Q}_p^\sigma[\mathbb{U}]$ is that $\mathcal{D}_A^i \not\models \mathbf{q}^\dagger$ does not imply $\mathcal{D}_A^i \not\models \mathbf{q}$ for non-peerless \mathbf{q} , as shown by the following example:

Example 15. Let $\mathbf{q} = XA^*AB^*A\perp^*AB^*AA^*BB^*X\emptyset^*$. For any data instance \mathcal{D}_\perp^3 satisfying (6)—for example, $\mathcal{D}_\perp^3 = XAA\sigma ABABX$ —we have $\mathcal{D}_\perp^3 \models \mathbf{q}$.

7 Characterisability in $\mathcal{Q}[\diamond]$

In the previous two sections, we have investigated characterisability of path-shaped queries. Here, we first justify that restriction by exhibiting two examples that show how temporal branching can destroy polynomial characterisability in $\mathcal{Q}[\diamond]$. Both examples make use of *unbalanced* queries, in which different branches have different length. We then show that this is no accident: one can at least partially restore polynomial characterisability for classes without unbalanced queries.

We start by observing that, without loss of generality, it is enough to consider conjunctions of path queries only:

Lemma 16. *For every $\mathbf{q} \in \mathcal{Q}[\diamond]$, one can compute in polynomial time an equivalent query of the form $\mathbf{q}_1 \wedge \dots \wedge \mathbf{q}_n$ with $\mathbf{q}_i \in \mathcal{Q}_p[\diamond]$, for $1 \leq i \leq n$.*

The first example showing non-polynomial characterisability is rather straightforward but requires unbounded branching and an unbounded number of atoms. We write queries $\mathbf{q} \in \mathcal{Q}_p^\sigma[\diamond]$ of the form

$$\mathbf{q} = \rho_0 \wedge \diamond(\rho_1 \wedge \diamond(\rho_2 \wedge \dots \wedge \diamond\rho_n)) \quad (7)$$

as words $\rho_0\rho_1 \dots \rho_n$ over 2^σ (omitting but not forgetting $\lambda_i^* = \emptyset^*$ from (5)) and also use $\rho_0\rho_1 \dots \rho_n$ to denote the data instance defined by \mathbf{q} .

Example 17. Consider $\mathbf{q}_n = s_1 \wedge \dots \wedge s_n$, where $n \geq 2$ and each s_i is a word repeating n times the sequence $A_1 \dots A_n$ (of singletons) with omitted A_i . Now, consider the queries $\mathbf{q}_n^p = \mathbf{q}_n \wedge \mathbf{p}$, where $\mathbf{p} = \diamond(A_{i_1} \wedge \dots \wedge \diamond(A_{i_2} \wedge \dots \wedge \diamond A_{i_n}))$ and $A_{i_1} \dots A_{i_n}$ is a permutation of $A_1 \dots A_n$. Then $\mathbf{q}_n^p \models \mathbf{q}_n$ and $\mathbf{q}_n \not\models \mathbf{q}_n^p$ as shown by the data instance $s_{i_1}s_{i_2} \dots s_{i_n}$. Moreover, if $\mathcal{D} \models \mathbf{q}_n$, $\mathcal{D} \not\models \mathbf{p}$ and $\mathbf{p}' \neq \mathbf{p}$, then $\mathcal{D} \models \mathbf{p}'$. It follows that, in any $E = (E^+, E^-)$ uniquely characterising \mathbf{q}_n , the set E^+ contains at least $n!$ data instances.

The class $\mathcal{Q}_{\leq n}[\diamond]$ of queries of *branching factor* at most n contains all queries in $\mathcal{Q}[\diamond]$ that are equivalent to a query of the form $\mathbf{q}_1 \wedge \dots \wedge \mathbf{q}_m$ with $m \leq n$ and $\mathbf{q}_i \in \mathcal{Q}_p[\diamond]$. We next provide an example of non-polynomial characterisability that requires four atoms and bounded branching only.

Example 18. Let $\sigma = \{A_1, A_2, B_1, B_2\}$, $\mathbf{q}_1 = \emptyset(s\sigma)^n s$, and $\mathbf{q}_2 = \emptyset\sigma^{2n+1}$, where $s = \{A_1, A_2\}\{B_1, B_2\}$. Consider the set P of 2^{n+1} -many queries of the form $\emptyset s_{i_1} \dots s_{i_{n+1}}$ with s_i either $\{A_1\}\{A_2\}$ or $\{B_1\}\{B_2\}$. Then $\mathbf{q}_1 \wedge \mathbf{q}_2 \not\models \mathbf{q}$ for any $\mathbf{q} \in P$ and, for any \mathcal{D} with $\mathcal{D} \models \mathbf{q}_1 \wedge \mathbf{q}_2$, there is at most one $\mathbf{q} \in P$ with $\mathcal{D} \not\models \mathbf{q}$ (the proof is rather involved). It follows that $\mathbf{q}_1 \wedge \mathbf{q}_2 \not\models \mathbf{q}_1 \wedge \mathbf{q}_2 \wedge \mathbf{q}$ for all $\mathbf{q} \in P$, but 2^{n+1} positive examples are needed to separate $\mathbf{q}_1 \wedge \mathbf{q}_2$ from all $\mathbf{q}_1 \wedge \mathbf{q}_2 \wedge \mathbf{q}$ with $\mathbf{q} \in P$.

We next identify polynomially characterisable classes of $\mathcal{Q}[\diamond]$ -queries, assuming as before that $\rho_n \neq \emptyset$ in any \mathbf{q} of the form (1). We call a query $\mathbf{q}_1 \wedge \dots \wedge \mathbf{q}_n \in \mathcal{Q}[\diamond]$ with $\mathbf{q}_1, \dots, \mathbf{q}_n \in \mathcal{Q}_p[\diamond]$ *balanced* if all \mathbf{q}_i have the same depth; further, we call it *simple* if, in each \mathbf{q}_i given by (1), $|\rho_j| = 1$ for all j . Let $\mathcal{Q}_b[\diamond]$ denote the class of queries in $\mathcal{Q}[\diamond]$ that are equivalent to a balanced query.

Theorem 19. (i) *The class of simple queries in $\mathcal{Q}_b[\diamond]$ is polynomially characterisable within $\mathcal{Q}_b[\diamond]$.*

(ii) *For any n , the class $\mathcal{Q}_b[\diamond] \cap \mathcal{Q}_{\leq n}[\diamond]$ is polynomially characterisable.*

Proof. Let $\mathbf{q} \in \mathcal{Q}_p^\sigma[\diamond]$. We start with a lemma on the existence of polynomial-size σ -data instances $\mathcal{D}_{\mathbf{q},k}$ such that $\mathcal{D}_{\mathbf{q},k} \not\models \mathbf{q}$ and $\mathcal{D}_{\mathbf{q},k} \models \mathbf{q}'$ for all $\mathbf{q}' \in \mathcal{Q}_p^\sigma[\diamond]$ with $\mathbf{q}' \not\models \mathbf{q}$ and $\text{tdp}(\mathbf{q}') \leq k$. Note that such $\mathcal{D}_{\mathbf{q},k}$ do not exist in general.

Example 20. Let $\mathbf{q} = A \wedge B$. Then $A \not\models \mathbf{q}$ and $B \not\models \mathbf{q}$ but there does not exist any $\mathcal{D}_{\mathbf{q},0}$ such that $\mathcal{D}_{\mathbf{q},0} \not\models \mathbf{q}$ and $\mathcal{D}_{\mathbf{q},0} \models A$ and $\mathcal{D}_{\mathbf{q},0} \models B$.

In the following lemma we therefore assume that \mathbf{q} does not speak about the initial timepoint.

Lemma 21. *Let $\mathbf{q} \in \mathcal{Q}_p^\sigma[\diamond]$ be of the form $\diamond\mathbf{q}'$ and let $k > 0$. Then one can construct in polynomial time a σ -data instance $\mathcal{D}_{\mathbf{q},k}$ such that $\mathcal{D}_{\mathbf{q},k} \not\models \mathbf{q}$ and $\mathcal{D}_{\mathbf{q},k} \models \mathbf{q}'$ for all $\mathbf{q}' \in \mathcal{Q}_p^\sigma[\diamond]$ with $\mathbf{q}' \not\models \mathbf{q}$ and $\text{tdp}(\mathbf{q}') \leq k$.*

Proof. Assuming that $\mathbf{q} = \diamond(\rho_1 \wedge \diamond(\rho_2 \wedge \dots \wedge \diamond\rho_n))$ with $\rho_i = \{A_1^i, \dots, A_{n_i}^i\}$ for $i \geq 1$, we set

$$\mathcal{D}_{\mathbf{q},k} = \sigma s_1^k \sigma \dots \sigma s_{n-1}^k \sigma s_n^k,$$

where $s_i = \sigma \setminus \{A_1^i\} \dots \sigma \setminus \{A_{n_i}^i\}$. One can show by induction that $\mathcal{D}_{q,k}$ is as required. \square

Using Lemma 21, for any $q \in \mathcal{Q}^\sigma[\diamond]$, one can construct a polynomial-size set of negative examples as follows. Suppose $q = q_1 \wedge \dots \wedge q_n \in \mathcal{Q}^\sigma[\diamond]$ with

$$q_i = \rho_0^i \wedge \diamond(\rho_1^i \wedge \diamond(\rho_2^i \wedge \dots \wedge \diamond \rho_{n_i}^i)).$$

Let $\rho = \bigwedge_{i=1}^n \rho_0^i$ and let q_i^- be q_i without the conjunct ρ_0^i , so Lemma 21 is applicable to q_i^- . Now let $E_{q_i^-,m}^-$ contain the σ -data instances $\mathcal{D}_{q_i^-,m}^-$ and $\sigma \setminus \{A\} \sigma^m$ for all $A \in \rho$.

Lemma 22. (i) For any $\mathcal{D} \in E_{q,m}^-$, we have $\mathcal{D} \not\models q$.

(ii) For any $q' \in \mathcal{Q}^\sigma[\diamond]$ with $q' \not\models q$ and $\text{tdp}(q') \leq m$, there exists $\mathcal{D} \in E_{q,m}^-$ with $\mathcal{D} \models q'$.

It follows from Lemma 22 that non-polynomial characterisability of $\mathcal{Q}[\diamond]$ -queries can only be caused by the need for super-polynomially-many positive examples. We now discuss the construction of positive examples in the proof of Theorem 19 (ii); part (i) is dealt with in the full version. Let $q = q_1 \wedge \dots \wedge q_m \in \mathcal{Q}_b^\sigma[\diamond] \cap \mathcal{Q}_{\leq n}^\sigma[\diamond]$ with $m \leq n$ and

$$q_i = \rho_0^i \wedge \diamond(\rho_1^i \wedge \diamond(\rho_2^i \wedge \dots \wedge \diamond \rho_N^i)).$$

For any map $f: \{1, \dots, m\} \rightarrow \{1, \dots, N\}$, construct a σ -data instance \mathcal{D}_f by inserting $\rho_{f(i)}^i$ into the data instance σ^N in position $f(i)$. Let E^+ contain the data instance $\rho \sigma^N$ for $\rho = \bigcup_{i=1}^m \rho_0^i$ and all the data instances \mathcal{D}_f . One can show that (E^+, E^-) characterises q in $\mathcal{Q}_b[\diamond] \cap \mathcal{Q}_{\leq n}[\diamond]$. \square

8 2D Temporal Instance Queries

Now we consider ‘two-dimensional’ query languages that combine instance queries (over the object domain) in the standard description logics \mathcal{EL} and \mathcal{ELI} (Baader et al. 2017) with the *LTL*-queries (over the temporal domain) considered above. Our aim is to understand how far the characterisability results of the previous sections can be generalised to the 2D case. A *relational signature* is a finite set $\Sigma \neq \emptyset$ of unary and binary predicate symbols. A Σ -data instance \mathcal{A} is a finite set of *atoms* $A(a)$ and $P(a, b)$ with $A, P \in \Sigma$ and *individual names* a, b . Let $\text{ind}(\mathcal{A})$ be the set of individual names in \mathcal{A} . We assume that $P^-(a, b) \in \mathcal{A}$ iff $P(b, a) \in \mathcal{A}$, calling P^- the *inverse* of P (with $P^{- -} = P$). Let $S := P \mid P^-$. *Temporal instance queries* are defined by the grammar

$$q := \top \mid \perp \mid A \mid \exists S.q \mid q_1 \wedge q_2 \mid \text{op } q \mid q_1 \cup q_2,$$

where $\text{op} \in \{\circ, \diamond, \diamond_r\}$. Such queries without temporal operators are called *ELI-queries*; those of them without inverses P^- are *EL-queries*. A *temporal Σ -data instance* \mathcal{D} is a finite sequence $\mathcal{A}_0, \dots, \mathcal{A}_n$ of Σ -data instances. We set $\text{ind}(\mathcal{D}) = \bigcup_{i=1}^n \text{ind}(\mathcal{A}_i)$. For any $\ell \in \mathbb{N}$ and $a \in \text{ind}(\mathcal{D})$, the *truth-relation* $\mathcal{D}, a, \ell \models q$ is defined by induction:

$$\mathcal{D}, a, \ell \models A \text{ iff } A(a) \in \mathcal{A}_i,$$

$$\mathcal{D}, a, \ell \models \exists S.q \text{ iff there is } b \in \text{ind}(\mathcal{A}_i) \text{ such that}$$

$$S(a, b) \in \mathcal{A}_i \text{ and } \mathcal{D}, b, \ell \models q,$$

with the remaining clauses being obvious generalisations of the *LTL* ones. An *example set* is a pair $E = (E^+, E^-)$

with finite sets E^+ and E^- of pointed temporal data instances \mathcal{D}, a such that $a \in \text{ind}(\mathcal{D})$. We say that q fits E if $\mathcal{D}^+, a^+, 0 \models q$ and $\mathcal{D}^-, a^-, 0 \not\models q$, for all $\mathcal{D}^+, a^+ \in E^+$ and $\mathcal{D}^-, a^- \in E^-$. As before, E *uniquely characterises* q if q fits it and every q' fitting E is logically equivalent to q .

We need the following result on the unique characterisability of \mathcal{ELI} -queries.

Theorem 23 (ten Cate and Dalmau 2021). *The class of \mathcal{ELI} -queries is polynomially characterisable.*

Theorem 24 is shown by constructing frontiers in the set of \mathcal{ELI} -queries partially ordered by entailment, where a set \mathcal{F} of \mathcal{ELI} -queries is called a *frontier* of an \mathcal{ELI} -query q if the following hold:

- $q \models q'$ and $q' \not\models q$, for all $q' \in \mathcal{F}$;
- if $q \models q''$ for some \mathcal{ELI} -query q'' , then $q'' \models q$ or there exists $q' \in \mathcal{F}$ with $q' \models q''$.

Theorem 24 (ten Cate and Dalmau 2021). *A frontier $\mathcal{F}(q)$ of any \mathcal{ELI} -query q can be computed in polynomial time.*

Theorem 23 follows from Theorem 24. For any \mathcal{ELI} -query q we denote by \hat{q} the tree-shaped data instance defined by q with designated root a . Then q is characterised by E with $E^+ = \{\hat{q}\}$ and $E^- = \{\hat{r} \mid r \in \mathcal{F}(q)\}$.

For any unrestricted temporal query language $\mathcal{Q}[\Phi]$ and $\mathcal{L} \in \{\mathcal{EL}, \mathcal{ELI}\}$, we denote by $\mathcal{Q}[\Phi] \otimes \mathcal{L}$ the set of all temporal instance queries with operators in Φ with (for \mathcal{ELI}) or without (for \mathcal{EL}) inverse predicates. We generalise the path-shaped queries $\mathcal{Q}_p[\Phi]$ as follows: $\mathcal{Q}_p[\Phi] \otimes \mathcal{L}$ denotes the class of queries q in $\mathcal{Q}[\Phi] \otimes \mathcal{L}$ such that, for any subquery $q_1 \wedge q_2$ of q , either q_1 or q_2 do not have an occurrence of any operator in Φ that is not in the scope of $\exists S$. To illustrate, $\exists S.\diamond A \wedge \diamond \exists S.A$ is in $\mathcal{Q}_p[\Phi] \otimes \mathcal{L}$, but $\diamond A \wedge \diamond \exists S.A$ is not. We make two observations about unique characterisability in these ‘full’ combinations.

Theorem 25. (i) $\mathcal{Q}[\circ, \diamond] \otimes \mathcal{EL}$ is uniquely characterisable.

(ii) $\mathcal{Q}_p[\circ] \otimes \mathcal{ELI}$ and $\mathcal{Q}_p[\diamond] \otimes \mathcal{ELI}$ are polynomially characterisable.

Here, (i) is shown similarly to Theorem 4 (it remains open whether it can be extended to $\mathcal{Q}[\circ, \diamond] \otimes \mathcal{ELI}$); (ii) is proved by generalising Theorem 24 to temporal data instances.

We now show that the application of the DL constructor $\exists P$ to temporal queries with both \circ and \diamond destroys polynomial characterisability. Denote by $\mathcal{EL}(\mathcal{Q}_p[\circ, \diamond])$ the class of queries in $\mathcal{Q}_p[\circ, \diamond] \otimes \mathcal{EL}$ that contain no $\exists P$ in the scope of a temporal operator.

Theorem 26. $\mathcal{EL}(\mathcal{Q}_p[\circ, \diamond])$ is not polynomially characterisable.

Proof. Consider the queries $q_n = \exists P.q_1^n \wedge \dots \wedge \exists P.q_n^n$, in which each q_i^n corresponds to the regular expression

$$\underbrace{BB\emptyset^*A}_1 \dots \underbrace{BB\emptyset^*A\emptyset^*B\emptyset^*A}_{i-1} \underbrace{\emptyset^*B\emptyset^*A}_{i} \underbrace{BB\emptyset^*A}_{i+1} \dots \underbrace{BB\emptyset^*A\emptyset^*}_{n}$$

(with omitted $\perp^* = \varepsilon$ in BB). One can show that any unique characterisation of q_n contains at least 2^n positive examples to separate it from all queries $q_n \wedge \exists P.s$ with

$$s = \mathbf{o}_1(B \wedge \diamond(A \wedge \mathbf{o}_2(B \wedge \diamond(A \wedge \dots \wedge \mathbf{o}_n(B \wedge \diamond A) \dots))))),$$

where \mathbf{o}_i is \circ or $\diamond \circ$ if $i > 1$, and blank or \diamond if $i = 1$. \square

The situation changes drastically if we do not admit temporal operators in the scope of $\exists P$. We start by investigating the class $\mathcal{Q}_p[\circ, \diamond_r](\mathcal{ELI})$ of queries of the form

$$\mathbf{q} = \mathbf{r}_0 \wedge \mathbf{o}_1(\mathbf{r}_1 \wedge \mathbf{o}_2(\mathbf{r}_2 \wedge \dots \wedge \mathbf{o}_n \mathbf{r}_n)),$$

where the \mathbf{r}_i are \mathcal{ELI} -queries and $\mathbf{o}_i \in \{\circ, \diamond_r\}$. We can generalise the CQ-representation, the normal form, and the notion of lone conjunct from $\mathcal{Q}_p[\circ, \diamond_r]$ to $\mathcal{Q}_p[\circ, \diamond_r](\mathcal{ELI})$ in a straightforward way. To formulate conditions (n1)–(n4), we replace the set inclusions ' $\rho_i \subseteq \rho_j$ ' by entailment ' $\mathbf{r}_i \models \mathbf{r}_j$ '. For example, (n4) becomes

(n4') $\mathbf{r}_0^{i+1} \not\models \mathbf{r}_{k_i}^i$ if $i > 0$, \mathbf{q}_i is primitive and R_{i+1} is \leq .

The condition for lone conjuncts now requires that \mathbf{r} is not equivalent to any $\mathbf{q}_1 \wedge \mathbf{q}_2$ with \mathcal{ELI} -queries $\mathbf{q}_1, \mathbf{q}_2$ such that $\mathbf{q}_i \not\models \mathbf{r}$ for $i = 1, 2$. Then one can show again that every $\mathcal{Q}_p[\circ, \diamond_r](\mathcal{ELI})$ -query is equivalent to a query in normal form, which can be computed in polynomial time.

Theorem 27. *The statements of Theorem 8 (i)–(iv) also hold if one replaces $\mathcal{Q}_p[\circ, \diamond_r]$ by $\mathcal{Q}_p[\circ, \diamond_r](\mathcal{ELI})$.*

The proof generalises the example set defined in Theorem 8 using the frontiers provided by Theorem 24 as a *black box*. Indeed, in the definition of examples replace ρ_i by $\hat{\mathbf{r}}_i$, the data instance corresponding to the \mathcal{ELI} -query \mathbf{r}_i , and replace ' $\rho \setminus \{A\}$ for $A \in \rho$ ' by 'the data instance corresponding to a query in $\mathcal{F}(\mathbf{r})$ '. We choose a single individual, a , as the root of these data instances. For example, item (a) becomes:

(a') replacing some \mathbf{r}_j^i by the data instance corresponding to a query in $\mathcal{F}(\mathbf{r}_j^i)$ or removing the whole $\mathbf{r}_j^i = \emptyset$ for $i \neq 0$ and $j \neq 0$ from some \mathbf{q}_i .

Next consider the class $\mathcal{Q}_p[\cup](\mathcal{L})$ of queries of the form

$$\mathbf{q} = \mathbf{r}_0 \wedge (\mathbf{l}_1 \cup (\mathbf{r}_1 \wedge (\mathbf{l}_2 \cup (\dots (\mathbf{l}_n \cup \mathbf{r}_n) \dots))), \quad (8)$$

where \mathbf{r}_i is an \mathcal{L} -query and \mathbf{l}_i is either an \mathcal{L} -query or \perp , for $\mathcal{L} \in \{\mathcal{EL}, \mathcal{ELI}\}$. For the same reason as in the 1D case, we fix a finite signature Σ of predicate symbols. Denote by $\mathcal{L}(\Sigma)$ and $\mathcal{Q}_p[\cup](\mathcal{L})$ the set of queries in \mathcal{L} and $\mathcal{Q}_p^\Sigma[\cup](\mathcal{L})$, respectively, with predicate symbols in Σ . Aiming to generalise Theorem 14, we again translate set-inclusion to entailment, so the *peerless queries* $\mathcal{P}^\Sigma[\cup](\mathcal{L})$ take the form (8) such that either $\mathbf{l}_i = \perp$ or $\mathbf{l}_i \not\models \mathbf{r}_i$ and $\mathbf{r}_i \not\models \mathbf{l}_i$.

Theorem 28. *Let Σ be a finite relational signature. Then $\mathcal{P}^\Sigma[\cup](\mathcal{EL})$ is polynomially characterisable within $\mathcal{Q}_p^\Sigma[\cup](\mathcal{EL})$, while $\mathcal{P}^\Sigma[\cup](\mathcal{ELI})$ is exponentially characterisable within $\mathcal{Q}_p^\Sigma[\cup](\mathcal{ELI})$.*

To prove Theorem 28, we generalise the example set from the proof of Theorem 14. The positive examples are straightforward: simply replace ρ_i and λ_i by the data instances corresponding to \mathbf{r}_i and \mathbf{l}_i (and choose a single root individual). For the negative examples, we have to generalise the construction of σ , $\sigma \setminus \{A\}$, and $\sigma \setminus \{A, B\}$. For σ this is straightforward as its role can now be played by the Σ -data instance $\mathcal{A}_\Sigma = \{A(a), R(a, a) \mid A, R \in \Sigma\}$ for which $\mathcal{A}_\Sigma \models \mathbf{q}(a)$ for all $\mathbf{q} \in \mathcal{ELI}(\Sigma)$. For $\sigma \setminus \{A\}$ and $\sigma \setminus \{A, B\}$, we require *split partners* defined as follows. Let Q be a finite set of $\mathcal{L}(\Sigma)$ -queries. A set $\mathcal{S}(Q)$ of pointed Σ -data instances (\mathcal{A}, a) is called a *split partner* of Q in $\mathcal{L}(\Sigma)$ if the following conditions are equivalent for all $\mathcal{L}(\Sigma)$ -queries \mathbf{q}' :

- $\mathcal{A} \models \mathbf{q}'(a)$ for some $(\mathcal{A}, a) \in \mathcal{S}(Q)$;
- $\mathbf{q}' \not\models \mathbf{q}$ for all $\mathbf{q} \in Q$.

Example 29. The split partner $\mathcal{S}(\{A\})$ of $\{A\}$ in $\mathcal{EL}(\Sigma)$ is the singleton set containing $\mathcal{A}_{\Sigma}^{-A}$ defined as

$$\{B(a), R(a, b), R(b, b), B'(b) \mid B \in \Sigma \setminus \{A\}, R, B' \in \Sigma\}.$$

Theorem 30. *Let $n > 0$ be fixed. For every set Q of $\mathcal{EL}(\Sigma)$ -queries with $|Q| \leq n$, one can compute in polynomial time a split partner $\mathcal{S}(Q)$ of Q in $\mathcal{EL}(\Sigma)$. For \mathcal{ELI} , one can compute a split partner in exponential time.*

The proof, given in the full version, requires (as does \mathcal{A}_Σ) the construction of non-tree-shaped data instances. The exponential upper bound for \mathcal{ELI} follows from results on generalised dualities for homomorphisms between relational structures (Foniok, Nesetril, and Tardif 2008; Nesetril and Tardif 2005). It remains open whether a polynomial construction is possible for our particular case.

We obtain the negative examples for \mathbf{q} of the form (8) by taking the following data instances \mathcal{D} :

(n₀') \mathcal{A}_Σ^n and $\mathcal{A}_\Sigma^{n-i} \mathcal{A}_\Sigma^i$, for $\mathcal{A} \in \mathcal{S}(\{\mathbf{r}_i\})$;

(n₁') $\mathcal{D} = \hat{\mathbf{r}}_0 \dots \hat{\mathbf{r}}_{i-1} \mathcal{A} \hat{\mathbf{r}}_i \dots \hat{\mathbf{r}}_n$ such that $\mathcal{D} \not\models \mathbf{q}$, where $\mathcal{A} \in \mathcal{S}(\{\mathbf{l}_i, \mathbf{r}_i\}) \cup \mathcal{S}(\{\mathbf{l}_i\}) \cup \mathcal{S}(\{\mathbf{r}_i\}) \cup \{\mathcal{A}_\Sigma\}$;

(n₂') for all i and $\mathcal{A} \in \mathcal{S}(\{\mathbf{l}_i\}) \cup \{\mathcal{A}_\Sigma\}$, some data instance

$$\mathcal{D}_{\mathcal{A}}^i = \hat{\mathbf{r}}_0 \dots \hat{\mathbf{r}}_{i-1} \mathcal{A} \hat{\mathbf{r}}_i \hat{\mathbf{l}}_{i+1}^{k_{i+1}} \hat{\mathbf{r}}_{i+1} \dots \hat{\mathbf{l}}_n^{k_n} \mathbf{r}_n,$$

if any, such that $\max(\mathcal{D}_{\mathcal{A}}^i) \leq (n+1)^2$ and $\mathcal{D}_{\mathcal{A}}^i \not\models \mathbf{q}^\dagger$ for \mathbf{q}^\dagger obtained from \mathbf{q} by replacing all \mathbf{l}_j , for $j \leq i$, with \perp .

We illustrate the construction by generalising Example 2.

Example 31. For $\mathbf{q} = \circ A$ and any relational signature Σ containing A , we obtain, after removing redundant data instances, E^+ with $\emptyset\{A(a)\}$ and E^- with $\mathcal{A}_\Sigma \mathcal{A}_\Sigma^{-A} \{A(a)\}$.

We finally generalise Theorem 19 (ii) (part (i) is not interesting since simple queries do not generalise to any new class of \mathcal{ELI} -queries). Query classes such as $\mathcal{Q}[\diamond](\mathcal{EL})$ are defined in the obvious way by replacing in $\mathcal{Q}[\diamond]$ -queries conjunctions of atoms by \mathcal{EL} -queries.

Theorem 32. *The class $\mathcal{Q}_b[\diamond](\mathcal{EL}) \cap \mathcal{Q}_{\leq n}[\diamond](\mathcal{EL})$ is polynomially characterisable for any $n < \omega$.*

Again the positive and negative examples are obtained from the 1D case by replacing σ by \mathcal{A}_Σ and $\sigma \setminus \{A\}$ by appropriate split partners.

9 Applications to Learning

We apply our results on unique characterisability to exact learnability of temporal instance queries. Given a class \mathcal{Q} of such queries, we aim to identify a *target query* $\mathbf{q} \in \mathcal{Q}$ using queries to an oracle. The learner knows \mathcal{Q} and the signature σ (Σ in the 2D case) of \mathbf{q} . We allow only one type of queries, called *membership queries*, in which the learner picks a σ -data instance \mathcal{D} and asks the oracle whether $\mathcal{D} \models \mathbf{q}$ holds. (In the 2D case, the learner picks a pointed Σ -data instance (\mathcal{D}, a) and asks whether $\mathcal{D}, a, 0 \models \mathbf{q}$ holds.) The oracle answers 'yes' or 'no' truthfully. The class \mathcal{Q} is (*polynomial*

time) learnable with membership queries if there exists an algorithm that halts for any $\mathbf{q} \in \mathcal{Q}$ and computes (in polynomial time in the size of \mathbf{q} and σ/Σ), using membership queries, a query $\mathbf{q}' \in \mathcal{Q}$ that is equivalent to \mathbf{q} . By default, the learner does not know $|\mathbf{q}|$ in advance but reflecting Theorem 8 (iii), we also consider the case when $|\mathbf{q}|$ is known (which is common in active learning).

Obviously, unique characterisability is a necessary condition for learnability with membership queries. Conversely, if there is an algorithm that computes, for every $\mathbf{q} \in \mathcal{Q}$, an example set that uniquely characterises \mathbf{q} within $\mathcal{Q}^{sig(\mathbf{q})}$, then \mathcal{Q} is learnable with membership queries: enumerate $\mathcal{Q}^{sig(\mathbf{q})}$ starting with the smallest query \mathbf{q} , compute a characterising set E for \mathbf{q} and check using membership queries whether \mathbf{q} is equivalent to the target query. Eventually the algorithm will terminate with a query that is equivalent to the target query. As all our positive results on unique characterisability provide algorithms computing example sets, we directly obtain learnability with membership queries. Moreover, if the examples sets are computed in exponential time, then we obtain an exponential-time learning algorithm: in the enumeration above only $|sig(\mathbf{q})|^{|\mathbf{q}|}$ queries are checked before the target query is found. Unfortunately, we cannot infer polynomial-time learnability from polynomial characterisability in this way. In fact, monotone DNFs (disjunctions of conjunctions of propositional atoms) are a class of queries, which is polynomially characterisable (simply take the models corresponding to the disjuncts as positive examples and the models obtained from them by removing an atom as negative examples) but not polynomial-time learnable with membership queries (Angluin 1987b).

A detailed analysis of polynomial-time learnability using membership queries is beyond the scope of this paper. Instead, we focus on one main result, the polynomial-time learnability of $\mathcal{Q}_p[\circ, \diamond](\mathcal{ELI})$.

Theorem 33. (i) *The class of safe queries in $\mathcal{Q}_p[\circ, \diamond_r](\mathcal{ELI})$ is polynomial-time learnable with membership queries.*

(ii) *The class $\mathcal{Q}_p[\circ, \diamond_r](\mathcal{ELI})$ is polynomial-time learnable with membership queries if the learner knows the size of the target query in advance.*

(iii) *The class $\mathcal{Q}_p[\circ, \diamond](\mathcal{ELI})$ is polynomially-time learnable with membership queries.*

Proof. We consider the 1D case without \mathcal{ELI} -queries first.

(i) Our proof strategy is to construct a query \mathbf{q}' that agrees with \mathbf{q} on the positive and negative examples for \mathbf{q}' from Theorem 8. The algorithm proceeds by computing a data instance \mathcal{D} . Our aim is to arrive at \mathcal{D}_b through iterations of steps, from which the required query can be ‘read off’.

Step 1. First, identify the number of \circ and \diamond in \mathbf{q} by asking membership queries of the form σ^k incrementally, starting from $k = 1$, and then set $b = \min\{k \mid \sigma^k \models \mathbf{q}\} + 1$ and $\mathcal{D}_0 = \sigma^b$. Initialise $\mathcal{D} = \mathcal{D}_0$.

Step 2. Suppose that a data instance \mathcal{D}' is obtained from \mathcal{D} by applying one of the rules (a)–(e) of Theorem 8. If $\mathcal{D}' \models \mathbf{q}$ then replace \mathcal{D} with \mathcal{D}' . Repeat as long as possible. One can show that the number of applications

of each rule is bounded by a polynomial in $|\sigma|$ and the size of \mathbf{q} , and so **step 2** finishes in polynomial time.

Step 3. Suppose \mathcal{D} contains $\emptyset^b \rho_0^i \emptyset^b$ and $|\rho_0^i| \geq 2$. Since rule (a) does not apply, every homomorphism $h: \mathbf{q} \rightarrow \mathcal{D}$ sends some t_1, \dots, t_l to ρ_0^i , for $l \geq 1$. As \mathbf{q} does not contain lone conjuncts, \mathbf{q} contains singleton primitive blocks at positions t_1, \dots, t_l . Suppose $\rho_0^i = \{A_1, \dots, A_{|\rho_0^i|}\}$ and let $w = \{A_1\} \emptyset^b \{A_2\} \emptyset^b \dots \{A_{|\rho_0^i|}\} \emptyset^b$ (the order in which $A_1, \dots, A_{|\rho_0^i|}$, the elements of ρ_0^i , are enumerated does not matter, we fix any one). Let \mathcal{D}_k^i be obtained from \mathcal{D} by replacing $\emptyset^b \rho_0^i \emptyset^b$ with $\emptyset^b(w)^k$. Notice that, for $k = |\mathbf{q}|$, we have $\mathcal{D}_k^i \models \mathbf{q}$; however, the algorithm is not given this k . Instead, the algorithm incrementally iterates starting from $k = 1$ until $\mathcal{D}_k^i \models \mathbf{q}$. Since $k \leq |\mathbf{q}|$, this takes polynomially-many iterations. Let \mathcal{D}' be obtained from \mathcal{D}_k^i by removing primitive blocks as long as $\mathcal{D}' \models \mathbf{q}$. Notice that rules (a)–(e) do not apply to \mathcal{D}' . Replace \mathcal{D} with \mathcal{D}' . Repeat **step 3** as long as possible. Since no new lone conjuncts are introduced, the process finishes after polynomially-many steps.

Step 4. At this point of computation, the algorithm has identified all blocks of \mathbf{q} but not the sequences of \diamond and \diamond_r between them. They can be easily determined based on the positive and negative examples \mathcal{D}_i and \mathcal{D}_i^- .

The proof of (ii) is similar, with a modified **step 3**. Finally, (iii) is a consequence of (ii) as the size of the query \mathbf{q} does not exceed $n = |\sigma|b$.

We obtain a learning algorithm for $\mathcal{Q}_p[\circ, \diamond_r](\mathcal{ELI})$ by combining the learning algorithm above with the learning algorithm for \mathcal{ELI} -queries by ten Cate and Dalmau (2021) using the positive and negative examples given in Theorem 27. Note that the data instance \mathcal{A}_Σ is now used instead of σ and that one has to ‘unfold’ such non tree-shaped data instances into tree-shaped ones. \square

10 Conclusions

In this paper, we have considered temporal instance queries with *LTL* operators and started investigating their unique (polynomial) characterisability and exact learnability using membership queries. We have obtained both positive and negative results, depending on the available temporal operators and the allowed interaction between the temporal and object dimensions in queries. The results indicate that finding complete classifications of 1D and 2D temporal queries according to (polynomial) characterisability and learnability could be a very difficult task. In particular, interesting open problems include the polynomial characterisability of full $\mathcal{Q}_p^\sigma[\cup]$, more general criteria of polynomial characterisability for temporal branching queries and other temporal operators, the existence of polynomial split partners for \mathcal{ELI} -queries, and the polynomial-time learnability of $\mathcal{Q}_p^\sigma[\cup]$ and 2D extensions. From a conceptual viewpoint, it would be of interest to develop a framework that spells out explicitly the conditions that non-temporal queries should satisfy so that their combination with *LTL*-queries preserves polynomial characterisability and polynomial-time learnability.

Acknowledgments

This research was supported by the EPSRC UK grants EP/S032207 and EP/S032282 for the joint project ‘quant^{MD}: Ontology-Based Management for Many-Dimensional Quantitative Data’.

References

- Aarts, F., and Vaandrager, F. 2010. Learning i/o automata. In *International Conference on Concurrency Theory*, 71–85. Springer.
- Alexe, B.; ten Cate, B.; Kolaitis, P. G.; and Tan, W. C. 2011. Characterizing schema mappings via data examples. *ACM Trans. Database Syst.* 36(4):23.
- Angluin, D.; Frazier, M.; and Pitt, L. 1992. Learning conjunctions of Horn clauses. *Mach. Learn.* 9:147–164.
- Angluin, D. 1987a. Learning regular sets from queries and counterexamples. *Inf. Comput.* 75(2):87–106.
- Angluin, D. 1987b. Queries and concept learning. *Mach. Learn.* 2(4):319–342.
- Arenas, M., and Diaz, G. I. 2016. The exact complexity of the first-order logic definability problem. *ACM Trans. Database Syst.* 41(2):13:1–13:14.
- Artale, A.; Kontchakov, R.; Kovtunova, A.; Ryzhikov, V.; Wolter, F.; and Zakharyashev, M. 2017. Ontology-mediated query answering over temporal data: A survey (invited talk). In *Proc. of TIME 2017*, volume 90 of *LIPICs*, 1:1–1:37. Schloss Dagstuhl - Leibniz-Zentrum für Informatik.
- Baader, F.; Horrocks, I.; Lutz, C.; and Sattler, U. 2017. *An Introduction to Description Logics*. Cambridge University Press.
- Baader, F.; Borgwardt, S.; and Lippmann, M. 2015. Temporal query entailment in the description logic SHQ. *J. Web Semant.* 33:71–93.
- Barceló, P., and Romero, M. 2017. The complexity of reverse engineering problems for conjunctive queries. In *Proc. of ICDT*, 7:1–7:17.
- Borgwardt, S., and Thost, V. 2015. Temporal query answering in the description logic EL. In *Proc. of IJCAI 2015*, 2819–2825. AAAI Press.
- Camacho, A., and McIlraith, S. A. 2019. Learning interpretable models expressed in linear temporal logic. In *Proc. of ICAPS 2018*, 621–630. AAAI Press.
- Cassel, S.; Howar, F.; Jonsson, B.; and Steffen, B. 2016. Active learning for extended finite state machines. *Formal Aspects Comput.* 28(2):233–263.
- Chomicki, J., and Toman, D. 2018. *Temporal Logic in Database Query Languages*. New York, NY: Springer New York. 3992–3998.
- Fijalkow, N., and Lagarde, G. 2021. The complexity of learning linear temporal formulas from examples. *CoRR* abs/2102.00876.
- Foniok, J.; Nesetril, J.; and Tardif, C. 2008. Generalised dualities and maximal finite antichains in the homomorphism order of relational structures. *Eur. J. Comb.* 29(4):881–899.
- Funk, M.; Jung, J. C.; Lutz, C.; Pulcini, H.; and Wolter, F. 2019. Learning description logic concepts: When can positive and negative examples be separated? In *Proc. of IJCAI*, 1682–1688.
- Funk, M.; Jung, J. C.; and Lutz, C. 2021. Actively learning concepts and conjunctive queries under ELr-ontologies. In *Proc. of IJCAI 2021*, 1887–1893. ijcai.org.
- Gutiérrez-Basulto, V.; Jung, J. C.; and Sabellek, L. 2018. Reverse engineering queries in ontology-enriched systems: The case of expressive Horn description logic ontologies. In *Proc. of IJCAI-ECAI*.
- Hodkinson, I. M.; Wolter, F.; and Zakharyashev, M. 2000. Decidable fragment of first-order temporal logics. *Ann. Pure Appl. Log.* 106(1-3):85–134.
- Howar, F., and Steffen, B. 2018. Active automata learning in practice - an annotated bibliography of the years 2011 to 2016. In *Machine Learning for Dynamic Software Analysis: Potentials and Limits, International Dagstuhl Seminar 16172*, volume 11026 of *Lecture Notes in Computer Science*, 123–148. Springer.
- Jung, J. C.; Lutz, C.; Pulcini, H.; and Wolter, F. 2020. Logical separability of incomplete data under ontologies. In *Proc. of KR 2020*, 517–528.
- Kolaitis, P. G. 2011. Schema Mappings and Data Examples: Deriving Syntax from Semantics (Invited Talk). In *Proc. of FSTTCS 2011*, volume 13 of *Leibniz International Proceedings in Informatics (LIPIcs)*, 25–25. Dagstuhl, Germany: Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik.
- Konev, B.; Lutz, C.; Ozaki, A.; and Wolter, F. 2017. Exact learning of lightweight description logic ontologies. *J. Mach. Learn. Res.* 18:201:1–201:63.
- Konev, B.; Ozaki, A.; and Wolter, F. 2016. A model for learning description logic ontologies based on exact learning. In *Proc. of AAAI*, 1008–1015. AAAI Press.
- Lehmann, J., and Hitzler, P. 2010. Concept learning in description logics using refinement operators. *Machine Learning* 78:203–250.
- Lemieux, C.; Park, D.; and Beschastnikh, I. 2015. General ltl specification mining (t). In *Proc. of ASE*, 81–92. IEEE.
- Lutz, C.; Wolter, F.; and Zakharyashev, M. 2008. Temporal description logics: A survey. In *Proc. of TIME 2008*, 3–14. IEEE Computer Society.
- Martins, D. M. L. 2019. Reverse engineering database queries from examples: State-of-the-art, challenges, and research opportunities. *Inf. Syst.* 83:89–100.
- Neider, D., and Gavran, I. 2018. Learning linear temporal properties. In *Proc. of FMCAD 2018*, 1–10. IEEE.
- Nesetril, J., and Tardif, C. 2005. Short answers to exponentially long questions: Extremal aspects of homomorphism duality. *SIAM J. Discret. Math.* 19(4):914–920.
- Schild, K. 1993. Combining terminological logics with tense logic. In *Proc. of the 6th Portuguese Conf. on Progress in Artificial Intelligence, EPIA’93*, volume 727 of *Lecture Notes in Computer Science*, 105–120. Springer.

Shahbaz, M., and Groz, R. 2009. Inferring mealy machines. In *International Symposium on Formal Methods*, 207–222. Springer.

ten Cate, B., and Dalmau, V. 2021. Conjunctive queries: Unique characterizations and exact learnability. In *Proc. of ICDT 2021*, volume 186 of *LIPICs*, 9:1–9:24. Schloss Dagstuhl - Leibniz-Zentrum für Informatik.

ten Cate, B.; Kolaitis, P. G.; Qian, K.; and Tan, W. 2018. Active learning of GAV schema mappings. In *Proc. of PODS 2018*, 355–368. ACM.

ten Cate, B.; Dalmau, V.; and Kolaitis, P. G. 2013. Learning schema mappings. *ACM Trans. Database Syst.* 38(4):28:1–28:31.

A Proofs for Section 4

We prove the claims made in the introduction and Example 1.

(1) The query $\mathbf{q} = \diamond_r(A \wedge B)$ is not uniquely characterisable within $\mathcal{Q}_p[\diamond_r]$. Indeed, consider the queries $\mathbf{q}_1 = \diamond_r(A \wedge \diamond_r B)$ and $\mathbf{q}_i = \diamond_r(A \wedge \diamond_r(B \wedge \diamond_r \mathbf{q}_{i-1}))$. Clearly, $\mathbf{q} \models \mathbf{q}_i$ and $\mathbf{q}_i \not\models \mathbf{q}$ for all $i \geq 1$. Suppose \mathbf{q} fits (E^+, E^-) and n is the length of the longest example in E^- . Then \mathbf{q}_{n+1} also fits (E^+, E^-) as $\mathcal{D}, 0 \not\models \diamond_r(A \wedge B)$, and so $\mathcal{D}, 0 \not\models \mathbf{q}_{n+1}$, for any $\mathcal{D} \in E^-$.

(2) The query $\mathbf{q} = \perp \cup A$ (i.e., $\circ A$) is not uniquely characterisable within $\mathcal{Q}^{\{A\}}[\cup]$. For suppose \mathbf{q} fits (E^+, E^-) and n is the length of the longest example in E^- . Consider $\mathbf{q}' = (\circ^{n+1} A) \cup A$. Clearly, $\mathbf{q}' \not\models \mathbf{q}$ and $E \models \mathbf{q}'$ ($\mathcal{D} \not\models \mathbf{q}'$, for any $\mathcal{D} \in E^-$, because $\mathcal{D}, 1 \not\models A$ and $\mathcal{D}, 1 \not\models \circ^{n+1} A$).

(3) While the query $\diamond_r(A \wedge B)$ is not characterisable, there is a polynomial f such that for all $n \in \mathbb{N}$, it is characterisable within $\mathcal{Q}_p^n[\circ, \diamond_r]$ by some E_n of size $\leq f(n)$. Take $E_n = (E^+, E^-)$ with $E^+ = \{(\{A, B\}), (\emptyset, \{A, B\})\}$ and $E^- = \underbrace{\{(\{A\}, \{B\}), \dots, \{A\}, \{B\})\}}_{n \text{ times}}$. If $\mathbf{q}' \in \mathcal{Q}_p^n[\circ, \diamond_r]$

fits E then we can assume without loss of generality that it does not use \circ , as $(\{A, B\}) \models \mathbf{q}'$. This means that \mathbf{q}' is of the form $\rho_0 \wedge \diamond_r(\rho_1 \wedge \diamond_r(\rho_2 \wedge \dots \wedge (\diamond_r \rho_m) \dots))$, with $m < n$. Moreover, $\rho_0 = \emptyset$ as $(\emptyset, \{A, B\}) \models \mathbf{q}'$. And since $(\{A, B\}) \models \mathbf{q}'$, we must have $\rho_i \subseteq \{A, B\}$ for all i . Finally, as $\underbrace{\{(\{A\}, \{B\}), \dots, \{A\}, \{B\})\}}_{n \text{ times}} \not\models \mathbf{q}'$, there must be i

such that $\rho_i = \{A, B\}$. Thus, $\mathbf{q}' \equiv \diamond_r(A \wedge B)$.

B Proofs for Section 5

We show Theorem 8. We use the fact that $\mathcal{D} \models \mathbf{q}$ iff there is a homomorphism h from the set $\text{var}(\mathbf{q})$ of variables in \mathbf{q} to $[0, \max(\mathcal{D})]$, i.e., $h(t_0) = 0$, $A(h(t)) \in \mathcal{D}$ if $A(t) \in \mathbf{q}$, $h(t') = h(t) + 1$ if $\text{suc}(t, t') \in \mathbf{q}$, and $h(t) R h(t')$ if $R(t, t') \in \mathbf{q}$ for $R \in \{<, \leq\}$.

Theorem 8. (i) A query $\mathbf{q} \in \mathcal{Q}_p[\circ, \diamond_r]$ is uniquely characterisable within $\mathcal{Q}_p[\circ, \diamond_r]$ iff \mathbf{q} is safe.

(ii) Those queries that are uniquely characterisable within $\mathcal{Q}_p[\circ, \diamond_r]$ are actually polynomially characterisable within $\mathcal{Q}_p[\circ, \diamond_r]$.

(iii) The class $\mathcal{Q}_p[\circ, \diamond_r]$ is polynomially characterisable for bounded query size.

(iv) The class $\mathcal{Q}_p[\circ, \diamond]$ is polynomially characterisable.

Proof. (i) Assume that \mathbf{q} is given. To show (\Leftarrow) , we may assume \mathbf{q} in normal form (4) does not contain lone conjuncts. Let b be the number of \circ and \diamond in \mathbf{q} plus 1. Consider the example set $E = (E^+, E^-)$ constructed in the main paper. Then E is polynomial in $|\mathbf{q}|$ and $\mathcal{D} \models \mathbf{q}$ for all $\mathcal{D} \in E^+$. Using the condition that \mathbf{q} is in normal form, we show the following:

Claim 1. (i) There is only one homomorphism $h: \mathbf{q} \rightarrow \mathcal{D}_b$, and it maps isomorphically each \mathbf{q}_i onto $\bar{\mathbf{q}}_i$.

(ii) $\mathcal{D}_i^- \not\models \mathbf{q}$, for any \mathcal{R}_i different from \leq .

(iii) If \mathcal{D}'_b is obtained from \mathcal{D}_b by replacing some $\bar{\mathbf{q}}_i$ with $\bar{\mathbf{q}}'_i$ such that $\bar{\mathbf{q}}'_i, \ell \not\models \mathbf{q}_i$ for any $\ell \leq \max(\bar{\mathbf{q}}'_i)$, then $\mathcal{D}'_b \not\models \mathbf{q}$, and so $\mathcal{D} \not\models \mathbf{q}$, for all $\mathcal{D} \in E^-$.

Proof of claim. (i) As \mathbf{q} is in normal form and the gaps between $\bar{\mathbf{q}}_i$ and $\bar{\mathbf{q}}_{i+1}$ are not shorter than any block in \mathbf{q} , each block \mathbf{q}_i in \mathbf{q} is mapped by h to a single block $\bar{\mathbf{q}}_j$ of \mathcal{D}_b . The function $f: [0, n] \rightarrow [0, n]$ defined by taking $f(i) = j$ is such that $f(0) = 0$, $i < j$ implies $f(i) \leq f(j)$, and block \mathbf{q}_i is satisfied in $\bar{\mathbf{q}}_{f(i)}$. It also follows from the definition of normal form that if $f(i) = i$, then h isomorphically maps \mathbf{q}_i onto $\bar{\mathbf{q}}_i$ and $f(i-1) < i$ and $f(i+1) > i$. To show that $f(i) = i$ for all i , we first observe that $f(1) \geq 1$ and $f(j) = j$, for $j = \max\{i \mid f(i) \geq i\}$, from which $f(j-1) < j$ and $f(j+1) > j$. Then we can proceed in the same way inductively by considering f restricted to the smaller intervals $[j, n]$ and $[0, j]$.

(ii) Suppose \mathcal{R}_i is not \leq but there is a homomorphism $h: \mathbf{q} \rightarrow \mathcal{D}_i^-$. Consider the location of $h(s_0^i) = \ell$. Suppose ℓ is in $\bar{\mathbf{q}}_i$. Since $\rho_{k+1}^{i+1} \neq \emptyset$ and by the construction of \mathcal{D}_i^- , $h(s_0^{i+1})$ lies in some $\bar{\mathbf{q}}_j$ with $j > i+1$. But then there is a homomorphism $h': \mathbf{q} \rightarrow \mathcal{D}_b$ different from the one in (i), which is impossible. We arrive to the same contradiction if we assume that ℓ lies in $\bar{\mathbf{q}}_j$ with $j < i$ or $j > i$.

(iii) is proved analogously. \square

Thus, \mathbf{q} fits E . Suppose now \mathbf{q}' is any $\mathcal{Q}_p[\circ, \diamond, \diamond_r]$ -query in normal form. For $t \in \text{var}(\mathbf{q}')$, denote by τ_t the set of atoms A with $A(t) \in \mathbf{q}'$ and call it the *type* of t in \mathbf{q}' . Similarly, for $\ell \in [0, \max(\mathcal{D}_b)]$, denote by ρ_ℓ the set of atoms A with $A(\ell) \in \mathcal{D}_b$ and call it the *type* of ℓ in \mathcal{D}_b . A homomorphism $h: \mathbf{q}' \rightarrow \mathcal{D}_b$ is *block surjective* if every point in every block $\bar{\mathbf{q}}_i$ of \mathcal{D}_b is in the range $\text{ran}(h)$ of h ; it is *type surjective* if $\rho_\ell = \bigcup_{h(t)=\ell} \tau_t$ for all $\ell \in \text{ran}(h)$. The following claim follows immediately from the definitions:

Claim 2. (i) If there is a homomorphism $h: \mathbf{q}' \rightarrow \mathcal{D}_b$ that is not block or type surjective, then $\mathcal{D} \models \mathbf{q}'$ for some $\mathcal{D} \in E^-$ obtained from \mathcal{D}_b by (a).

(ii) If there exist a homomorphism $h: \mathbf{q}' \rightarrow \mathcal{D}_b$ and $(t < t') \in \mathbf{q}'$ or $(t \leq t') \in \mathbf{q}'$ such that $h(t) \neq h(t')$ and $h(t), h(t') \in \bar{\mathbf{q}}_i$, for some block $\bar{\mathbf{q}}_i$, then $\mathcal{D} \models \mathbf{q}'$ for some $\mathcal{D} \in E^-$ obtained from \mathcal{D}_b by (b).

Suppose now that $h: \mathbf{q}' \rightarrow \mathcal{D}_b$ is a block and type surjective homomorphism, $(t \leq t') \in \mathbf{q}'$ and $h(t) = h(t') = \ell$ lies in $\bar{\mathbf{q}}_i$. Then $h^{-1}(\ell) = \{t_1, \dots, t_k\}$ with $k \geq 2$ and

$(t_j \leq t_{j+1}) \in \mathbf{q}'$, $1 \leq j < k$. By (n3) and (n4), $\tau_{t_j} \neq \emptyset$ for at least one t_j , and so $\rho_\ell \neq \emptyset$. Consider possible locations of ℓ in $\bar{\mathbf{q}}_i$.

Case 1: ℓ has both a left and a right neighbour in $\bar{\mathbf{q}}_i$. Then there is $\mathcal{D} \in E^-$ obtained by (c)—i.e., by replacing the appropriate ρ_i^j with $\rho_i^j \emptyset^b \rho_i^j$ —and a homomorphism $h': \mathbf{q}' \rightarrow \mathcal{D}$, which ‘coincides’ with h except that $h'(t_1)$ is the point with the first ρ_i^j and $h'(t_j)$, for $j = 2, \dots, k$, is the point with the second ρ_i^j .

Case 2: ℓ has no neighbours in $\bar{\mathbf{q}}_i$, so this block is primitive and ρ_ℓ is a singleton (as \mathbf{q} has no lone conjuncts by our assumption). Then t_1 is the last variable in its block in \mathbf{q}' , t_k is the first variable in its block in \mathbf{q}' , and the t_i with $1 < i < k$, if any, are all primitive blocks. But then the types τ_{t_i} and $\tau_{t_{i+1}}$ are not comparable with respect to \subseteq , contrary to ρ_ℓ being a singleton. Thus, Case 2 cannot happen.

Case 3: ℓ has a left neighbour in $\bar{\mathbf{q}}_i$ but no right neighbour. As h is type surjective and in view of (n3), $\tau_{t_1} \subsetneq \rho_\ell$. Let $A \in \rho_\ell \setminus \tau_{t_1}$ and let $\mathcal{D} \in E^-$ be obtained by the first part of (d) by replacing $\rho_{k_i}^j$ with $\rho_{k_i}^j \setminus \{A\} \emptyset^b \rho_{k_i}^j$. Then there is a homomorphism $h': \mathbf{q}' \rightarrow \mathcal{D}$ that sends t_1 to the point with $\rho_{k_i}^j \setminus \{A\}$ and the remaining t_j to the point with $\rho_{k_i}^j$.

Case 4: ℓ has a right neighbour in $\bar{\mathbf{q}}_i$, $i \neq 0$, but no left neighbour. This case is dual to Case 3 and we use the second part of (d).

Case 5: $\ell = 0$. If $\bar{\mathbf{q}}_0$ is primitive, then all of the t_i are primitive blocks in \mathbf{q}' . By (n3), $\tau_{t_2} \not\subseteq \tau_{t_1}$; by type surjectivity, $\tau_{t_1} \subsetneq \rho_\ell$, and so there is $A \in \rho_\ell \setminus \tau_{t_1}$. By the first part of (e), we have $\mathcal{D} \in E^-$ obtained by replacing ρ_0^0 with $\rho_0^0 \setminus \{A\} \emptyset^b \rho_0^0$. Then there is a homomorphism $h': \mathbf{q}' \rightarrow \mathcal{D}$ that sends t_1 to the point with $\rho_0^0 \setminus \{A\}$ and the remaining t_j to the point with ρ_0^0 . Finally, if $\bar{\mathbf{q}}_0$ is not primitive, the second part of (e) gives $\mathcal{D} \in E^-$ by replacing ρ_0^0 in \mathcal{D}_b with $\rho_0^0 \emptyset^b \rho_0^0$. We obtain a homomorphism from \mathbf{q}' to \mathcal{D} by sending t_1 to the first ρ_0^0 and the remaining t_j to the second ρ_0^0 .

It remains to consider the case when there is a homomorphism $h: \mathbf{q}' \rightarrow \mathcal{D}_b$ that is an isomorphism between the blocks in \mathbf{q}' and the blocks in \mathcal{D}_b , and so the difference between \mathbf{q}' and \mathbf{q} can only be in the sequences of \diamond and \diamond_r between blocks. To be more precise, \mathbf{q} is of the form (4),

$$\mathbf{q}' = \mathbf{q}_0 \mathcal{R}'_1 \mathbf{q}_1 \dots \mathcal{R}'_n \mathbf{q}_n \quad (9)$$

and $\mathcal{R}_i \neq \mathcal{R}'_i$ for some i . Four cases are possible:

- (i) $\mathcal{R}_i = (r_0 \leq r_1)$ and $\mathcal{R}'_i = (s_0 < s_1) \dots (s_{l-1} < s_l)$, for $l \geq 1$. In this case, $\mathcal{D}_i \not\models \mathbf{q}'$, for $\mathcal{D}_i \in E^+$.
- (ii) $\mathcal{R}_i = (r_0 < r_1) \dots (r_{k-1} < r_k)$, $\mathcal{R}'_i = (s_0 < s_1) \dots (s_{l-1} < s_l)$, for $l > k$. Then again $\mathcal{D}_i \not\models \mathbf{q}'$.
- (iii) $\mathcal{R}_i = (r_0 < r_1) \dots (r_{k-1} < r_k)$, $\mathcal{R}'_i = (s_0 \leq s_1)$, for $k \geq 1$. In this case $\mathcal{D}_i^- \models \mathbf{q}'$, for $\mathcal{D}_i^- \in E^-$.
- (iv) $\mathcal{R}_i = (r_0 < r_1) \dots (r_{k-1} < r_k)$ and $\mathcal{R}'_i = (s_0 < s_1) \dots (s_{l-1} < s_l)$, for $l < k$. Then again $\mathcal{D}_i^- \models \mathbf{q}'$.

(\Rightarrow) Suppose \mathbf{q} in normal form (4) does contain a lone conjunct $\mathbf{q}_i = \rho$. Let ρ^- be the last type of the block \mathbf{q}_{i-1} and let ρ^+ be the first type of the block \mathbf{q}_{i+1} . Then ρ is a

disjoint union of some nonempty τ and τ' such that at least one of the queries \mathbf{s}'_1 or \mathbf{s}''_1 below is in normal form:

$$\begin{aligned} \mathbf{s}'_1 &= \mathbf{q}_0 \mathcal{R}_1 \dots \mathcal{R}_i \tau (\leq) \tau' \mathcal{R}_{i+1} \dots \mathcal{R}_n \mathbf{q}_n, \\ \mathbf{s}''_1 &= \mathbf{q}_0 \mathcal{R}_1 \dots \mathcal{R}_i \tau (\leq) \tau' (\leq) \tau \mathcal{R}_{i+1} \dots \mathcal{R}_n \mathbf{q}_n \end{aligned}$$

For example, if $\rho^- = \{A, A'\}$, $\rho = \{A, B\}$, $\rho^+ = \{A, B'\}$ and \mathcal{R}_i and \mathcal{R}_{i+1} are both \leq , we take $\tau = \{B\}$, $\tau' = \{A\}$, for which \mathbf{s}'_1 is not in normal form, while \mathbf{s}''_1 is. Pick one of \mathbf{s}'_1 and \mathbf{s}''_1 , which is in normal form, and denote it by \mathbf{s}_1 . For $n \geq 2$, let \mathbf{s}_n be the query obtained from \mathbf{s}_1 by duplicating n times the part $\tau(\leq)\tau'$ in \mathbf{s}_1 and inserting \leq between the copies. It is readily seen that \mathbf{s}_n is in normal form. Clearly, $\mathbf{q} \models \mathbf{s}_n$ and, similarly to the proof of Claim 1, one can show that $\mathbf{s}_n \not\models \mathbf{q}$, for any $n \geq 1$.

Now suppose $E = (E^+, E^-)$ characterises \mathbf{q} and let $n = \max\{\max(\mathcal{D}) \mid \mathcal{D} \in E^-\} + 1$. Then $E \models \mathbf{s}_n$, which is impossible. Indeed, consider any $\mathcal{D} \in E^-$. To show that $\mathcal{D} \not\models \mathbf{s}_n$, suppose otherwise. Then there is a homomorphism $h: \mathbf{s}_n \rightarrow \mathcal{D}$. By the pigeonhole principle, h maps some variables of types τ and τ' in \mathbf{s}_n to the same point in \mathcal{D} . But then h can be readily modified to obtain a homomorphism $h': \mathbf{q} \rightarrow \mathcal{D}$, contrary to $E^- \not\models \mathbf{q}$.

(ii) follows from the proof of (i) as (E^+, E^-) is of polynomial size.

(iii) We aim to characterize \mathbf{q} in normal form (4) which may contain lone conjuncts within the class of queries in $\mathcal{Q}_p[\diamond, \diamond]$ of size at most n , where n is the size of \mathbf{q} . The set E^+ of positive examples is defined as before and we extend the set of rules (a) to (e) in the definition of E^- as follows: if $\mathbf{q}_i = \rho(s)$ with $\rho = \{A_1, \dots, A_k\}$ is a block in \mathcal{D}_b with ρ a lone conjunct in \mathbf{q} , then

- (f) replace ρ with $(\rho \setminus \{A_1\} \emptyset^b \dots \emptyset^b \rho \setminus \{A_k\})^n$.

For the proof that (E^+, E^-) characterizes \mathbf{q} within the class of queries of size at most n , observe that with the exception of *Case 2* the proof of (i) still goes through. In *Case 2*, however, we can now apply the assumption that the size of \mathbf{q}' is bounded by n as we then obtain a data instance $\mathcal{D} \in E^-$ and a homomorphism $h': \mathbf{q}' \rightarrow \mathcal{D}$.

(iv) Assume \mathbf{q} in $\mathcal{Q}_p[\diamond, \diamond]$ is given. The proof of (i) shows that (E^+, E^-) , defined in the same way as in (i) except that the rules (c), (d), and (e) are not used to construct E^- , characterizes \mathbf{q} within $\mathcal{Q}_p[\diamond, \diamond]$ even if \mathbf{q} contains lone conjuncts. \square

C Proofs for Section 6

Theorem 9. Any $\mathcal{Q}_p^\sigma[\mathbf{U}]$ -queries $\mathbf{q} \not\equiv \mathbf{q}'$ can be separated by some \mathcal{D} with $\max(\mathcal{D}) \leq O((\min\{tdp(\mathbf{q}), tdp(\mathbf{q}')\})^2)$.

Proof. Let $\mathbf{q} = \rho_0 \lambda_1^* \rho_1 \dots \rho_n \emptyset^*$ and $\mathbf{q}' = \tau_0 \mu_1^* \tau_1 \dots \tau_k \emptyset^*$. If $n < k$, then $\rho_0 \rho_1 \dots \rho_n$ separates \mathbf{q} from \mathbf{q}' . Suppose $n = k$. If $\rho_i \subsetneq \tau_i$, for some $i \leq n$, then again $\rho_0 \rho_1 \dots \rho_n$ separates \mathbf{q} from \mathbf{q}' . So suppose $\rho_i = \tau_i$ for all $i \leq n$.

Let $\mathbf{q} \not\equiv \mathbf{q}'$. Then there is $\mathcal{D} = \rho_0 \lambda_1^{k_1} \rho_1 \dots \lambda_n^{k_n} \rho_n \emptyset^{k_{n+1}}$ separating \mathbf{q} from \mathbf{q}' . We show that $k_i \leq n+1$, for all $i \leq n$. To see this, we convert $\mathfrak{A}_{\mathbf{q}'}$ to a DFA $\mathfrak{B}_{\mathbf{q}'}$ using the subset construction and observe that whenever there are transitions

$Q_1 \rightarrow_{\lambda_i} \dots \rightarrow_{\lambda_i} Q_{n+1} \rightarrow_{\lambda_i} Q_{n+2}$ in $\mathfrak{B}_{q'}$ (with $Q_j \subseteq [0, n+1]$), then $Q_{n+1} = Q_{n+2}$ because, by the structure of $\mathfrak{A}_{q'}$, we have $\delta'_{\alpha^{n+1}} = \delta'_{\alpha^{n+2}}$, for any α , where δ'_w is the transition function on the states of $\mathfrak{B}_{q'}$ corresponding to the word w . \square

Lemma 12. *For any queries q and q' as above, either (i) each $\lambda_i \neq \perp$ subsumes μ_j occurring in some matching pair (λ_k, μ_j) or (ii) q and q' are separated by a data instance of the form \mathcal{D}_q^i or $\mathcal{D}_{q'}^j$. Also, if q is peerless, λ_i can only subsume μ_j in the matching pair (λ_i, μ_j) with $i \geq j$, in which case $\mu_j = \rho_j = \dots = \rho_{i-1} = \lambda_i$.*

Proof. If $\lambda_i \neq \perp$ does not subsume any μ_j , then, as we know, $\mathcal{D}_q^i \models q$ and $\mathcal{D}_q^i \not\models q'$. So suppose λ_i subsumes some μ_j . Then either (λ_i, μ_j) is a matching pair or $\mu_j \subsetneq \lambda_i$. Note that the latter is impossible if q is peerless. If μ_j does not subsume any λ_l , then $\mathcal{D}_{q'}^j \models q'$ and $\mathcal{D}_{q'}^j \not\models q$. Otherwise, we consider λ_l subsumed by μ_j , etc. Since $\lambda_l \subsetneq \lambda_i$, sooner or later this process will terminate. \square

Lemma 34. *Suppose that $q = \rho_0 \lambda_1^* \rho_1 \lambda_2^* \dots \lambda_n^* \rho_n \emptyset^*$ and $q' = \rho_0 \rho_1 \lambda_2^* \dots \lambda_n^* \rho_n \emptyset^*$ with $q \not\models q'$. Then there is a data instance \mathcal{D} of the form $\rho_0 \lambda_1 \rho_1 \lambda_2^{k_2} \dots \lambda_n^{k_n} \rho_n$ such that $\mathcal{D} \models q$ and $\mathcal{D} \not\models q'$.*

Proof. Take any $\mathcal{D} = \rho_0 \lambda_1^{k_1} \rho_1 \lambda_2^{k_2} \dots \lambda_n^{k_n} \rho_n$ with $\mathcal{D} \models q$ and $\mathcal{D} \not\models q'$. Let \mathcal{D}_i be \mathcal{D} with $k_1 = i$. Choose $\mathcal{D}_i \not\models q'$ with $\mathcal{D}_i \models q'$, for all $l < i$. Consider some $\mathcal{D}' = \rho_0 \rho_1 u$ with $\mathcal{D}' \in \mathcal{D}_{i-1}$ and set $\mathcal{D}'' = \rho_0 \lambda_1 \rho_1 u$. If $\mathcal{D}'' \models q'$, then $\mathcal{D}_i \models q'$ as $\mathcal{D}'' \in \mathcal{D}_i$, which is impossible. Thus, \mathcal{D}'' is the data instance we need. \square

Theorem 14. $\mathcal{P}^\sigma[\mathbb{U}]$ is polynomially characterisable within $\mathcal{Q}_p^\sigma[\mathbb{U}]$.

Proof. We show that any $q = \rho_0 \lambda_1^* \rho_1 \lambda_2^* \dots \lambda_n^* \rho_n \emptyset^*$ in $\mathcal{P}^\sigma[\mathbb{U}]$ is characterised by $E = (E^+, E^-)$, where E^+ contains all data instances of the following forms:

- (p₀) $\rho_0 \dots \rho_n$,
- (p₁) $\rho_0 \dots \rho_{i-1} \lambda_i \rho_i \dots \rho_n = \mathcal{D}_q^i$,
- (p₂) $\rho_0 \dots \rho_{i-1} \lambda_i^k \rho_i \dots \rho_{j-1} \lambda_j \rho_j \dots \rho_n = \mathcal{D}_{i,k}^j$, for $i < j$ and $k = 1, 2$;

and E^- has all instances that are not in $L(q)$ of the forms:

- (n₀) σ^n and $\sigma^{n-i} \sigma \setminus \{A\} \sigma^i$, for $A \in \rho_i$,
- (n₁) $\rho_0 \dots \rho_{i-1} \sigma \setminus \{A, B\} \rho_i \dots \rho_n$, for $A \in \lambda_i \cup \{\perp\}$, $B \in \rho_i \cup \{\perp\}$,
- (n₂) for all i and $A \in \lambda_i \cup \{\perp\}$, some data instance

$$\mathcal{D}_A^i = \rho_0 \dots \rho_{i-1} (\sigma \setminus \{A\}) \rho_i \lambda_{i+1}^{k_{i+1}} \dots \lambda_n^{k_n} \rho_n, \quad (10)$$

if any, such that $\max(\mathcal{D}_A^i) \leq (n+1)^2$ and $\mathcal{D}_A^i \not\models q^\dagger$ for q^\dagger obtained from q by replacing λ_j , for all $j \leq i$, with \perp . Note that $\mathcal{D}_A^i \not\models q$ for peerless q .

By definition, q fits E and $|E|$ is polynomial in $|q|$. Suppose $q' = \tau_0 \mu_1^* \tau_1 \dots \mu_m^* \tau_m \emptyset^*$ also fits E . By (p₀) and (n₀), we have $n = m$ and $\rho_i = \tau_i$, for $i \leq n$. Consider the maximal i with $\lambda_i \neq \mu_i$ (if there is no such, $q \equiv q'$).

Case 1: $\mu_i \neq \perp$, $\mu_i \not\subseteq \lambda_i$. By Lemma 12, if μ_i does not subsume any λ_j , then n_1 separates q and q' . So suppose μ_i subsumes λ_j . As q is peerless, $j > i$ and we have $\rho_i \subseteq \mu_i$. By the minimality of q' , there is \mathcal{D} such that $\mathcal{D} \models \rho_{i-1} \mu_i^* \rho_i \dots \mu_n^* \rho_n \emptyset^*$ but $\mathcal{D} \not\models \rho_{i-1} \rho_i \dots \mu_n^* \rho_n \emptyset^*$. By Lemma 34, we can choose $\mathcal{D} = \rho_{i-1} \mu_i \rho_i u$. Consider $\mathcal{D}' = \rho_{i-1} (\sigma \setminus \{A\}) \rho_i u$, where $A \in \lambda_i \setminus \mu_i$ if $\lambda_i \neq \perp$, or $A = \perp$ otherwise. Since $\rho_i \subseteq \mu_i \subseteq \sigma \setminus \{A\}$ and $\lambda_i \not\subseteq \sigma \setminus \{A\}$, we have $\rho_0 \dots \rho_{i-2} \mathcal{D}' \models \rho_0 \dots \rho_{i-1} (\sigma \setminus \{A\})^* \rho_i \lambda_{i+1}^* \dots \lambda_n^* \rho_n$ but $\mathcal{D}' \not\models \rho_0 \dots \rho_{i-1} \rho_i \lambda_{i+1}^* \dots \lambda_n^* \rho_n \emptyset^*$, for otherwise $\mathcal{D} \models \rho_{i-1} \rho_i \dots \mu_n^* \rho_n$ as $\lambda_j = \mu_j$ for all $j > i$. Therefore, there is $\mathcal{D}_A^i \in n_2$ with $\mathcal{D}_A^i \models q'$.

Case 2: $\lambda_i \neq \perp$, $\mu_i \subsetneq \lambda_i$. By Lemma 12, if μ_i does not subsume any λ_j , then n_1 separates q and q' . So suppose μ_i subsumes λ_j . As q is peerless, $j > i$ and we have $\rho_i \subseteq \mu_i$. But then $\rho_i \subseteq \lambda_i$, which is a contradiction.

Case 3: $\lambda_i \neq \perp$, $\mu_i = \perp$. Find the maximal $j < i$ such that λ_i subsumes μ_j . Then $\mu_j = \rho_j = \dots = \rho_{i-1} = \lambda_i$ by Lemma 12.

Suppose $\lambda_{j'} = \perp$ for all $j' \in [j+1, i-1]$. By Lemma 12, $\mu_{j'} = \perp$ for all $j' \in [j+1, i-1]$. For if $\mu_{j'} \neq \perp$, then it either does not subsume anything, or it subsumes λ_i . In the latter case, we have $\rho_{j'-1} \subsetneq \mu_{j'}$ and either there is $j'' \in (j, j')$ with $\rho_{j''} \not\subseteq \mu_{j''}$, in which case $\mu_{j''}$ does not subsume anything, or we violate the minimality condition. The queries $\rho_{j-1} \mu_j^* \rho_j \dots \rho_{i-1} \rho_i \dots$ and $\rho_{j-1} \rho_j \dots \rho_{i-1} \mu_j^* \rho_i \dots$ are obviously equivalent, so we can change this part of q' and look for the next place where the queries differ.

Now assume that there is $\lambda_{j'} \neq \perp$ with $j' \in [j+1, i-1]$. Suppose there is no data instance in (p₁), (p₂) or (n₁) separating q and q' . Then $\mu_{j'} = \lambda_{j'}$. If $j' > j+1$, consider $\mathcal{D}_{j',1}^i$. Suppose $\mathcal{D}_{j',1}^i \models q'$. Then there is $\mathcal{D}' = \rho_0 \dots \rho_{m-1} \mu_m \rho_m \dots \rho_{l-1} \mu_l \rho_l \dots \rho_n \in \mathcal{D}_{j',k}^i$. By analysing all possible orderings of i, j, m, l we see that it is only possible when $\mu_{j'-1} \subseteq \lambda_j$. Then we have $\lambda_{j'-1} = \mu_{j'-1}$. After that, we consider $\mathcal{D}_{j'-1,1}^i$ and so on. We now have $\mu_{j+1} \subseteq \mu_{j+2} \subseteq \dots \subseteq \lambda_{j'}$.

Consider $\mathcal{D}_{j+1,1}^i$. If $\mathcal{D}_{j+1,1}^i \models q'$, then, in view of $\mu_j \not\subseteq \lambda_{j+1}$, we have $\rho_{j-1} \subseteq \lambda_{j+1}$.

Considering $\mathcal{D}_{j+1,2}^i$, we see that either $\rho_{j-2} \subseteq \lambda_{j+1}$ and $\mu_{j-1} \subseteq \rho_j$ or $\rho_{j-2} \subseteq \rho_j$ and $\mu_{j-1} \subseteq \lambda_{j+1}$ with μ_{j-1} and ρ_{j-2} incomparable, since otherwise we will violate the peerlessness of q .

If $\rho_{j-2} \subseteq \lambda_{j+1}$ and $\mu_{j-1} \subseteq \rho_j$, then there is $j_1 < j-1$ such that $\mu_{j_1} \subseteq \rho_{j_1} \subseteq \dots \subseteq \rho_{j-1} \subseteq \lambda_{j+1}$. This μ_{j_1} is paired with some λ_l for $l > j$. By peerlessness, we have $l = j$ and $\mu_j = \rho_{j_1} = \dots = \rho_{j-1} = \lambda_j \subseteq \lambda_{j+1}$. It follows that $\lambda_{j-1} = \mu_{j-1}$ and we can repeat a previous argument and move further.

If $\rho_{j-2} \subseteq \rho_j$ and $\mu_{j-1} \subseteq \lambda_{j+1}$, then we have $\lambda_j = \rho_{j-1} = \mu_{j-1}$ and either $\mu_{j-2} \subseteq \rho_{j-1}$ or $\rho_{j-3} \subseteq \rho_{j-1}$.

Suppose $\mu_{j-2} \subseteq \rho_{j-1}$. Then $\rho_{j-3} \subseteq \rho_{j-2}$ and there is $j_2 \leq j-3$ such that $\mu_{j_2} \subseteq \rho_{j_2} \subseteq \dots \subseteq \rho_{j-2}$. We know that μ_{j_2} subsumes some λ_i ; by peerlessness it can only be λ_{j-1} with $\mu_{j_2} = \rho_{j_2} = \dots = \rho_{j-2} = \lambda_{j-1}$, and $\rho_{j-2} \not\subseteq \mu_{j-2}$. So we repeat a previous argument and move further.

Suppose $\rho_{j-3} \subseteq \rho_{j-1}$. In this case we have either $\mu_{j-3} \subseteq \rho_{j-2}$ or $\rho_{j-4} \subseteq \rho_{j-2}$ and we move further.

Either way we cannot stop moving further and the process cannot terminate. Therefore, there is a data instance from (p_1) , (p_2) or (n_1) separating \mathbf{q} and \mathbf{q}' .

This completes the proof of the theorem. \square

D Proofs for Section 7

We start by giving the proof of Lemma 16.

Lemma 16. *For every $\mathbf{q} \in \mathcal{Q}[\diamond]$ one can compute in polynomial time an equivalent query of the form $\mathbf{q}_1 \wedge \dots \wedge \mathbf{q}_n$ with $\mathbf{q}_i \in \mathcal{Q}_p[\diamond]$ for $i \leq n$.*

Proof. It is sufficient to observe that a query \mathbf{q} of the form

$$\rho_0 \wedge \diamond(\rho_1 \wedge \bigwedge_{i=1}^n \diamond \mathbf{q}_i)$$

with $\mathbf{q}_1, \dots, \mathbf{q}_n \in \mathcal{Q}[\diamond]$ is equivalent to

$$\mathbf{q}' = \rho_0 \wedge \bigwedge_{i=1}^n \diamond(\rho_1 \wedge \diamond \mathbf{q}_i)$$

To see this, observe that $\mathbf{q} \models \mathbf{q}'$ is trivial. For the converse direction consider any data instance \mathcal{D} with $\mathcal{D}, 0 \models \mathbf{q}'$. Then take the minimum ℓ_0 of all $\ell > 0$ such that $\mathcal{D}, \ell \models \rho_1 \wedge \diamond \mathbf{q}_i$. We have $\mathcal{D}, \ell_0 \models \rho_1 \wedge \diamond \mathbf{q}_i$ for all $i \leq n$ and so $\mathcal{D}, 0 \models \mathbf{q}$. \square

Details for Example 18. Recall that

$$\mathbf{q}_1 = \emptyset(s, \sigma)^n s, \quad \mathbf{q}_2 = \emptyset \sigma^{2n+1}$$

where

$$s = \{A_1, A_2\}\{B_1, B_2\}$$

and consider the set P of queries of the form

$$\emptyset s_1 \dots s_{n+1}$$

with s_i either $\{A_1\}\{A_2\}$ or $\{B_1\}\{B_2\}$. P contains 2^{n+1} queries.

Lemma 35. $\mathbf{q}_1 \wedge \mathbf{q}_2 \not\models \mathbf{q}$ for any $\mathbf{q} \in P$. For any data instance \mathcal{D} with $\mathcal{D} \models \mathbf{q}_1 \wedge \mathbf{q}_2$ there is at most one $\mathbf{q} \in P$ with $\mathcal{D} \models \mathbf{q}$.

Proof. We first construct for every $\mathbf{q} = \emptyset s_1 \dots s_{n+1} \in P$ a data instance \mathcal{D} with $\mathcal{D} \models \mathbf{q}_1 \wedge \mathbf{q}_2$ and $\mathcal{D} \models \mathbf{q}$. The data instance $\mathcal{D}_\mathbf{q}$ is defined by taking the data instance $\bar{\mathbf{q}}_1 = \rho_0, \dots, \rho_{3n+2}$ and replacing ρ_i by σ in it if, for some $j \leq n$:

- $i = 3j + 1$ and $s_{j+1} = \{A_1\}\{A_2\}$; or
- $i = 3j + 2$ and $s_{j+1} = \{B_1\}\{B_2\}$.

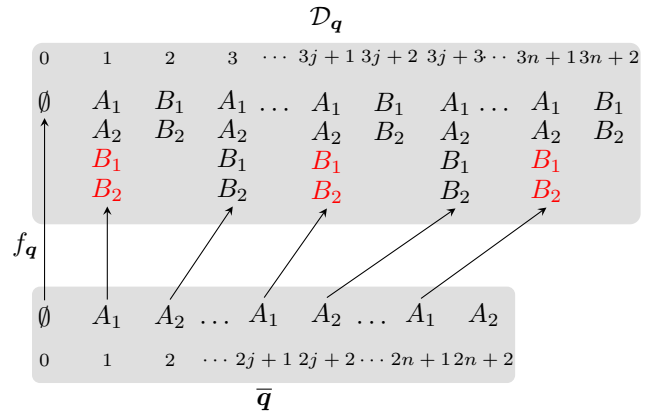


Figure 1: Definition of $\mathcal{D}_\mathbf{q}$ (additions to $\bar{\mathbf{q}}_1$ are indicated in red and $f_\mathbf{q}$ for $\mathbf{q} = \emptyset(\{A_1\}\{A_2\})^{n+1}$).

The data instance $\mathcal{D}_\mathbf{q}$ with $\mathbf{q} = \emptyset(\{A_1\}\{A_2\})^{n+1}$ is depicted in Figure 1. It also illustrates that $\mathcal{D}_\mathbf{q} \not\models \mathbf{q}$.

We have $\mathcal{D}_\mathbf{q} \models \mathbf{q}_1 \wedge \mathbf{q}_2$. We next show that $\mathcal{D}_\mathbf{q} \not\models \mathbf{q}$ for all $\mathbf{q} \in P$. We introduce a method to prove this that regards the query \mathbf{q} as a CQ. Namely, for $\mathbf{q} = \rho_0 \dots \rho_n$ and $\mathcal{D} = \delta_0 \dots \delta_m$ we have $\mathcal{D} \models \mathbf{q}$ iff the following inductive definition of a partial assignment f from $[0, n]$ to $[0, m]$ has domain $[0, n]$:

- $f(0) = 0$ if $\rho_0 \subseteq \delta_0$, otherwise $f(0)$ is undefined;
- if $f(i)$ is defined then let $f(i+1)$ be the minimal $j > f(i)$ such that $j \leq m$ and $\rho_{i+1} \subseteq \delta_j$, if such a j exists. Otherwise $f(i+1)$ is undefined.

If the domain of f does not equal $[0, n]$, then the definition either fails at $i_0 = 0$ or at some induction step $i_0 = i + 1$. In either case we say that *the definition of a satisfying assignment fails at i_0* .

We return to the proof that $\mathcal{D}_\mathbf{q} \not\models \mathbf{q}$. We show that the definition of an assignment fails at $2n+2$. In fact, the partial assignment $f_\mathbf{q}$ in $\mathcal{D}_\mathbf{q}$ is defined as follows

- $f_\mathbf{q}(0) = 0$,
- $f_\mathbf{q}(2i+1) = 3i+1$ if $s_{i+1} = \{A_1\}\{A_2\}$, for $0 \leq i \leq n$;
- $f_\mathbf{q}(2i+1) = 3i+2$ if $s_{i+1} = \{B_1\}\{B_2\}$, for $0 \leq i \leq n$;
- $f_\mathbf{q}(2i+2) = 3i+3$, for $0 \leq i < n$,

and $f_\mathbf{q}(2n+2)$ remains undefined.

Next we show that there does not exist a database \mathcal{D} with $\mathcal{D} \models \mathbf{q}_1 \wedge \mathbf{q}_2$ such that $\mathcal{D} \not\models \mathbf{q}$ and $\mathcal{D} \not\models \mathbf{q}'$ for distinct \mathbf{q}, \mathbf{q}' in P . Assume for a proof by contradiction that $\mathcal{D} = \delta_0 \dots \delta_m$ is such a data instance.

Let $\mathcal{D}_\mathbf{q} = \delta'_0 \dots \delta'_{3n+2}$ and $\mathcal{D}_{\mathbf{q}'} = \delta''_0 \dots \delta''_{3n+2}$ be the data instances defined above by extending $\rho_0 \dots \rho_{3n+2}$ with σ depending on \mathbf{q} and \mathbf{q}' , respectively.

As $\mathcal{D} \models \mathbf{q}_1$ there exists a bijective function g from some subset $A \subseteq [1, m]$ to $[1, 3n+2]$ satisfying the following

- $i < j$ iff $g(i) < g(j)$;
- $\delta_i \supseteq \rho_{g(i)}$ for all $i \in A$.

As $\mathcal{D} \models \mathbf{q}_2$ and $\mathbf{q} \neq \mathbf{q}'$,

- there exists $i \in [1, m] \setminus A$ such that $\delta_i = \sigma$; or
- there exists $i \in A$ such that $\delta_i = \sigma$ and $\delta'_{g(i)} \neq \sigma$; or
- there exists $i \in A$ such that $\delta_i = \sigma$ and $\delta''_{g(i)} \neq \sigma$.

We may thus assume w.l.o.g. that \mathcal{D} is obtained from $\rho_0 \dots \rho_{3n+2}$ by (1) inserting σ after some ρ_i with $i > 0$ or (2) by replacing some ρ_i with $i > 0$ and $\delta'_i \neq \sigma$ by σ or (3) by replacing some ρ_i with $i > 0$ and $\delta''_i \neq \sigma$ by σ .

Case 2 and Case 3 are equivalent, so we only only consider Case 1 and Case 2.

Case 1. One can show that then $\mathcal{D} \models q''$ for all $q'' = \emptyset s''_1 \dots s''_{n+1} \in P$. We do this for q .

Assume first that σ is inserted in $\rho_0 \dots \rho_m$ directly before timepoint 1. We define an assignment f_q^* by mapping 1 to the new node and (we assume that the original nodes are still numbered as before and not shifted):

- $f_q^*(0) = 0$,
- $f_q^*(2i+2) = 3i+1$ if $s_{i+1} = \{A_1\}\{A_2\}$, for $0 \leq i \leq n$;
- $f_q^*(2i+2) = 3i+2$ if $s_{i+1} = \{B_1\}\{B_2\}$, for $0 \leq i \leq n$;
- $f_q^*(2(i+1)+1) = 3i+3$, for $0 \leq i < n$.

Hence $\mathcal{D} \models q$. If σ is inserted later, the argument is similar. In this case the assignment is defined by first taking f_q as defined above and then from the point where σ is inserted f_q^* .

Case 2. This is proved similarly to Case 1. \square

We next provide further details of the proof of Theorem 19. We first provide the missing proof of Lemma 22.

Lemma 22. (i) For any $\mathcal{D} \in E_{q,m}^-$, we have $\mathcal{D} \not\models q$.

(ii) For any $q' \in \mathcal{Q}^\sigma[\diamond]$ with $q' \not\models q$ and $\text{tdp}(q') \leq m$, there exists $\mathcal{D} \in E_{q,m}^-$ with $\mathcal{D} \models q'$.

Proof. $\mathcal{D} \not\models q$ for all $\mathcal{D} \in E_{q,m}^-$ holds by definition and Lemma 21.

Now assume $q' \in \mathcal{Q}^\sigma[\diamond]$ with $q' \not\models q$ and $\text{tdp}(q') \leq m$. q' is equivalent to a conjunction $q'_1 \wedge \dots \wedge q'_k$ with $q'_i \in \mathcal{Q}_p[\diamond]$ and

$$q'_i = \tau_0^i \wedge \diamond(\tau_1^i \wedge \diamond(\tau_2^i \wedge \dots \wedge \diamond\tau_{\ell_i}^i))$$

If there exists $A \in \rho$ with A not in any τ_0^i , then $\sigma^A \sigma^m \models q'$, as required. Otherwise, $q' \not\models \bigwedge_{i=1}^k q'_i$. Pick an i with $q' \not\models q'_i$. Then $q'_j \not\models q'_i$ for all j . Hence $\mathcal{D}_{q'_i, m}^- \models q'_j$ for all j , and we obtain $\mathcal{D}_{q'_i, m}^- \models q'$. \square

We come to the construction of positive examples required for the proof of Part 1 of Theorem 19. For

$$q = A_0 \wedge \diamond(A_1 \wedge \diamond(A_2 \wedge \dots \wedge \diamond A_n))$$

we call any $q_{|k} := A_1, \dots, A_k$ with $k \leq n$ a *prefix* of q (note that A_0 is not taken into account).

Lemma 36. Let $w = A_1, \dots, A_k$ be a sequence of atomic concepts in a signature σ and assume $k < n$. Then one can construct in polynomial time a σ -data instance $\mathcal{D}_{w,n}$ such that for all simple $q \in \mathcal{Q}_p^\sigma[\diamond]$ with $\text{tdp}(q) = n$: $q_{|k} = w$ iff $\mathcal{D}_{w,n} \not\models q$.

Proof. Assume w and n are given. Then $\mathcal{D}_{w,n}$ defined as

$$\sigma(\sigma \setminus \{A_1\})^n \sigma(\sigma \setminus \{A_2\})^n \dots (\sigma \setminus \{A_k\})^n \sigma^{n-k-1}$$

is as required. \square

Assume $q = q_1 \wedge \dots \wedge q_m$ is balanced, simple, $n = \text{tdp}(q)$, and

$$q_i = A_0^i \wedge \diamond(A_1^i \wedge \diamond(A_2^i \wedge \dots \wedge \diamond A_n^i))$$

with $A_0 = A_0^1 = \dots = A_0^m$. Then let E^+ contain the σ -data instances

$$\{A_0\}\sigma^n, \quad \{A_0^1\} \dots \{A_n^1\} \dots \{A_0^m\} \dots \{A_n^m\},$$

and $\mathcal{D}_{w,n}$, for $w \in I$, where I is the set of all A_1^i, \dots, A_k^i, A such that $A \in \sigma$, $1 \leq i \leq m$, and A_1^i, \dots, A_k^i, A is not a prefix of any q_j . Clearly $\mathcal{D} \models q$ for all $\mathcal{D} \in E^+$.

Lemma 37. (E^+, E^-) characterises q within the class of balanced queries in $\mathcal{Q}_b[\diamond]$.

Proof. First observe that

$$\{A_0^1\} \dots \{A_n^1\} \dots \{A_0^m\} \dots \{A_n^m\} \not\models q'$$

for any q' that uses symbols not in σ or that is not simple. It remains to show that if $q \not\models q'$ and $q' \models q$, where $q' = q'_1 \wedge \dots \wedge q'_m$ is simple, balanced, and uses symbols in σ only, then there exists $\mathcal{D} \in E^+$ with $\mathcal{D} \not\models q'$. Recall that $n = \text{tdp}(q)$. If $\text{tdp}(q') > n$, then $\{A_0\}\sigma^n \not\models q'$ and we are done. As $q' \models q$, we then have $\text{tdp}(q') = n$. Then there exists q'_i with $q \not\models q'_i$. If q_i does not start with A_0 , then $\{A_0\}\sigma^n \not\models q'_i$ and we are done. Otherwise let k be maximal such that there exists q_j with $q_{j|k} = q'_{i|k}$. Then $k < n$ and $w = q'_{i|k+1} \in I$. Hence $\mathcal{D}_{w,n} \not\models q'_i$. \square

We finally show that the positive examples defined in the main paper are as required for Part 2 of Theorem 19. Assume $q = q_1 \wedge \dots \wedge q_m \in \mathcal{Q}_b^\sigma[\diamond] \cap \mathcal{Q}_{\leq n}^\sigma[\diamond]$ and that $q \not\models q'$ with $q' \in \mathcal{Q}_b[\diamond] \cap \mathcal{Q}_{\leq n}[\diamond]$. Assume $q' = q'_1 \wedge \dots \wedge q'_m$. If q' uses a symbol not in σ or $\text{tdp}(q') > N$, then $\rho\sigma^N \not\models q'$ and we are done. If the initial conjunct of some q'_i contains an $A \notin \rho$, then $\rho\sigma^N \not\models q'$, and we are done. If $\text{tdp}(q') < N$, then $q' \not\models q$, and so there exists $\mathcal{D} \in E^-$ with $\mathcal{D} \models q'$. Hence we may assume that $\text{tdp}(q') = N$. Take a j with $q \not\models q'_j$ and assume that

$$q'_j = \tau_0 \wedge \diamond(\tau_1 \wedge \diamond(\tau_2 \wedge \dots \wedge \diamond\tau_N))$$

Define $f: \{1, \dots, m\} \rightarrow \{1, \dots, N\}$ by taking for every i a ρ_j^i such that $\tau_j \not\subseteq \rho_j^i$ and setting $f(i) = j$. Such j exist since $q \not\models q'_j$ and since we excluded any other reason for non-entailment already. Then $\mathcal{D}_f \not\models q'_j$, as required.

E Proofs for Section 8

We begin by introducing some notation. A *pointed* atemporal Σ -data instance is a pair (\mathcal{A}, a) with \mathcal{A} an atemporal Σ -data instance and $a \in \text{ind}(\mathcal{A})$. We associate with \mathcal{A} the undirected graph

$$G_{\mathcal{A}} = (\text{ind}(\mathcal{A}), \bigcup_{P \in \sigma} \{\{a, b\} \mid P(a, b) \in \mathcal{A}\})$$

and call \mathcal{A} *acyclic* if $G_{\mathcal{A}}$ is acyclic and $P(a, b) \in \mathcal{A}$ implies $Q(a, b) \notin \mathcal{A}$ for any $Q \neq P$ and $Q(b, a) \notin \mathcal{A}$ for any Q . \mathcal{A} is *connected* if $G_{\mathcal{A}}$ is connected. We sometimes call acyclic and connected data instances *tree-shaped*. We assume the standard representation of an \mathcal{ELI} -query q as a set of atoms of the form $A(x), P(x, y)$ with x, y variables and a *distinguished variable*, also called the *answer variable* of q . For instance, $r = B \wedge \exists P. \exists P^- . A$ is represented as $\{B(x), P(x, y), P(z, y), A(z)\}$ with the distinguished variable x .

Every \mathcal{ELI} -query q defines a pointed data instance $\hat{q} = (q, x)$ (with $\text{ind}(q)$ being the variables), where q is tree shaped. Conversely, every pointed database (\mathcal{A}, a) with tree-shaped \mathcal{A} defines an \mathcal{ELI} -query.

Let \mathcal{A} and \mathcal{B} be data instances. We call a mapping h from $\text{ind}(\mathcal{A})$ to $\text{ind}(\mathcal{B})$ a *homomorphism* from \mathcal{A} to \mathcal{B} , in symbols $h : \mathcal{A} \rightarrow \mathcal{B}$, if

- $A(a) \in \mathcal{A}$ implies $A(h(a)) \in \mathcal{B}$;
- $P(a, b) \in \mathcal{A}$ implies $P(h(a), h(b)) \in \mathcal{B}$.

We call h a *homomorphism* from pointed (\mathcal{A}, a) to pointed (\mathcal{B}, b) if it is a homomorphism from \mathcal{A} to \mathcal{B} and $h(a) = b$. We write $(\mathcal{A}, a) \rightarrow (\mathcal{B}, b)$ if there exists a homomorphism from (\mathcal{A}, a) to (\mathcal{B}, b) .

Lemma 38. *For all \mathcal{ELI} -queries q_1 and q_2 , we have $q_1 \models q_2$ iff $\hat{q}_2 \rightarrow \hat{q}_1$.*

We call a pointed data instance (\mathcal{A}, a) *core* if every homomorphism $h : (\mathcal{A}, a) \rightarrow (\mathcal{A}, a)$ is an isomorphism. Pointed structures (\mathcal{A}, a) and (\mathcal{B}, b) are *homomorphically equivalent* if $(\mathcal{A}, a) \rightarrow (\mathcal{B}, b)$ and $(\mathcal{B}, b) \rightarrow (\mathcal{A}, a)$.

Theorem 39. *For every tree-shaped (\mathcal{A}, a) , one can construct in polynomial time a core that is tree-shaped and homomorphically equivalent to (\mathcal{A}, a) .*

We have defined frontiers within the set of \mathcal{ELI} -queries partially ordered by entailment. It is sometimes more convenient to define frontiers on the class of tree-shaped data instance partially ordered by homomorphisms. A set \mathcal{F} of tree-shaped (\mathcal{A}', a') is called a *frontier* of tree-shaped (\mathcal{A}, a) if

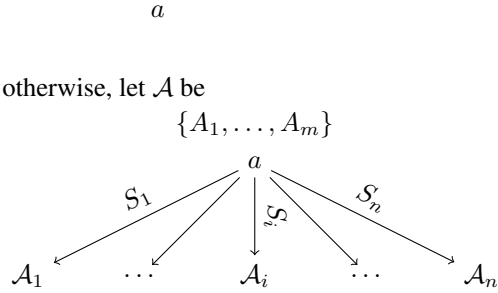
- $(\mathcal{A}', a') \rightarrow (\mathcal{A}, a)$ for all $(\mathcal{A}', a') \in \mathcal{F}$;
- $(\mathcal{A}, a) \not\rightarrow (\mathcal{A}', a')$ for all $(\mathcal{A}', a') \in \mathcal{F}$;
- if $(\mathcal{B}, b) \rightarrow (\mathcal{A}, a)$, then either $(\mathcal{A}, a) \rightarrow (\mathcal{B}, b)$ or there exists $(\mathcal{A}', a') \in \mathcal{F}$ with $(\mathcal{B}, b) \rightarrow (\mathcal{A}', a')$.

The frontier of (\mathcal{A}, a) is denoted by $\mathcal{F}(\mathcal{A}, a)$.

E.1 Proof of Theorem 23

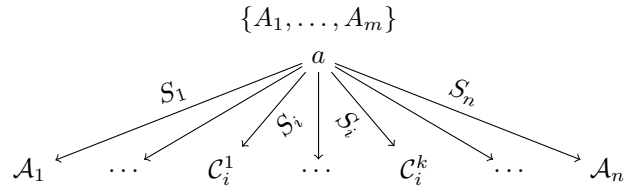
We describe the construction of the frontier $\mathcal{F}(\mathcal{A}, a)$ for a tree-shaped pointed data instance \mathcal{A}, a . (This construction is a minor adaptation of the construction in (ten Cate and Dalmau 2021), where the proof of its correctness is given.) In this case, we can view (\mathcal{A}, a) as a labelled directed tree with the set of nodes $\text{ind}(\mathcal{A})$ labelled with (possibly empty) sets of concept names $\{A_1, \dots, A_n\}$, edges (b, c) labelled with either P or P^- , for a role name $P \in \Sigma$, and rooted at a . First, we construct *pre-frontier*, which is a set of structures $\mathcal{P}(\mathcal{A}, a)$ defined inductively as follows:

- if a is a leave of \mathcal{A} and a has no concept name labels, then $\mathcal{P}(\mathcal{A}, a) = \emptyset$;
- if a is a leave of \mathcal{A} labelled with $\{A_1, \dots, A_m\}$, then $\mathcal{P}(\mathcal{A}, a) = \{(B^1, a), \dots, (B^m, a)\}$, where B^i is $\{A_1, \dots, A_{i-1}, A_{i+1}, A_m\}$

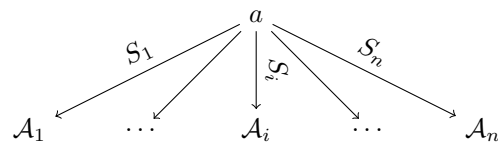


otherwise, let \mathcal{A} be

where each \mathcal{A}_i has a root a_i and let $\mathcal{P}(\mathcal{A}_i, a_i) = \{(B_i^1, a_i), \dots, (B_i^k, a_i)\}$, for $k \geq 0$. We construct another set $\{(C_i^j, c_i^j)\}$ of pointed structures, where the domain of C_i^j is equal to the set of sequences bj where b is in the domain of B_i^j , and $S(bj, cj)$ is in C_i^j iff $S(b, c)$ is in B_i^j (and similarly for concept assertions $A(b)$). We set $c_i^j = a_i.j$. We note that every element d in the domain of B_i^j or C_i^j has the form $a.j_1 \dots j_\ell$, for $\ell \geq 0$ and a from the domain of \mathcal{A} . We write $\text{orig}(d) = a$ to indicate that a is the *original* of d . Then, $\mathcal{P}(\mathcal{A}, a)$ is defined as $\{(\mathcal{D}_1, a), \dots, (\mathcal{D}_n, a)\} \cup \{(\mathcal{E}_1, a), \dots, (\mathcal{E}_m, a)\}$, where \mathcal{D}_i is the structure:



and each subtree C_i^j is rooted at c_i^j , and \mathcal{E}_i is $\{A_1, \dots, A_{i-1}, A_{i+1}, \dots, A_m\}$



Finally, the frontier $\mathcal{F}(\mathcal{A}, a)$ is $\{(\mathcal{F}_1, a), \dots, (\mathcal{F}_n, a)\}$, where each (\mathcal{F}_i, a) is obtained from $(\mathcal{D}_i, a) \in \mathcal{F}(\mathcal{A}, a)$ in the following way. Let the domain of \mathcal{D}_i be $\{a, d_1, \dots, d_k\}$. We construct a j -th copy \mathcal{A}^j of \mathcal{A} , where domain of the structure \mathcal{A}^j is the set of $b - j$, for b in the domain of \mathcal{A} and $S(b - j, c - j)$ is in \mathcal{A}^j iff $S(b, c)$ is in \mathcal{A} (and similarly for concept assertions $A(b)$; note that we use negative numbers). Then, \mathcal{F}_i is the union of $\mathcal{D}_i, \mathcal{A}^1, \dots, \mathcal{A}^k$ (note that they are all pairwise disjoint), where we add $S(b - j, d_j)$ to \mathcal{F}_i iff $S(b, c)$ is in \mathcal{A} and $\text{orig}(d_j) = c$.

E.2 Proof of Theorem 25 (i)

Let $\text{tdp}(q)$ be the maximum number of nested temporal operators in q (assuming that no subquery of q starting with

a temporal operator is equivalent to \top). Let $rdp(q)$ be the length n of the longest sequence $\exists S_1 \dots \exists S_n$ such that $\exists S_{i+1}$ is in the scope of $\exists S_i$ but not in the scope of a temporal operator in the scope of $\exists S_i$. The set of $\mathcal{Q}[\circ, \diamond] \otimes \mathcal{EL}$ -queries q with $tdp(q) \leq d$ and $rdp(q) \leq r$, for some fixed $d, r < \omega$, contains finitely-many, say N_{dr} -many, non-equivalent queries. It is easy to construct data instances \mathcal{D}^d and \mathcal{D}^{dr} such that

- $\mathcal{D}^d, a, 0 \models q$ iff $tdp(q) \leq d$, for all q ;
- $\mathcal{D}^{dr}, a, 0 \models q$ iff $rdp(q) \leq r$, for all q with $tdp(q) = d$.

It is also not hard to show that any pair of nonequivalent queries $q, q' \in \mathcal{Q}[\circ, \diamond] \otimes \mathcal{EL}$ with $tdp(q) = tdp(q')$ and $rdp(q) = rdp(q')$ is distinguished by a data instance \mathcal{D} with $\max(\mathcal{D}) \leq N_{dr}$ and $|ind(\mathcal{D})| \leq N_{dr}$.

Now, given a q with $tdp(q) = d$ and $rdp(q) = r$, we construct an example set $(E, a, 0)$ with

$$E^+ = \{\mathcal{D} \mid \mathcal{D}, a, 0 \models q, \max(\mathcal{D}) \leq N_{dr}, |ind(\mathcal{D})| \leq N_{dr}\},$$

$$E^- = \{\mathcal{D} \mid \mathcal{D}, a, 0 \not\models q, \max(\mathcal{D}) \leq N_{dr}, |ind(\mathcal{D})| \leq N_{dr}\}.$$

It follows from the observations above that this example set uniquely characterises q .

E.3 Proof of Theorem 25 (ii)

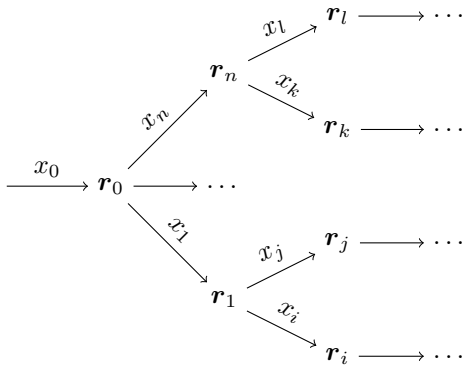
We observe that the queries $q \in \mathcal{Q}[\circ, \diamond] \otimes \mathcal{EL}$ can be viewed as the formulas $q(x)$ of the monodic fragment of temporal FO (Hodkinson, Wolter, and Zakharyashev 2000) (see also (Schild 1993)). In particular, the queries $q \in \mathcal{Q}_p[\diamond] \otimes \mathcal{EL}$ naturally correspond to the formulas of the form:

$$q_0(x_0) = \exists x_1, \dots, x_s$$

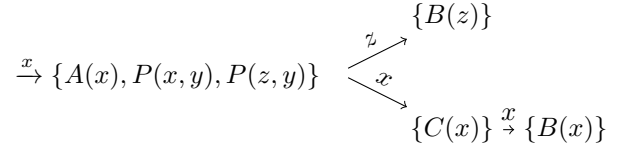
$$\left(\bigwedge_{j \in I_1} A_j(x_{i_j}) \wedge \bigwedge_{k, \ell \in I_2} P_{k, \ell}(x_{i_k}, x_{i_\ell}) \wedge \bigwedge_{m \in I_3} \diamond q_m(x_{i_m}) \right),$$

where $s \geq 0$, $0 \leq i_j, i_k, i_\ell, i_m \leq s$, q_m are of the same form as q_0 , and $x_{i_m} \neq x_{i_{m'}}$ for all $m, m' \in I_3$ with $m \neq m'$ (i.e., we do not allow more than one \diamond -subformula for with the same free variable). Moreover, the set of atoms $r_0 = \{A_j(x_{i_j} \mid j \in I_1)\} \cup \{P_{k, \ell}(x_{i_k}, x_{i_\ell}) \mid k, \ell \in I_2\}$ is an \mathcal{EL} -query with a distinguished variable x_0 (i.e., \hat{r} it is tree-shaped pointed data instance).

Thus, the $q \in \mathcal{Q}_p[\diamond] \otimes \mathcal{EL}$ queries q can be represented as trees:



where the nodes are labelled with \mathcal{EL} -queries r_i and the edges between the nodes are labelled with (query) variables, constructed as follows. The nodes of this tree, for a query q_0 , are from the set $\{q_0\} \cup \{q_i \mid \diamond q_i \text{ is a subformula of } q_0\}$. The label of each q_i is the \mathcal{EL} -query r_i constructed from the nontemporal atoms of q_i ; see r_0 for q_0 above. (We always omit the nodes on pictures as above; they are clear from the context.) The set of edges is defined as $\{(q_i, q_j) \mid \diamond q_j \text{ occurs in } q_i \text{ not in scope of a } \diamond\}$. The label of the edge (q_i, q_j) is the (unique) free variable of q_j . By the construction of q_i , we note that there are no outgoing edges of a node with the same label. For example, the query $q = \diamond(C \wedge \diamond B) \wedge A \wedge \exists P. \exists P^- . \diamond B$ is represented as follows:



For every $q \in \mathcal{Q}_p[\diamond] \otimes \mathcal{EL}$ (in what follows we do not distinguish between q and its tree representation) and every node r_i in it, we define $dep(r_i)$ to be equal to the length of the path between the root r and r_i . With such a q , we associate the data instance $\mathcal{D}^q = \mathcal{D}_0 \dots \mathcal{D}_n$ with $n = \max\{dep(r_i) \mid r_i \text{ in } q\}$, $ind(\mathcal{D}^q)$ equal to the set of variables of q , and $\mathcal{D}_i = \{A(x) \mid A(x) \in r_i \text{ for some } r_i \text{ in } q\} \cup \{P(x, y) \mid P(x, y) \in r_i \text{ for some } r_i \text{ in } q\}$. We also set $a^q = x$, for the distinguished variable x of q . On the other hand, we associate with q the atemporal pointed database $(\mathcal{A}^q, (x, 0))$, where x is the distinguished variable of q , with $ind(\mathcal{A}^q) = \{(y, m) \mid y \text{ occurs in some } r_i \text{ in } q \text{ and } dep(r_i) = m\}$, $A((y, m)) \in \mathcal{A}^q$ iff $A(y) \in r_i$ and $dep(r_i) = m$, and $P((y, m), (z, m)) \in \mathcal{A}^q$ iff $P(y, z) \in r_i$ and $dep(r_i) = m$. Moreover, for a fresh role name T , we add $T((y, m), (y, n))$ to \mathcal{A}^q for all $(y, m), (y, n) \in ind(\mathcal{A}^q)$, such that $n = m + 1$. We observe that \mathcal{A}^q is tree-shaped. Also, every connected (taking into account T -atoms) substructure of \mathcal{A}^q induces a query $q' \in \mathcal{Q}_p[\diamond] \otimes \mathcal{EL}$. Denote by $cl(\mathcal{A}^q)$ the structure that extends \mathcal{A}^q by adding the T -transitive closure.

Lemma 40. *Let $q, q' \in \mathcal{Q}_p[\diamond] \otimes \mathcal{EL}$ and assume that q (respectively, q') has a distinguished variable x (y). Then $q \models q'$ iff $\mathcal{D}^q, a^q \models q'$ iff $(cl(\mathcal{A}^q), x) \rightarrow (cl(\mathcal{A}^q), y)$.*

The following result shows the important property of the core of \mathcal{A}^q (note that (i) does not follow from Theorem 23):

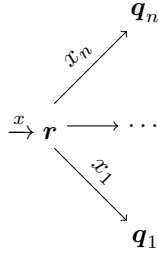
Theorem 41. *Let $q \in \mathcal{Q}_p[\diamond] \otimes \mathcal{EL}$. Then*

- (i) a pointed substructure $(\mathcal{C}, (x, 0))$ of $(cl(\mathcal{A}^q), (x, 0))$ that is a core is computable in polynomial time;
- (ii) there is a query $q' \in \mathcal{Q}_p[\diamond] \otimes \mathcal{EL}$ such that $cl(\mathcal{A}^q) = \mathcal{C}$.

Frontier for 2D queries. We now define a *frontier* $\mathcal{F}(q, r, x)$ for a node r with a distinguished variable x in a query $q \in \mathcal{Q}_p[\diamond] \otimes \mathcal{EL}$ by induction on the construction of (the tree of) q :

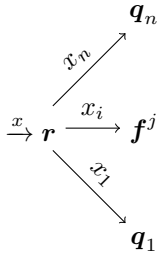
- if q is $\overset{x}{\rightarrow} r$ (i.e., an \mathcal{ELI} -query r with a distinguished variable x), then we set $\mathcal{F}(q, r, x)$ equal to $\mathcal{F}(r, x)$, which takes the form $\{\overset{x}{\rightarrow} r^1, \dots, \overset{x}{\rightarrow} r^k\}$.

- otherwise, let q be

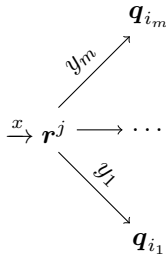


In this case, we set $\mathcal{F}(q, r, x)$ equal to $Q_1 \cup \dots \cup Q_n \cup Q'$, where

- assuming $\mathcal{F}(q_i, r_i, x_i) = \{\overset{x_i}{\rightarrow} f^1, \dots, \overset{x_i}{\rightarrow} f^\ell\}$, we set $Q_i = \{\overset{x}{\rightarrow} q_i^1, \dots, \overset{x}{\rightarrow} q_i^\ell\}$ and q_i^j :

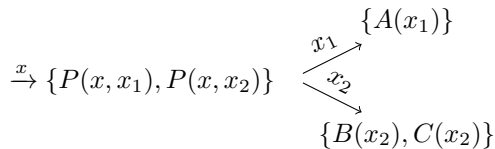


- assuming $\mathcal{F}(r, x) = \{\overset{x}{\rightarrow} r^1, \dots, \overset{x}{\rightarrow} r^k\}$, we set $Q' = \{\overset{x}{\rightarrow} q^1, \dots, \overset{x}{\rightarrow} q^k\}$ and q^j :



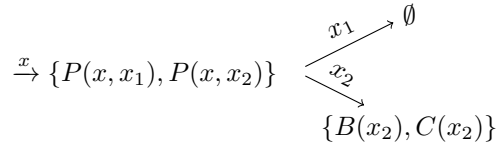
where $\{y_1, \dots, y_m\}$ is the domain of r^j and $orig(y_\ell) = x_{i_\ell}$.

Example 42. Let $q = \exists P.\diamond A \wedge \exists P.\diamond(B \wedge C)$, i.e.,

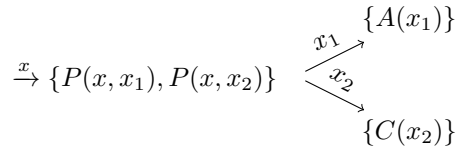
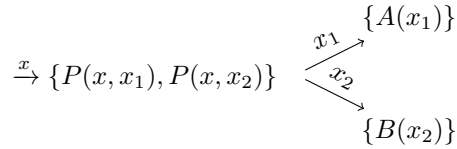


Clearly, q is core. Then, $\mathcal{F}(q, \{A(x_1)\}, x_1) = \{(\emptyset, x_1)\}$ and $\mathcal{F}(q, \{B(x_2), C(x_2)\}, x_2) = \{(\{B(x_2)\}, x_2), (\{C(x_2)\}, x_2)\}$. Therefore, $\mathcal{F}(q, \{P(x, x_1), P(x, x_2)\}, x) = Q_1 \cup Q_2 \cup Q'$, where:

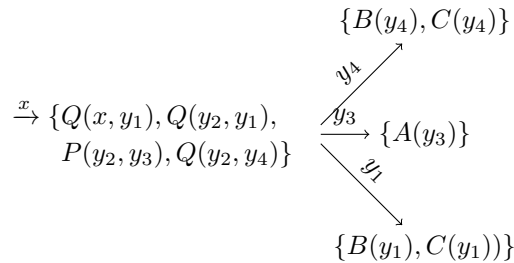
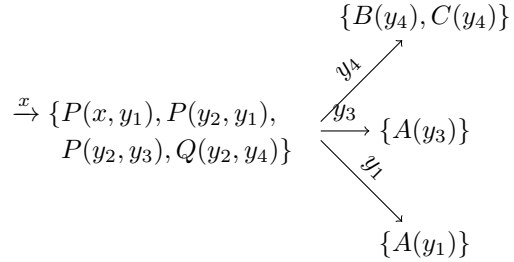
- Q_1 contains a single structure:



- Q_2 contains two structures:



- Q' contains two structures:



We observe the following important property, which can be proved by induction on the construction of (the tree of) q using Lemma 40 and Theorem 41:

Theorem 43. Let $q \in \mathcal{Q}_p[\diamond] \otimes \mathcal{ELI}$ be such that $cl(\mathcal{A}^q)$ is a core and r be the root of q with a distinguished variable x . Then (i) $\mathcal{F}(q, r, x) \subseteq \mathcal{Q}_p[\diamond] \otimes \mathcal{ELI}$; (ii) $\mathcal{F}(q, r, x)$ is polynomial in the size of q ; (iii) the following properties hold:

- $q \models q'$ for every $q' \in \mathcal{F}(q, r, x)$;
- $q' \not\models q$ for every $q' \in \mathcal{F}(q, r, x)$;
- if $q \models q''$, then either $q'' \models q$ or there exists $q' \in \mathcal{F}(q, r, x)$ such that $q' \models q''$, for any $q'' \in \mathcal{Q}_p[\diamond] \otimes \mathcal{ELI}$.

Constructing examples. Given a $q \in \mathcal{Q}_p[\diamond] \otimes \mathcal{ELI}$, we take q' from Theorem 41 and set $E_q^+ = \{\mathcal{D}^{q'}, a^{q'}\}$. We take $E_q^- = \{\mathcal{D}^{q''}, a^{q''} \mid q'' \in \mathcal{F}(q', r, x')\}$, where r is the root of q' with distinguished variable x' . Finally, using

Lemma 40 and Theorems 41 and 43, we obtain the following:

Theorem 44. (E_q^+, E_q^-) uniquely characterises q , for any $q \in \mathcal{Q}_p[\diamond] \otimes \mathcal{ELI}$.

For $q \in \mathcal{Q}_p[\circ] \otimes \mathcal{ELI}$ the construction of $\mathcal{F}(q, r)$ is even easier. Indeed, it is sufficient to take $cl(\mathcal{A}^q) = \mathcal{A}^q$. This completes the proof of Theorem 25 (ii).

E.4 Proof of Theorem 26

For $n \geq 1$, $1 \leq i \leq n$, $j \leq i$, define $\mathcal{Q}_p[\mathcal{EL}/\circ, \diamond]$ -queries $q_{i,j}^n$ recursively by taking:

$$\begin{aligned} q_{ii}^n &= r_i^n, & q_{i,j-1}^n &= B \wedge \circ (B \wedge \diamond (A \wedge \circ q_{i,j}^n)), \\ r_n^n &= \diamond (B \wedge \diamond A), & r_l^n &= \diamond (B \wedge \diamond (A \wedge \circ s_{n-l})), \quad l < n, \\ s_1 &= B \wedge \circ (B \wedge \diamond A), & s_{i+1} &= B \wedge \circ (B \wedge \diamond (A \wedge \circ q_i)) \end{aligned}$$

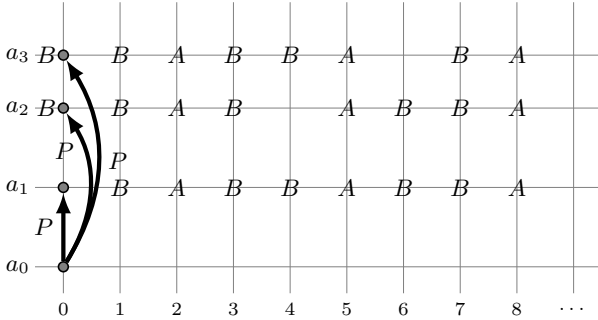
and set $q_i^n = q_{i,1}^n$. Minimal models of q_i^n look as follows:

$$\underbrace{BB\emptyset^*A}_{1} \dots \underbrace{BB\emptyset^*A}_{i-1} \underbrace{\emptyset^*B\emptyset^*A}_{i} \underbrace{BB\emptyset^*A}_{i+1} \dots \underbrace{BB\emptyset^*A}_{n}$$

Consider the $\mathcal{Q}_p[\mathcal{EL}/\circ, \diamond]$ -query $q = \exists P.q_1^n \wedge \dots \wedge \exists P.q_n^n$. We claim that any unique characterisation of q contains at least 2^n positive examples (in E^+). Indeed, let \mathcal{Q} be the set of all queries of the form $q \wedge \exists P.s$ with

$$s = o_1(B \wedge \diamond (A \wedge o_2(B \wedge \diamond (A \wedge \dots \wedge o_n(B \wedge \diamond A) \dots))))),$$

where each o_i is either \circ or $\diamond\circ$ if $i > 1$, and either blank or \diamond if $i = 1$. Observe that $|\mathcal{Q}| = 2^n$ and $q' \models q$ for each $q' \in \mathcal{Q}$. On the other hand, $q \not\models q'$. Indeed, let I_0 be the set of indices i such that $o_i = \circ$ and $i > 1$ or o_i is blank and $i = 1$. Let $I_1 = [1, n] \setminus I_0$. Define a data instance \mathcal{D}_{I_0} with $ind(\mathcal{D}_{I_0}) = \{a_0, \dots, a_n\}$ and $\max(\mathcal{D}_{I_0}) = 3n - 1$ as follows. To begin with, \mathcal{D}_{I_0} contains $P(a_0, a_1, 0), \dots, P(a_0, a_n, 0)$. For $i \in [1, n]$ and $k \in [0, 3n - 1]$, \mathcal{D}_{I_0} contains $A(a_i, k)$ if $k \equiv 2 \pmod{3}$, and $B(a_i, k)$ if $k \notin [3(i-1), 3i)$ and $k \not\equiv 2 \pmod{3}$. Also, for $i \in I_0$ and $k \in [3(i-1), 3i)$, \mathcal{D}_{I_0} contains $B(a_i, k)$ if $k \equiv 1 \pmod{3}$. Finally, for each $i \in I_1$ and $k \in [3(i-1), 3i)$, it has $B(a_i, k)$ if $k \equiv 0 \pmod{3}$. For example, \mathcal{D}_{I_0} for $n = 3$ and $I_0 = \{1, 3\}$ is shown below:



Then $\mathcal{D}_{I_0}, a_0, 0 \models q$ but $\mathcal{D}_{I_0}, a_0, 0 \not\models q'$. To illustrate, take $s = B \wedge \diamond (A \wedge \diamond \circ (B \wedge \diamond (A \wedge \circ (B \wedge \diamond A))))$, for which $q' = q \wedge \exists P.s \in \mathcal{Q}$ with $I_0 = \{1, 3\}$. One can readily see that q' cannot be satisfied at $a_0, 0$ in \mathcal{D}_{I_0} depicted above.

To complete the proof that $|E^+| \geq 2^n$, we prove the following property (cf. Example 17):

$$\mathcal{D} \models q \ \& \ \mathcal{D} \not\models q' \ \& \ (q' \neq q'') \Rightarrow \mathcal{D} \models q''$$

for all data instances \mathcal{D} and $q', q'' \in \mathcal{Q}$. Indeed, take an arbitrary q such that $\mathcal{D} \models q$ and let $\{a_0, \dots, a_n\} \subseteq ind(\mathcal{D})$ be the elements such that $P(a_0, a_i, 0)$ is in \mathcal{D} and $\mathcal{D}, a_i, 0 \models q_i^n$, for each $1 \leq i \leq n$. Let $(K_{i,1}, \dots, K_{i,n})$ be a vector of pairs/triples over \mathbb{N} such that, for $j \neq i$, $K_{i,j} = (b, b', a)$ and $K_{i,i} = (b, a)$, where:

- $b < b' < a, b < a$ in every $K_{i,j}$;
- b from $K_{i,j}$ is strictly greater than a from $K_{i,j-1}$;
- a from $K_{i,j}$ is strictly smaller than b from $K_{i,j+1}$;
- $\mathcal{D}, a_i, b \models B, \mathcal{D}, a_i, b' \models B, \mathcal{D}, a_i, a \models A$ for all $(b, b', a) \in K_{i,j}$;
- $\mathcal{D}, a_i, b \models B, \mathcal{D}, a_i, a \models A$ for all $(b, a) \in K_{i,i}$.

Note that, for each $1 \leq i \leq n$, a vector $(K_{i,1}, \dots, K_{i,n})$ as above exists. Take $q' \in \mathcal{Q}$ such that $\mathcal{D} \not\models q'$ and let I_0 be the set corresponding to q' . We observe that $\mathcal{D}, a_i, 0 \not\models s$, for all $1 \leq i \leq n$. Therefore, for any $i \in I_0$ and any vector $(K_{i,1}, \dots, K_{i,n})$, we have that $b > a + 1$ for a from $K_{i,i-1}$ and b from $K_{i,i}$, if $i > 0$, and $b > 0$ for b from $K_{i,i}$, if $i = 1$. Similarly, for any $i \in [1, n] \setminus I_0$ and any vector $(K_{i,1}, \dots, K_{i,n})$, we have that $b = a + 1$ for a from $K_{i,i-1}$ and b from $K_{i,i}$, if $i > 0$, and $b = 0$ for b from $K_{i,i}$, if $i = 1$. Now, take an arbitrary $q'' = q \wedge \exists P.s'' \in \mathcal{Q}$ such that $q'' \neq q'$ and let J_0 be the set corresponding to it. Clearly, $J_0 \neq I_0$. Suppose first that there exists $i \in J_0 \setminus I_0$. Therefore, $i \in [1, n] \setminus I_0$ and then, as it is easy to see, $\mathcal{D}, a_i, 0 \models s''$. Thus, $\mathcal{D}, a_0, 0 \models q''$ as was required. Now, if there is $i \in I_0 \setminus J_0$, the proof is analogous and left to the reader.

E.5 Proof of Theorem 27

Recall that we consider the class $\mathcal{Q}_p[\circ, \diamond_r](\mathcal{ELI})$ of queries of the form

$$q = r_0 \wedge o_1(r_1 \wedge o_2(r_2 \wedge \dots \wedge o_n r_n)), \quad (11)$$

where the r_i are \mathcal{ELI} -queries and $o_i \in \{\circ, \diamond, \diamond_r\}$. We first generalise the normal form introduced for $\mathcal{Q}_p[\diamond_r, \circ]$. Any q in $\mathcal{Q}_p[\circ, \diamond_r](\mathcal{ELI})$ can be represented as a sequence

$$r_0(t_0), R_1(t_0, t_1), \dots, r_{n-1}(t_{n-1}), R_n(t_{n-1}, t_n), r_n(t_n),$$

where $R_i \in \{suc, <, \leq\}$ and r_i is an \mathcal{ELI} -query. As before, we divide q into blocks q_i such that

$$q = q_0 \mathcal{R}_1 q_1 \dots \mathcal{R}_n q_n \quad (12)$$

with $\mathcal{R}_i = R_1^i(t_0^i, t_1^i) \dots R_{n_i}^i(t_{n_i-1}^i, t_{n_i}^i)$, for $R_j^i \in \{<, \leq\}$,

$$q_i = r_0^i(s_0^i) suc(s_0^i, s_1^i) r_1^i(s_1^i) \dots suc(s_{k_i-1}^i, s_{k_i}^i) r_{k_i}^i(s_{k_i}^i)$$

and $s_{k_i}^i = t_0^{i+1}, t_{n_i}^i = s_0^i$. If $k_i = 0$, the block q_i is primitive. A primitive block $q_i = r_0^i(s_0^i)$ with $i > 0$ such that r_0^i is not equivalent to a conjunction of \mathcal{ELI} -queries that are not equivalent to r_0^i is called a lone conjunct. Now, we say that q is in normal form if the following conditions hold:

- (n1') $r_0^i \neq \top$ if $i > 0$, and $r_{k_i}^i \neq \top$ if $i > 0$ or $k_i > 0$ (thus, of all the first/last r in a block only r_0^0 can be trivial);
- (n2') each \mathcal{R}_i is either a single $t_0^i \leq t_1^i$ or a sequence of $<$;

(n3') $r_{k_i}^i \not\models r_0^{i+1}$ if q_{i+1} is primitive and R_{i+1} is \leq ;

(n4') $r_0^{i+1} \not\models r_{k_i}^i$ if $i > 0$, q_i is primitive and R_{i+1} is \leq .

Lemma 45. *Every query in $\mathcal{Q}_p[\circ, \diamond_r](\mathcal{ELI})$ is equivalent to a query in normal form that can be computed in polynomial time.*

We call a query in $\mathcal{Q}_p[\circ, \diamond_r](\mathcal{ELI})$ *safe* if it is equivalent to a query in normal form in $\mathcal{Q}_p[\circ, \diamond_r](\mathcal{ELI})$ without lone conjuncts.

Theorem 27. (i) *A query $q \in \mathcal{Q}_p[\circ, \diamond_r](\mathcal{ELI})$ is uniquely characterisable within $\mathcal{Q}_p[\circ, \diamond_r](\mathcal{ELI})$ iff q is safe.*

(ii) *Those queries that are uniquely characterisable within $\mathcal{Q}_p[\circ, \diamond_r](\mathcal{ELI})$ are actually polynomially characterisable within $\mathcal{Q}_p[\circ, \diamond_r](\mathcal{ELI})$.*

(iii) *The class $\mathcal{Q}_p[\circ, \diamond_r](\mathcal{ELI})$ is polynomially characterisable for bounded query size.*

(iv) *The class $\mathcal{Q}_p[\circ, \diamond](\mathcal{ELI})$ is polynomially characterisable.*

We show how the construction of positive and negative examples provided in the proof of Theorem 8 can be generalised.

Suppose q in normal form (12) does not contain lone conjuncts. Let b be again the number of \circ and \diamond in q plus 1. We construct $E = (E^+, E^-)$ characterising q as follows. Pick an individual name a and let, for every \mathcal{ELI} -query r in q , \hat{r} denote the pointed tree-shaped data instance defined by r (note that we take the same individual a for every r).

For each block q_i in (12), we take two temporal data instances

$$\begin{aligned} \bar{q}_i &= \hat{r}_0^i \dots \hat{r}_{k_i}^i \\ \bar{q}_i \times \bar{q}_{i+1} &= \hat{r}_0^i \dots \hat{r}_{k_i}^i \cup \hat{r}_0^{i+1} \dots \hat{r}_{k_{i+1}}^{i+1}. \end{aligned}$$

The set E^+ contains the data instances given by

- $\mathcal{D}_b = \bar{q}_0 \emptyset^b \dots \bar{q}_i \emptyset^b \bar{q}_{i+1} \dots \emptyset^b \bar{q}_n$,
- $\mathcal{D}_i = \bar{q}_0 \emptyset^b \dots \bar{q}_i \times \bar{q}_{i+1} \dots \emptyset^b \bar{q}_n$ if \mathcal{R}_{i+1} is \leq ,
- $\mathcal{D}_i = \bar{q}_0 \emptyset^b \dots \bar{q}_i \emptyset^{n_{i+1}} \bar{q}_{i+1} \dots \emptyset^b \bar{q}_n$ otherwise.

The set E^- contains all data instances of the form

- $\mathcal{D}_i^- = \bar{q}_0 \emptyset^b \dots \bar{q}_i \emptyset^{n_{i+1}-1} \bar{q}_{i+1} \dots \emptyset^b \bar{q}_n$ if $n_{i+1} > 1$;
- $\mathcal{D}_i^- = \bar{q}_0 \emptyset^b \dots \bar{q}_i \times \bar{q}_{i+1} \dots \emptyset^b \bar{q}_n$ if \mathcal{R}_{i+1} is a single $<$,

and also the data instances obtained from \mathcal{D}_b by

- (a) replacing $\hat{r}_j^i \neq \emptyset$ by an element of $\mathcal{F}(\hat{r}_j^i)$ or removing the whole $\hat{r}_j^i = \emptyset$, for $i \neq 0$ and $j \neq 0$, from some \bar{q}_i ;
- (b) replacing $\bar{q}_i = \hat{r}_0^i \dots \hat{r}_l^i \hat{r}_{l+1}^i \dots \hat{r}_{k_i}^i$ ($k_i > 0$) by $\bar{q}'_i \emptyset^b \bar{q}''_i$, where $\bar{q}'_i = \hat{r}_0^i \dots \hat{r}_l^i$, $\bar{q}''_i = \hat{r}_{l+1}^i \dots \hat{r}_{k_i}^i$ and $l \geq 0$;
- (c) replacing some $\hat{r}_l^i \neq \emptyset$, $0 < l < k_i$, by $\hat{r}_l^i \emptyset^b \hat{r}_l^i$;
- (d) replacing $\hat{r}_{k_i}^i$ ($k_i > 0$) with $\mathcal{A} \emptyset^b \hat{r}_{k_i}^i$, for some $\mathcal{A} \in \mathcal{F}(\hat{r}_{k_i}^i)$, or replacing \hat{r}_0^i ($k_i > 0$) with $\hat{r}_0^i \emptyset^b \mathcal{A}$, for some $\mathcal{A} \in \mathcal{F}(\hat{r}_0^i)$;
- (e) replacing $\hat{r}_0^0 \neq \emptyset$ with $\mathcal{A} \emptyset^b \hat{r}_0^0$, for $\mathcal{A} \in \mathcal{F}(\hat{r}_0^0)$, if $k_0 = 0$, and with $\hat{r}_0^0 \emptyset^b \hat{r}_0^0$ if $k_0 > 0$.

Let q be of the form (12). We generalise the notion of a homomorphism as follows. Let $\mathcal{D} = \mathcal{A}_0, \dots, \mathcal{A}_n$ and $a \in \text{ind}(\mathcal{D})$. A mapping h from the set $\text{var}(q)$ of variables in q to $[0, \max(\mathcal{D})]$ is a *generalised homomorphism* from q to \mathcal{D} for a if $h(t_0) = 0$, $\mathcal{A}_{h(t)} \models r(a)$ if $r(t)$, $h(t') = h(t) + 1$ if $\text{suc}(t, t') \in q$, and $h(t) R h(t')$ if $R(t, t') \in q$ for $R \in \{<, \leq\}$. Then one can show that $\mathcal{D}, a, 0 \models q$ if there exists a generalised homomorphism from q to \mathcal{D} for a .

It is now almost trivial to extend the proof of Theorem 8 to a proof of Theorem 27 by replacing homomorphisms by generalised homomorphisms. For example, block surjectivity is generalised in a straightforward way as follows: a generalised homomorphism $h: q' \rightarrow \mathcal{D}_b$ is *block surjective* if every point in every block \bar{q}_i of \mathcal{D}_b is in the range $\text{ran}(h)$ of h . To define type surjectivity let, for $\ell \in \text{ran}(h)$, r_ℓ denote the conjunction of all \mathcal{ELI} -queries r' with $r'(t)$ in q' such that $h(t) = \ell$. Clearly $\mathcal{A}_\ell \models r_\ell(a)$. h is *type surjective* if $\mathcal{A} \not\models r_\ell(a)$, for every $(\mathcal{A}, a) \in \mathcal{F}(\hat{r})$.

Theorem 30. *Let $n > 0$ be fixed. For every set Q of $\mathcal{EL}(\Sigma)$ -queries with $|Q| \leq n$, one can compute in polynomial time a split partner $\mathcal{S}(Q)$ of Q in $\mathcal{EL}(\Sigma)$. For \mathcal{ELI} , one can compute a split partner in exponential time.*

Proof. We prove the statement for $n = 1$, the generalisation is straightforward. Let $Q = \{q\}$. The construction is by induction over the depth of q . Assume $\text{depth}(q) = 0$. Thus $q = \bigwedge_{i=1}^k A_i$ with A_i atomic concepts. Then let for $i \leq k$:

$$\begin{aligned} \mathcal{A}_{A_i} &= \{B(a) \mid B \in \Sigma \setminus \{A\}\} \cup \\ &\quad \{R(a, b), R(b, b) \mid R \in \Sigma\} \cup \\ &\quad \{B(b) \mid B \in \Sigma\} \end{aligned}$$

and set $\mathcal{S}(q) = \{(\mathcal{A}_{A_i}, a) \mid 1 \leq i \leq k\}$. We show that $\mathcal{S}(q)$ is as required. Assume

$$q' = \bigwedge_{i=1}^{m_1} B_i \wedge \bigwedge_{i=1}^{m_2} \exists R_i \cdot q_i$$

If $q' \not\models q$, then there exist A_i with $A_i \notin \{B_i \mid 1 \leq i \leq m_1\}$. Then $\mathcal{A}_{A_i} \models q'(a)$, as required. Conversely, if $\mathcal{A}_{A_i} \models q'(a)$ for some A_i , then $A_i \notin \{B_i \mid 1 \leq i \leq m_1\}$. Hence $q' \not\models q$, as required.

Assume now that $\text{depth}(q) = n + 1$ and that split partners $\mathcal{S}(q')$ have been defined for queries q' of depth $\leq n$. Assume

$$q = \bigwedge_{i=1}^{n_1} A_i \wedge \bigwedge_{i=1}^{n_2} \exists S_i \cdot q_i.$$

Then assume $\mathcal{S}(q_i) = \{(\mathcal{A}_1, a_1), \dots, (\mathcal{A}_{k_i}, a_{k_i})\}$ and let c_1, \dots, c_{k_i} be fresh individuals. Define for $i \leq n_2$ the data instance

$$\begin{aligned} \mathcal{A}_i &= \{B(a) \mid B \in \Sigma\} \cup \\ &\quad \{R(a, b), S(b, b), B(b) \mid R \in \Sigma \setminus \{S_i\}, B, S \in \Sigma\} \cup \\ &\quad \{S_i(a, c_j) \mid 1 \leq j \leq k_i\} \cup \\ &\quad \mathcal{A}_1(c_1/a_1) \cup \dots \cup \mathcal{A}_{k_i}(c_{k_i}/a_{k_i}) \end{aligned}$$

with $\mathcal{A}(c/a)$ the result of replacing a by c in \mathcal{A} . Let $\mathcal{S}(q)$ be the union of $\mathcal{S}(\bigwedge_{i=1}^{n_1} A_i)$ and $\{(\mathcal{A}_i, a) \mid 1 \leq i \leq n_2\}$.

We show that $\mathcal{S}(\mathbf{q})$ is as required. Assume

$$\mathbf{q}' = \bigwedge_{i=1}^{m_1} B_i \wedge \bigwedge_{i=1}^{m_2} \exists R_i. \mathbf{q}'_i$$

If $\mathbf{q}' \not\models \mathbf{q}$, then either there exists A_i with $A_i \notin \{B_i \mid 1 \leq i \leq m_1\}$ or there exists $\exists S_i. \mathbf{q}_i$ such that for every $R_j. \mathbf{q}_j$ with $S_i = R_j$, we have $\mathbf{q}'_j \not\models \mathbf{q}_i$. In the former case $\mathcal{A}_{A_i} \models \mathbf{q}'(a)$ and in the latter case, by induction hypothesis, $\mathcal{A}_i \models \mathbf{q}'(a)$. Conversely, if $\mathcal{A} \not\models \mathbf{q}'(a)$ for some $\mathcal{A}, a \in \mathcal{S}(\mathbf{q})$, then either there exists A_i with $A_i \notin \{B_i \mid 1 \leq i \leq m_1\}$ or there exists $\exists S_i. \mathbf{q}_i$ such that for every $R_j. \mathbf{q}_j$ with $S_i = R_j$, there exists $\mathcal{A}, a \in \mathcal{S}(\mathbf{q}_i)$ such that $\mathcal{A} \not\models \mathbf{q}'_i(a)$. By induction hypothesis, $\mathbf{q}'_i \not\models \mathbf{q}_i$. But then $\mathbf{q}' \not\models \mathbf{q}$, as required. \square

E.6 Proof of Theorem 28

Suppose we are given a $\mathcal{P}^\Sigma[\mathcal{U}](\mathcal{EL})$ -query

$$\mathbf{q} = \mathbf{r}_0 \wedge (\mathbf{l}_1 \cup (\mathbf{r}_1 \wedge (\mathbf{l}_2 \cup (\dots (\mathbf{l}_n \cup \mathbf{r}_n) \dots))))).$$

To show that it is uniquely characterised by the polynomial-size example set with $(\mathbf{p}'_0) - (\mathbf{p}'_2)$ and $(\mathbf{n}'_0) - (\mathbf{n}'_2)$, consider any $\mathcal{Q}^\Sigma[\mathcal{U}](\mathcal{EL})$ -query

$$\mathbf{q}' = \mathbf{r}'_0 \wedge (\mathbf{l}'_1 \cup (\mathbf{r}'_1 \wedge (\mathbf{l}'_2 \cup (\dots (\mathbf{l}'_m \cup \mathbf{r}'_m) \dots))))).$$

such that $\mathbf{q} \not\models \mathbf{q}'$.

We now define a map f that reduces the 2D case to the 1D case. Consider the alphabet

$$\Gamma = \{\mathbf{r}_0, \dots, \mathbf{r}_n, \mathbf{l}_1, \dots, \mathbf{l}_n, \mathbf{r}'_0, \dots, \mathbf{r}'_m, \mathbf{l}'_1, \dots, \mathbf{l}'_m\} \setminus \{\perp\}$$

in which we regard the \mathcal{EL} -queries $\mathbf{r}_i, \mathbf{l}_i, \mathbf{r}'_j, \mathbf{l}'_j$ as symbols.

Let $\hat{\Gamma} = \{(\hat{\mathbf{a}}, a) \mid \mathbf{a} \in \Gamma\}$, that is, $\hat{\Gamma}$ consists of the pointed databases corresponding to the \mathcal{EL} -queries $\mathbf{a} \in \Gamma$.

For any \mathcal{EL} instance query \mathbf{a} , we set

$$f(\mathbf{a}) = \{\mathbf{b} \in \Gamma \mid (\hat{\mathbf{a}}, a) \models \mathbf{b}\}.$$

Similarly, for any \mathcal{EL} pointed data instance (\mathcal{A}, a) , we set

$$f(\mathcal{A}, a) = \{\mathbf{b} \in \Gamma \mid (\mathcal{A}, a) \models \mathbf{b}\}$$

and, for any temporal data instance $\mathcal{D} = (\delta_0, \dots, \delta_k)$ with \mathcal{EL} pointed data instances (δ_i, a_i) , set

$$f(\mathcal{D}) = (f(\delta_0, a_0), \dots, f(\delta_k, a_k)),$$

which is an *LTL*-data instance over the signature Γ . Finally, we define a query

$$f(\mathbf{q}) = \rho_0 \wedge (\lambda_1 \cup (\rho_1 \wedge (\lambda_2 \cup (\dots (\lambda_n \cup \rho_n) \dots))))$$

by taking $\rho_i = f(\mathbf{r}_i)$ and $\lambda_i = f(\mathbf{l}_i)$, and similarly for \mathbf{q}' .

By definition, $f(\mathbf{q})$ is a $\mathcal{P}^\Gamma[\mathcal{U}]$ -query (indeed, since $\hat{\mathbf{r}}_i, a \not\models \mathbf{l}_i$, we have $\mathbf{l}_i \in f(\mathbf{l}_i) \setminus f(\mathbf{r}_i)$, and since $\hat{\mathbf{l}}_i, a \not\models \mathbf{r}_i$, we have $\mathbf{r}_i \in f(\mathbf{r}_i) \setminus f(\mathbf{l}_i)$), and $f(\mathbf{q}')$ is a $\mathcal{Q}_p^\Gamma[\mathcal{U}]$ -query such that $f(\mathbf{q}) \not\models f(\mathbf{q}')$: it follows immediately from the definition that, for any data instance \mathcal{D} , we have $\mathcal{D} \models \mathbf{q}$ iff $f(\mathcal{D}) \models f(\mathbf{q})$ and similarly for \mathbf{q}' . By Theorem 14, \mathbf{q} and \mathbf{q}' are separated by the corresponding example set with $(\mathbf{p}_0) - (\mathbf{p}_2)$ and $(\mathbf{n}_0) - (\mathbf{n}_2)$. Notice that the positive examples from $(\mathbf{p}_0) - (\mathbf{p}_2)$ are exactly the f -images of the examples $(\mathbf{p}'_0) - (\mathbf{p}'_2)$. So if $f(\mathbf{q})$ and $f(\mathbf{q}')$ are separated by some \mathcal{D} from

$(\mathbf{p}_0) - (\mathbf{p}_2)$, the corresponding member of $(\mathbf{p}'_0) - (\mathbf{p}'_2)$ separates \mathbf{q} and \mathbf{q}' .

So suppose $f(\mathbf{q})$ and $f(\mathbf{q}')$ are separated by some \mathcal{D} from $(\mathbf{n}_0) - (\mathbf{n}_2)$. If $\mathcal{D} = \Gamma^n$, then it means that the temporal depth of $f(\mathbf{q}')$ is less than the temporal depth of $f(\mathbf{q})$, so $m < n$, and the queries \mathbf{q} and \mathbf{q}' are separated by \mathcal{A}_Σ^n .

Suppose $\mathcal{D} = \Gamma^{n-i}(\Gamma \setminus \{\mathbf{a}\})\Gamma^i$. Since $\mathcal{D} \not\models f(\mathbf{q})$, we have $f(\mathbf{r}_i) \not\subseteq \Gamma \setminus \{\mathbf{a}\}$, and so $\hat{\mathbf{r}}_i, a \models \mathbf{a}$ but $f(\mathbf{r}'_i) \subseteq \Gamma \setminus \{\mathbf{a}\}$. Then $\hat{\mathbf{r}}'_i, a \not\models \mathbf{a}$, and so $\mathbf{r}'_i \not\models \mathbf{r}_i$. Therefore, there is $(\mathcal{A}, a) \in \mathcal{S}(\{\mathbf{r}_i\})$ such that $\mathcal{A} \models \mathbf{r}'_i(a)$ and $\mathcal{A}_\Sigma^{n-i} \mathcal{A} \mathcal{A}_\Sigma^i$ separates \mathbf{q} and \mathbf{q}' .

The cases when \mathcal{D} is from (\mathbf{n}_1) or (\mathbf{n}_2) are treated in a similar manner.

If $\mathbf{q} \in \mathcal{P}^\Sigma[\mathcal{U}](\mathcal{EL})$, the characterisation is exponential as the size of $(\mathbf{n}_0) - (\mathbf{n}_2)$ is exponential because the exponential size of constructed split partners for \mathcal{EL} -queries.

F Proofs for Section 9

We provide further details of the proof of Theorem 33.

Theorem 33. (i) *The class of safe queries in $\mathcal{Q}_p[\square, \diamond_r](\mathcal{EL})$ is polynomial-time learnable with membership queries.*

(ii) *The class $\mathcal{Q}_p[\square, \diamond_r](\mathcal{EL})$ is polynomial-time learnable with membership queries if the learner knows the size of the target query in advance.*

(iii) *The class $\mathcal{Q}_p[\square, \diamond](\mathcal{EL})$ is polynomially-time learnable with membership queries.*

We start by completing the proof for the 1D case. For (i), it remains to consider **step 4**. At that point of the computation, the algorithm has identified all blocks of \mathbf{q} but not the sequences of \diamond and \diamond_r between them. Suppose that $\mathcal{D} = \bar{\mathbf{q}}_0 \emptyset^b \dots \bar{\mathbf{q}}_i \emptyset^b \bar{\mathbf{q}}_{i+1} \dots \emptyset^b \bar{\mathbf{q}}_n$. We construct \mathbf{q}' from the blocks of \mathcal{D} by selecting \mathcal{R}_{i+1} as follows: If $\mathcal{D}_i \models \mathbf{q}$ for $\mathcal{D}_i = \bar{\mathbf{q}}_0 \emptyset^b \dots \bar{\mathbf{q}}_i \bowtie \bar{\mathbf{q}}_{i+1} \dots \emptyset^b \bar{\mathbf{q}}_n$ then set \mathcal{R}_{i+1} to be \leq . Otherwise, let n_{i+1} be the smallest number such that $\mathcal{D}_i \models \mathbf{q}$ for $\mathcal{D}_i = \bar{\mathbf{q}}_0 \emptyset^b \dots \bar{\mathbf{q}}_i \emptyset^{n_{i+1}} \bar{\mathbf{q}}_{i+1} \dots \emptyset^b \bar{\mathbf{q}}_n$. Set \mathcal{R}_{i+1} to be a sequence of $<$ of length n_{i+1} . It is now easy to see that \mathbf{q}' fits the example set (E^+, E^-) and so $\mathbf{q}' \equiv \mathbf{q}$, as required.

For (ii) we have to replace **step 3** by a computation step that ensures that after applying (f) one does not obtain a data instance \mathcal{D} with $\mathcal{D} \models \mathbf{q}$. Recall that for bound $n = |\mathbf{q}|$ and $\mathbf{q}_i = \rho(s)$ a lone conjunct in \mathbf{q} with $\rho = \{A_1, \dots, A_k\}$ the rule (f) is defined as follows:

(f) replace ρ with $(\rho \setminus \{A_1\}) \emptyset^b \dots \emptyset^b \rho \setminus \{A_k\}^n$.

is satisfied. Now, after having computed \mathcal{D} in **step 2**, and a data instance \mathcal{D}' is obtained from \mathcal{D} by applying rule (f) such that $\mathcal{D}' \models \mathbf{q}$, then replace \mathcal{D} with \mathcal{D}' and return to **step 2**. If no such \mathcal{D}' exists, proceed to **step 4**. To bound the number of applications of rule (f) notice that at the end of **step 2**, the number of time points in \mathcal{D} other than \emptyset does not exceed $|\mathbf{q}|$. Indeed, any time point in \mathcal{D} not in the range of some $h : \mathbf{q} \rightarrow \mathcal{D}$ would be eliminated by rule (a). We obtain a polynomial bound on the number of applications by observing that each application of rule (f) removes a symbol from ρ .

We now move to the 2D case. The proof extends the argument given above for the 1D case and uses the example set $E = (E^+, E^-)$ defined in the proof of Theorem 27. Intuitively, whenever in the argument above we replace ρ by $\rho \setminus \{A\}$ for a set ρ of atoms, we now replace \hat{r} by an element of the frontier $\mathcal{F}(\hat{r})$. There is only one difficulty: in the 1D case the algorithm starts with data instances $\sigma \cdots \sigma$ with σ the signature of the target query, whereas now we have to start with data instances $\mathcal{A}_\Sigma \cdots \mathcal{A}_\Sigma$ with Σ the signature of \mathbf{q} and $\mathcal{A}_\Sigma = \{R(a, a), A(a) \mid R, A \in \Sigma\}$. The atemporal data instance \mathcal{A}_Σ , however, is not tree-shaped, and we have not yet discussed frontiers for data instances that are not tree-shaped. Indeed, in (ten Cate and Dalmau 2021), frontiers are not only computed for tree-shaped data instances but for a generalisation called *c-acyclic* data instances with cycles through the distinguished node. We could at this point introduce the relevant machinery from (ten Cate and Dalmau 2021) and work with frontiers for c-acyclic data instances. Instead, we show that one can use the machinery we have introduced already and work with frontiers of tree-shaped data instances. But we require a straightforward intermediate step that transfers the data instance \mathcal{A}_Σ into a tree-shaped data instance \mathcal{D} with $\mathcal{D} \models \mathbf{q}$ using membership queries.

We adjust **step 1** of the learning algorithm from the 1D case as follows. Assume \mathbf{q} is the target query and let $\Sigma = \text{sig}(\mathbf{q})$. We aim to identify an initial temporal data instance $\mathcal{D}_0 = \mathcal{A}_0, \dots, \mathcal{A}_n$ with a designated individual a such that

- $\mathcal{D}_0, a, 0 \models \mathbf{q}$,
- if \mathcal{D}' is obtained from \mathcal{D}_0 by removing an atom then $\mathcal{D}', a, 0 \not\models \mathbf{q}$, and
- all \mathcal{A}_i are tree-shaped with distinguished node a .

By asking incrementally membership queries of the form ‘ $\mathcal{A}_\Sigma^k, a, 0 \models \mathbf{q}$?’, we identify the number of \circ and \diamond in \mathbf{q} . Let $b = \min\{k \mid (\mathcal{A}_\Sigma)^k, a, 0 \models \mathbf{q}\} + 1$.

Let $\mathcal{D}_0 = \mathcal{A}_0, \dots, \mathcal{A}_n$, where $n = b - 2$ and, initially, $\mathcal{A}_i = \mathcal{A}_\Sigma$ for $i = 1, \dots, n$. Before progressing, we make \mathcal{A}_i tree-shaped by applying the following unwind and minimise operations:

unwind Suppose that \mathcal{A}_i contains an atom $S(c, c)$. Then introduce fresh individuals c_R, c_{R^-} for every binary predicate R with $R(c, c) \in \mathcal{A}_i$, remove all $R(c, c)$ from \mathcal{A}_i , and add instead $R(c, c_R)$, $R(c_{R^-}, c)$, and $A(c_R)$, $A(c_{R^-})$, $R'(c_R, c_R)$, and $R'(c_{R^-}, c_{R^-})$ for all $A, R' \in \Sigma$. Let \mathcal{A}'_i denote the resulting data instance and set $\mathcal{D}' = \mathcal{A}_0, \dots, \mathcal{A}_{i-1}, \mathcal{A}'_i, \mathcal{A}_{i+1}, \dots, \mathcal{A}_n$.

minimise Remove exhaustively atoms from \mathcal{A}'_i as long as $\mathcal{D}', a, 0 \models \mathbf{q}$.

Observe that if $\mathcal{D} \models \mathbf{q}$ and \mathcal{D}' is obtained from \mathcal{D} by an unwinding step, then $\mathcal{D}' \models \mathbf{q}$. After the minimise step, the size of \mathcal{D}' does not exceed the size of \mathbf{q} . Therefore we can replace \mathcal{D}_0 with \mathcal{D}' . By applying the unwind-minimise steps exhaustively, we eventually eliminate all loops from \mathcal{A}_i . It remains to notice that the minimise step can be implemented by querying the membership oracle.

The remaining steps of the learning algorithm remain the same as before by replacing the removal of an atom by taking an element of the frontier of a tree-shaped data instance.

For (iii), notice that in part (ii) the learner actually only needs to know the temporal depth of the goal query, not the overall size. Hence (iii) also reduces to (ii).