

Containment of Expressive SPARQL Navigational Queries

Melisachew Wudage Chekol¹(✉) and Giuseppe Pirrò²

¹ Data and Web Science Group, University of Mannheim, Mannheim, Germany
mel@informatik.uni-mannheim.de

² Institute for High Performance Computing and Networking,
ICAR-CNR, Rende, Italy
pirro@icar.cnr.it

Abstract. Query containment is one of the building block of query optimization techniques. In the relational world, query containment is a well-studied problem. At the same time it is well-understood that relational queries are not enough to cope with graph-structured data, where one is interested in expressing queries that capture *navigation* in the graph. This paper contributes a study on the problem of query containment for an expressive class of navigational queries called Extended Property Paths (EPPs). EPPs are more expressive than previous navigational extension of SPARQL (e.g., nested regular expressions) as they allow to express path conjunction and path negation, among others. We attack the problem of EPPs containment and provide complexity bounds.

1 Introduction

Research in graph query languages has emerged as a consequence of the intrinsic limitations of relational query languages when it comes to the possibility to express recursion and navigation. This lead to the design of languages like Regular Path Queries (RPQs) and their extensions, including 2RPQs [2] that include the possibility to traverse paths backwards, Extended Regular Path Queries (ERPQs) [1] that offer higher expressive power by allowing to express queries that also capture graphs [11, 13]. Another well-studied class of queries is that of Nested Regular Expressions (NREs) [26], that were originally proposed as the navigational core of SPARQL. The current SPARQL 1.1 standard introduced Property Paths (PPs) as navigational core; PPs offer very limited expressive power due to the lack of features to express any type of test within a path. To cope with these issues and design a language as close as possible to the current W3C standard, the language of Extended Property Paths (EPPs) has been recently proposed [12]. EPPs are more expressive than NREs, PPs and other navigational extensions of SPARQL thanks to the possibility to express path conjunction, negation and more powerful types of tests (e.g., checking the values of nodes reached via a nested expression) while keeping query evaluation tractable.

Related Work. The problem of query containment is one of the main pillars of static analysis and query optimization. In what follows we will not consider containment under constraints as we tackle containment of EPPs without constraints.

In the relational world, query containment for conjunctive queries (CQs) is now well-understood; it is NP-complete for basic CQs [5] and union of CQs [29] while it is Π_2^P when considering arithmetic comparison and also bag semantics [6, 19].

For graph queries, containment of 2RPQs has been shown to be PSPACE-complete [1, 2]; the complexity jumps to EXPTIME-hard under the presence of functionality constraints and to 2EXPTIME when considering expressive description logics constraints [4]. The problem becomes exponential if the query on the right hand side has a tree structure (cf. for example, [3]); for *extended* RPQs containment is undecidable for Boolean queries over a fixed alphabet [1].

In the Semantic Web, the containment of PSPARQL has been studied in Chekol et al. [8], this work provides lower bounds for upper bound complexity results reported in that paper; for NRE, Reutter [28] has shown a PSPACE upper bound; for SPARQL PPs, Kostylev et al. [20] show that the containment ranges from EXSPACE-complete for OPT-free queries to PSPACE-complete if the right-hand side query is a pattern without projection. The study in [27] provides complexity analysis for several fragments of SPARQL: the results range from NP-completeness for AND-UNION queries to undecidability for the full SPARQL.

Another related language is XPath; the problem of XPath 2.0 query containment has been studied in ten Cate and Lutz [30]. They showed that the introduction of path intersection alone (i.e., the language CoreXPath (\cap) leads to 2EXPTIME-completeness. Finally, Kostylev et al. [22] studied the problem of containment of navigational XPath queries (i.e., the GXpath language), with results ranging from undecidability (when negation is considered) to EXP-TIME completeness for the positive fragment.

Contributions. We study query containment for EPPs, a significant extension of property paths (PPs), the current navigational core of SPARQL, and NREs for which containment has been already studied. We resort on two main ingredients: (i) an encoding of EPPs into the μ -calculus; (ii) the notion of (RDF) transition systems to check the validity test of μ -calculus formulae and provide an *upper bound* on the containment of EPPs and SPARQL with EPPs. For *lower bounds* we make connections between EPPs, PDL and XPath 2.0.

Automata theoretic notions and a reduction into validity test in a logic have been widely used to address the problem of query answering and containment [2, 3, 10, 15, 20, 24]. Contrary to the automata techniques, the logic based approaches are fairly implementable. In this respect, it has been shown [16] that logical combinators can provide an exponential gain in succinctness in terms of the size of a logical formula thus allowing to study containment for expressive query languages in exponential-time, even though their direct formulation into the underlying logic results in an exponential blow up of the formula size.

Outline. The remainder of the paper is organized as follows. We provide some background in Sect. 2. The language of EPPs is introduced in Sect. 3. Section 4 describes an encoding of EPPs into the μ -calculus. Section 5 discusses the containment of EPPs. We conclude in Sect. 6.

2 Preliminaries

This section provides some background about the machineries used in this paper. We start with RDF and SPARQL and then give a brief overview about μ -calculus that will be used to tackle the containment of EPPs.

RDF and SPARQL. An RDF triple¹ is a tuple of the form $\langle s, p, o \rangle \in \mathbf{I} \times \mathbf{I} \times \mathbf{I} \cup \mathbf{L}$, where \mathbf{I} (IRIs) and \mathbf{L} (literals) are countably infinite sets. The set of terms of an RDF graph will be $terms(G)$ while $nodes(G)$ will be the set of terms used as a subject or object of a triple. An RDF graph G is a set of triples. To query RDF data, a standard query language, called SPARQL, has been defined. The semantics of a SPARQL query [25] is defined in terms of solution mappings. A (solution) mapping m is a partial function $m: \mathcal{V} \rightarrow \mathbf{I} \cup \mathbf{L}$. The SPARQL semantics uses a function $\llbracket Q \rrbracket_G$ that evaluates a query Q on a graph G and gives a multiset (bag) of mappings in the general case. However, when considering SPARQL with patterns using recursive PPs (i.e., using $*$, $+$) the standard introduces auxiliary functions (called ALP) that return sets of mapping.

μ -calculus. The μ -calculus (\mathcal{L}_μ) is a logic obtained by adding fixpoint operators to ordinary modal logic [23]. For the purpose of this paper we will make usage of the μ -calculus with nominals and converse programs [31]. The syntax of the μ -calculus includes countable sets of *atomic propositions* AP , a set of *variables* Var , a set of *programs and their respective converses* $Prog = \{s, p, o, \bar{s}, \bar{p}, \bar{o}\}$ used to allow navigation in a graph. A μ -calculus formula φ is defined inductively as:

$$\varphi ::= \top \mid q \mid X \mid \neg\varphi \mid \varphi \wedge \psi \mid \varphi \vee \psi \mid \langle a \rangle\varphi \mid [a]\varphi \mid \mu X\varphi \mid \nu X\varphi$$

where $q \in AP$, $X \in Var$ and $a \in Prog$ is a transition program or its converse \bar{a} . The greatest fixpoint ν and least fixpoint operator μ introduce general and finite recursion in graphs, respectively. The semantics of the μ -calculus is given over a transition system (aka Kripke structure) $K = (S, R, L)$, where S is a non-empty set of nodes, $R : Prog \rightarrow 2^{S \times S}$ is the transition function, and $L : AP \rightarrow 2^S$ assigns a set of nodes to each atomic proposition or nominal where it holds, such that $L(p)$ is a *singleton* for each nominal p . For converse programs, R can be extended as $R(\bar{a}) = \{(s', s) \mid (s, s') \in R(a)\}$ where $s, s' \in S$.

Definition 1 (*Model of a formula*). For a sentence φ and a transition system $K = (S, R, L)$, K is model of φ , denoted $K \models \varphi$, if there exists $s \in S$ such that $K, s \models \varphi$ if and only if $s \in \llbracket \varphi \rrbracket^K$ – s is an element of the answer to the evaluation of φ over K . If a sentence has a model, then it is called satisfiable.

If a μ -calculus formula ψ appears under the scope of a least μ or greatest ν fixed point operator over all the programs $\{s, p, o, \bar{s}, \bar{p}, \bar{o}\}$ as, $\mu X.\psi \vee \langle s \rangle X \vee \langle p \rangle X \vee \dots$ or $\nu X.\psi \wedge \langle s \rangle X \wedge \langle p \rangle X \wedge \dots$, then, for legibility, we denote the formulae by $lfp(X, \psi)$ and $gfp(X, \psi)$, respectively.

¹ We do not consider bnodes.

RDF Transition System. An RDF transition system [9] is a labeled transition system, $K_G = (S = S' \cup S'', R, L)$, representation of an RDF graph G where two sets of nodes, S' and S'' , are introduced: one set S'' for each triple (called triple node) and the other set S' for each subject, predicate, and object of each triple. A triple node (the black node in the figure below) is connected to its subject, predicate, and object nodes. For instance, the RDF graph $\boxed{x} \xrightarrow{z} \boxed{y}$ can be turned into the RDF transition

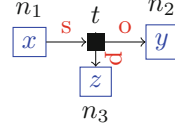


Fig. 1. Encoding of an RDF triple into a transition system.

system in Fig. 1, where $S' = \{n_1, n_2, n_3\}$, $S'' = \{t\}$, $R(s) = \{(n_1, t)\}$, $R(p) = \{(t, n_3)\}$, $R(o) = \{(t, n_2)\}$, and $L(x) = \{n_1\}$, $L(y) = \{n_2\}$, $L(z) = \{n_3\}$. Navigation from one node to another, in RDF transition systems, is done by using a set of transition programs $\{s, p, o\}$ and their converses. Since the μ -calculus with converse lacks functionality for number restrictions one cannot impose that each triple node is connected to exactly one node for each subject, predicate, and object node. However, one can impose a lighter restriction to achieve this by taking advantage of the technique introduced in [14] and adopted in [7]. Since it is not possible to ensure that there is only one successor, then we restrict all the successors to bear the same constraints; thus, they become interchangeable. This is achieved by rewriting the formulas using a function func such that all occurrences of $\langle a \rangle \varphi$ (existential formulas) are replaced by $\langle a \rangle \top \wedge [a] \text{func}(\varphi)$; in Definition 2, func is defined inductively on the structure of a μ -calculus formula. When checking for query containment, we assume that the formulas are rewritten using the function func .

Definition 2. func is inductively defined on the structure of a μ -calculus formula as follows:

$$\begin{array}{ll}
 \text{func}(\top) = \top & \text{func}(\perp) = \perp \\
 \text{func}(q) = q \quad q \in AP \cup \text{Nom} & \text{func}(X) = X \quad X \in \text{Var} \\
 \text{func}(\neg \varphi) = \neg \text{func}(\varphi) & \text{func}(\varphi \wedge \psi) = \text{func}(\varphi) \wedge \text{func}(\psi) \\
 \text{func}(\varphi \vee \psi) = \text{func}(\varphi) \vee \text{func}(\psi) & \text{func}(\langle a \rangle \varphi) = \langle a \rangle \top \wedge [a] \text{func}(\varphi) \quad a \in \text{Prog} \\
 \text{func}(\langle a \rangle \varphi) = \langle a \rangle \text{func}(\varphi) & \text{func}([a] \varphi) = [a] \text{func}(\varphi) \\
 \text{func}(\mu X. \varphi) = \mu X. \text{func}(\varphi) & \text{func}(\nu X. \varphi) = \nu X. \text{func}(\varphi)
 \end{array}$$

Figure 2 shows an RDF graph G and an excerpt of its corresponding RDF transition system K_G . Each RDF triple requires the introduction of a transition node (i.e., black nodes in the figure), where: the subject of the RDF triple has an incoming edge to the transition node labeled as s and the predicate and object have outgoing edges labeled as p and o , respectively. At this point, we shall define the model of formula ψ in terms of RDF transition systems.

Definition 3 (*RDF transition system model*). An RDF transition system K_G is considered as a model of formula ψ if there exists a node s in the transition system where ψ holds, i.e., $K_G, s \models \psi$. If a formula has a model, then it is called satisfiable.

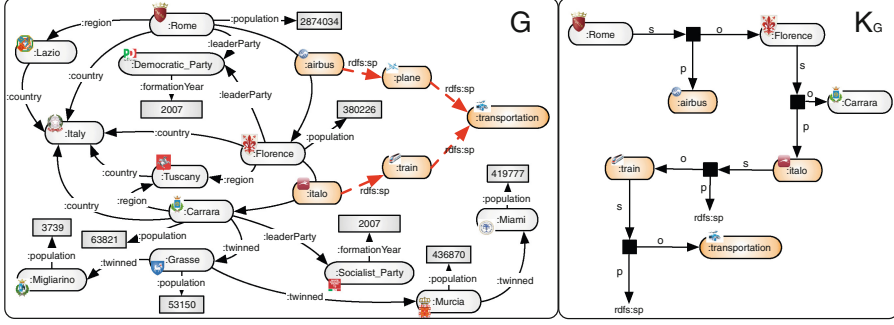


Fig. 2. An graph G and an excerpt of its corresponding RDF transition system K_G

The following lemma links RDF transition systems with transition systems when encoding in μ -calculus a query expressed in some RDF query language.

Lemma 1 (Chekol et al. [7]). *Let φ be a μ -calculus encoding of an EPPs query, φ is satisfied by some RDF transition system K_G if and only if $\text{func}(\varphi)$ is satisfied by some transition system K .*

The above lemma serves as a basis for the encoding of EPPs that will be encoded in μ -calculus formulas and then can be interpreted over RDF transition systems.

3 Extended Property Paths

We now introduce our navigational extension of SPARQL called Extended Property Paths (EPPs) [12].

Syntax. EPPs extend PP and NREs-like languages with path conjunction/negation, repetitions, and more types of tests. The syntax of EPPs is given below:

$$\begin{aligned}
 e &::= \text{'^'} e \mid e \text{'+'} \mid e \text{'?' } \mid e \text{'*'} \mid e \text{'/' } e \mid e \text{'|'} e \mid \\
 &\quad \text{'(' } e \text{')'} \mid [P]^1 \text{ test } [P]^2 \mid e \text{'\&'} e \mid e \text{'\sim'} e \mid e \text{'\{l, h\}'} \\
 \text{test} &::= \text{'!' } \text{test} \mid \text{test } \text{'\&\&'} \text{test} \mid \text{test } \text{'||'} \text{test} \mid \text{'(' test ')} \mid \text{base} \\
 \text{base} &::= \text{iri} \mid \text{'^'} \text{iri} \mid \text{'TP('P', 'e ')} \mid \text{'T('EEExp')'} \\
 P &::= \text{'_s'} \mid \text{'_p'} \mid \text{'_o'}
 \end{aligned}$$

¹Default is '_s' ; ²Default is '_o' . We refer to EPPs without path repetition, path complement as **cEPPs**.

EPPs introduce the following new features: path conjunction ($e_1 \& e_2$), path negation ($e_1 \sim e_2$), path repetitions between l and h times ($e \{l, h\}$), and different types of tests (**test**) within a path. EPPs allow to specify the *starting* and *ending* position (P) of a test; it is possible to test from each of the subject, predicate

and object positions in RDF triples, mapped in the EPPs syntax to the position symbols $_s$, $_p$ and $_o$, respectively. Positions do not need to be always specified; by default a test starts from the subject ($_s$) and ends on the object ($_o$) of the triple being traversed. There are different types of tests (**test**); a test can be a simple check for the existence of a IRI in forward/reverse direction, a *nested* EPP, i.e., $\text{TP}(\mathbf{P}, \mathbf{e})$, which corresponds to the evaluation of the expression \mathbf{e} starting from a position \mathbf{P} (of the last triple traversed) and returns true iff there exists at least one node that can be reached via \mathbf{e} . A test can also be of type **T**; here, **EExp** (not reported here for sake of space) extends the production [110] in the SPARQL grammar² which enables to use in EPPs tests available in SPARQL as built-in conditions. Tests can be combined via the logical operators AND ($\&$), OR ($\|$) and NOT ($!$).

Positions and Tests. EPPs tests can be coupled with positions. To formally explain the reasoning behind tests and positions, we make usage of a function $\Pi(\mathbf{P}, t)$, which projects the element in position \mathbf{P} of a triple t . If we have $t = \langle u_1, p_1, u_2 \rangle$ and the test $\text{T}(_p = p_1)$ then $\Pi(_p, \langle u_1, p_1, u_2 \rangle) = p_1$ that checks $p_1 = p_1$, and, in this case, returns true; however, it returns false for $\text{T}(_o = p_3)$.

Example 1 (Path Conjunction, Negation and Tests). Find pairs of cities located in the same country but not in the same region; such cities must be governed by the same political party, which has been founded before 2010.

NREs-based languages (and PPs) *cannot* express such request due to the lack of path conjunction/negation. With EPPs it can be expressed as shown in Fig. 3.

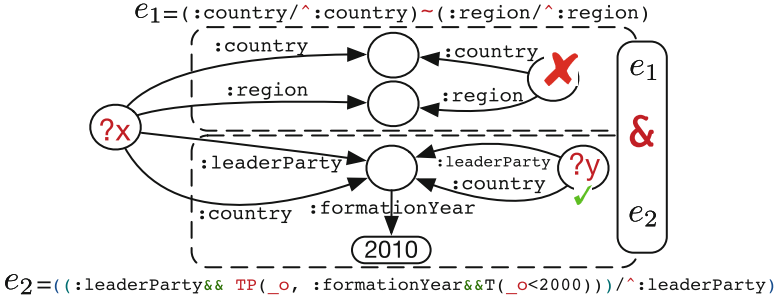


Fig. 3. EPP expression for Example 1.

The expression is the conjunction ($\&$) of the two sub-expressions e_1 and e_2 .

In the sub-expression e_1 , the symbol \wedge denotes backward navigation; from the object to the subject of a triple. Path negation \sim enables to discard from the set of cities in the same country (i.e., $:country/^:country$) those that are in the same region (i.e., $:region/^:region$). Path conjunction $\&$ enables to keep

² <http://www.w3.org/TR/sparql11-query/#rExpression> This is not considered when dealing with containment of EPPs.

from the set of nodes satisfying the first subexpression (e_1) those that also satisfy the second (e_2), i.e. governed by the same party, which has been founded before 2010. Note that PPs and NREs-based languages lack tests like **TP**, to check the existence of a path (via `:formationYear`) to nodes having value <2010 .

Example 2 (Positions and Tests). Consider again the expression in Example 1. The expression including both default positions³ and using positions to traverse backward edges is shown in Fig. 4.

As an example, the test **TP** (in Fig. 4) starts from the position `_o`, that is, the object of the last triple traversed.

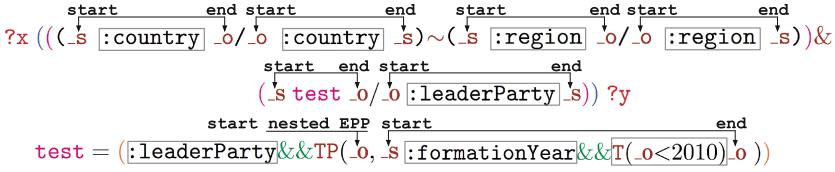


Fig. 4. An EPP expression with positions.

Definition 4. An EPP query q has the form $q ::= (\mathcal{V} \cup I, e, \mathcal{V} \cup I \cup L)$.

Semantics. For the purpose of this paper we will focus our attention on set semantics. This is in line with recent related work that studied the containment of PPs in SPARQL [20] where PPs were given a set-based semantics. Besides, the containment of conjunctive queries in the relational world when considering bag semantics is still an open issue⁴. Table 1 shows the semantics of EPPs. The semantics makes usage of two evaluation function; the first $\mathbf{E}[\cdot]^G$ is used to evaluate all forms of EPPs on an RDF graph G but tests; these latter are handled by the boolean evaluation function $\mathbf{E}_T[\cdot]_t^G$. Note that in this latter case tests consider a triple $t \in G$ and have a start and end position. By observing rule R9 one may notice that the final result is a set of pair of nodes in the graph G ; these pairs of nodes are obtained via the position mapping function Π that projects the two elements of the triple t appearing in position P_1 and P_2 if t satisfies the test.

4 Translating cEPPs into \mathcal{L}_μ

The goal of this section is to provide the first building block toward tackling the containment of EPP expressions; this concerns the encoding of **cEPPs** into \mathcal{L}_μ . We shall start by introducing the notion of path pattern translation.

³ These are automatically added during the parsing phase.

⁴ It is claimed that for basic CQs the complexity of query containment under bag semantics is in Π_2^P [6].

Table 1. Set-based semantics for EPPs.

R1	$\mathbf{E}[\hat{e}]^G := \{(u, v) : (v, u) \in \mathbf{E}[e]^G\}$
R2	$\mathbf{E}[e_1/e_2]^G := \{(u, v) : \exists w (u, w) \in \mathbf{E}[e_1]^G \wedge (w, v) \in \mathbf{E}[e_2]^G\}$
R3	$\mathbf{E}[(e)^*]^G := \{(u, u) \mid u \in \text{nodes}(G)\} \cup \bigcup_{i=1}^{\infty} \mathbf{E}[e_i]^G \mid$ $e_1 = e \wedge e_i = e_{i-1}/e$
R4	$\mathbf{E}[(e)^+]^G := \bigcup_{i=1}^{\infty} \mathbf{E}[e_i]^G \mid e_1 = e \wedge e_i = e_{i-1}/e$
R5	$\mathbf{E}[(e)?]^G := \{(u, u) \mid u \in \text{nodes}(G)\} \cup \mathbf{E}[e]^G$
R6	$\mathbf{E}[(e_1 e_2)]^G := \{(u, v) : (u, v) \in \mathbf{E}[e_1]^G \vee (u, v) \in \mathbf{E}[e_2]^G\}$
R7	$\mathbf{E}[(e_1 \& e_2)]^G := \{(u, v) : (u, v) \in \mathbf{E}[e_1]^G \wedge (u, v) \in \mathbf{E}[e_2]^G\}$
R8	$\mathbf{E}[(e_1 \sim e_2)]^G := \{(u, v) : (u, v) \in \mathbf{E}[e_1]^G \wedge (u, v) \notin \mathbf{E}[e_2]^G\}$
R9	$\mathbf{E}[P_1 \text{ test } P_2]^G := \{(\Pi(P_1, t), \Pi(P_2, t)) \mid \text{triple } t \in G \wedge \mathbf{E}_T[\text{test}]_t^G\}$
R10	$\mathbf{E}_T[u]_t^G := \Pi(\neg p, t) = u$
R11	$\mathbf{E}_T[T(\text{EExp})]_t^G := \text{EvalSPARQLBuilt-in}(\text{EExp}, t)$
R12	$\mathbf{E}_T[TP(p, e)]_t^G := \exists v : (\Pi(p, t), v) \in \mathbf{E}[e]^G$
R13	$\mathbf{E}_T[\text{test}_1 \& \& \text{test}_2]_t^G := \mathbf{E}_T[\text{test}_1]_t^G \wedge \mathbf{E}_T[\text{test}_2]_t^G$
R14	$\mathbf{E}_T[\text{test}_1 \mid \text{test}_2]_t^G := \mathbf{E}_T[\text{test}_1]_t^G \vee \mathbf{E}_T[\text{test}_2]_t^G$
R15	$\mathbf{E}_T[!\text{test}]_t^G := \neg \mathbf{E}_T[\text{test}]_t^G$

Definition 5 (*Path pattern translation*). Given a **cEPPs** query $q ::= (\alpha, e, \beta)$ where $\alpha \in \mathcal{V} \cup \mathbf{I}$ and $\beta \in \mathcal{V} \cup \mathbf{I} \cup \mathbf{L}$, its translation into \mathcal{L}_μ is given by the following (sub)-formula: $\mathbf{E}_\mu[q] = (\langle \bar{s} \rangle \alpha \wedge \langle p \rangle \mathbf{RE}_\mu(e, \beta))$, where $\alpha, \beta \in AP \cup \mathbf{T}$.

The formula states the existence of the encoding of the path pattern (via $\mathbf{E}_\mu[e]$) somewhere in a transition system and thus it is quantified by μ (least fixed point) so as to propagate the sub-formula to the entire transition system. μ encodes a reflexive transitive closure over all the programs and is denoted by $\text{lfp}(X, \mathbf{E}_\mu[q])$ [9]. The variables α and β are encoding as nominals whereas the IRIs in path expressions are encoded as atomic propositions. The encoding, shown in Table 2, makes usage of three functions

- \mathbf{RE}_μ : takes an **cEPPs** expression e and a \mathcal{L}_μ formula φ and builds inductively an encoding based on the structure of the path expression.
- T_μ : takes a path test expression and inductively builds an \mathcal{L}_μ formula. The test positions are translated into transition programs in \mathcal{L}_μ , i.e., $\neg s$ becomes the transition program \bar{s} whereas $\neg p$ and $\neg o$ become p and o programs, respectively.
- Γ : takes two (complex) **cEPPs** expressions and on the basis of the end position of the first and start position of the second returns transition programs expressed in \mathcal{L}_μ for the RDF transition system.

We now prove the correctness of the translation of **cEPPs** into \mathcal{L}_μ formulas as shown in the following Lemma.

Lemma 2. Let q be a **cEPPs** query, for every RDF transition system K_G whose associated RDF graph is G , it holds that $\mathbf{E}[q]^G \neq \emptyset$ iff $\mathbf{E}_\mu(q)^{K_G} \neq \emptyset$.

Table 2. Encoding of EPPs into \mathcal{L}_μ . The function Γ takes two (complex) EPP expressions and on the basis of the end position of the first and start position of the second returns transition programs expressed in \mathcal{L}_μ for the RDF transition system. Note that if $\mathbf{p} = \textcolor{teal}{s}$ in line (14), (resp., $\mathbf{p}_1 = \textcolor{teal}{s}$, $\mathbf{p}_2 = \textcolor{teal}{s}$ in line (9)) then the translation into \mathcal{L}_μ is the transition program \bar{s} .

(1) $\text{RE}_\mu(iri, \varphi)$	$= \begin{cases} \langle p \rangle iri \wedge \langle o \rangle \varphi & \text{if } \varphi \in AP \\ \langle p \rangle iri \wedge \varphi & \text{otherwise} \end{cases}$
(2) $\text{RE}_\mu(\wedge \mathbf{e}, \varphi)$	$= \text{RE}_\mu(\mathbf{e}, \langle \bar{s} \rangle \varphi)$
(3) $\text{RE}_\mu(\mathbf{e}_1 \mid \mathbf{e}_2, \varphi)$	$= \text{RE}_\mu(\mathbf{e}_1, \varphi) \vee \text{RE}_\mu(\mathbf{e}_2, \varphi)$
(4) $\text{RE}_\mu(\mathbf{e}_1 / \mathbf{e}_2, \varphi)$	$= \text{RE}_\mu(\mathbf{e}_1, \Gamma(\mathbf{e}_1, \mathbf{e}_2) \text{RE}_\mu(\mathbf{e}_2, \varphi))$
(5) $\text{RE}_\mu((\mathbf{e}_1 \& \mathbf{e}_2), \varphi)$	$= (\text{RE}_\mu(\mathbf{e}_1, \varphi) \wedge \text{RE}_\mu(\mathbf{e}_2, \varphi))$
(6) $\text{RE}_\mu(\mathbf{e}^?, \varphi)$	$= \text{RE}_\mu(\mathbf{e}, \varphi) \vee \langle \bar{s} \rangle \varphi$
(7) $\text{RE}_\mu(\mathbf{e}^+, \varphi)$	$= \mu X. \text{RE}_\mu(\mathbf{e}, \varphi) \vee \text{RE}_\mu(\mathbf{e}, \Gamma(\mathbf{e}, \mathbf{e})X)$
(8) $\text{RE}_\mu(\mathbf{e}^*, \varphi)$	$= \text{RE}_\mu(\mathbf{e}^+, \varphi) \vee \langle \bar{s} \rangle \varphi$
(9) $\text{RE}(\mathbf{p}_1 \text{ test } \mathbf{p}_2, \varphi)$	$= \langle \mathbf{p}_1 \rangle \top \wedge \text{T}_\mu(\text{test}) \wedge \langle \mathbf{p}_2 \rangle \top \wedge \varphi$
(10) $\text{T}_\mu(iri)$	$= \langle p \rangle iri$
(11) $\text{T}_\mu(\text{test}_1 \&\& \text{test}_2)$	$= \text{T}_\mu(\text{test}_1) \wedge \text{T}_\mu(\text{test}_2)$
(12) $\text{T}_\mu(\text{test}_1 \parallel \text{test}_2)$	$= \text{T}_\mu(\text{test}_1) \vee \text{T}_\mu(\text{test}_2)$
(13) $\text{T}_\mu(!\text{test})$	$= \neg \text{T}_\mu(\text{test})$
(14) $\text{T}_\mu(\text{TP}(\mathbf{p}, \mathbf{e}))$	$= \langle \mathbf{p} \rangle \text{RE}_\mu(\mathbf{e}, \top)$
(15) $\Gamma(\mathbf{e}_1, \mathbf{e}_2)$	$= \langle o \rangle \langle s \rangle$
(16) $\Gamma(\mathbf{e}_1, \wedge \mathbf{e}_2)$	$= \langle o \rangle \langle \bar{o} \rangle$
(17) $\Gamma(\wedge \mathbf{e}_1, \mathbf{e}_2)$	$= \langle o \rangle \langle s \rangle$
(18) $\Gamma(\wedge \mathbf{e}_1, \wedge \mathbf{e}_2)$	$= \langle o \rangle \langle \bar{o} \rangle$

Proof (\Rightarrow). Assume that there exists an RDF G such that the evaluation of q over G is nonempty, i.e., $\mathbf{E} \llbracket q \rrbracket G \neq \emptyset$. We can build such graph from the canonical instance of q and it can be produced using a function θ as shown below:

- if $(\alpha, iri, \beta) \in q$, then $\theta((\alpha, iri, \beta)) = (\alpha, y, \beta) \in G$,
- if $(\alpha, \mathbf{e}, \beta) \in q$, then $\theta((\alpha, \mathbf{e}, \beta)) \in G$,
- if $(\alpha, \wedge \mathbf{e}, \beta) \in q$, then $\theta((\alpha, \wedge \mathbf{e}, \beta)) = (\beta, \mathbf{e}, \alpha) \in G$,
- if $(\alpha, \mathbf{e}_1 / \mathbf{e}_2, \beta) \in q$, then $\theta((\alpha, \mathbf{e}_1, y)) \in G$ and $\theta((y, \mathbf{e}_2, \beta)) \in G$,
- if $(\alpha, \mathbf{e}_1 \& \mathbf{e}_2, \beta) \in q$, then $\theta((\alpha, \mathbf{e}_1, \beta)) \in G$ or $\theta((\alpha, \mathbf{e}_2, \beta)) \in G$,
- if $(\alpha, \mathbf{e}_1 \& \mathbf{e}_2, \beta) \in q$, then $\theta((\alpha, \mathbf{e}_1, \beta)) \in G$ and $\theta((\alpha, \mathbf{e}_2, \beta)) \in G$,
- if $(\alpha, \mathbf{e}^+, \beta) \in q$, then $\bigcup_{i=1}^n \theta((\alpha, \mathbf{e}^i, \beta)) \in G$, where \mathbf{e}^i denotes the composition of \mathbf{e} i times and $n \in \mathbb{N}$,
- if $(\alpha, \mathbf{e}^*, \beta) \in q$, then $(\alpha, u, \alpha) \in G$ or $(\beta, u, \beta) \in G$ or $\theta((\alpha, \mathbf{e}^+, \beta)) \in G$, where u is a fresh IRI,
- if $(\alpha, \mathbf{p}_1 \text{ iri } \mathbf{p}_2, \beta) \in q$, then $(\alpha, iri, \beta) \in G$,
- if $(\alpha, \mathbf{p}_1 \text{ test } \mathbf{p}_2, \beta) \in q$, then $\theta((\alpha, \mathbf{p}_1 \text{ test } \mathbf{p}_2, \beta)) \in G$,

- if $(\alpha, \mathbf{P}_1 \text{ test}_1 \mathbf{P}_2 \ \&\& \ \mathbf{P}_3 \text{ test}_2 \mathbf{P}_4, \beta) \in q$, then $\theta((\alpha, \mathbf{P}_1 \text{ test}_1 \mathbf{P}_2, \beta)) \in G$ and $\theta((\alpha, \mathbf{P}_3 \text{ test}_2 \mathbf{P}_4, \beta)) \in G$,
- if $(\alpha, \mathbf{P}_1 \text{ test}_1 \mathbf{P}_2 \parallel \mathbf{P}_3 \text{ test}_2 \mathbf{P}_4, \beta) \in q$, then $\theta((\alpha, \mathbf{P}_1 \text{ test}_1 \mathbf{P}_2, \beta)) \in G$ or $\theta((\alpha, \mathbf{P}_3 \text{ test}_2 \mathbf{P}_4, \beta)) \in G$,
- if $(\alpha, \mathbf{P}_1 \text{ !test } \mathbf{P}_2, \beta) \in q$, then we introduce a set of fresh IRI's iri_i that do not appear in test such that $(\alpha, iri_i, \beta) \in G$ and $\theta(\alpha, \text{test}, \beta) \notin G$,
- if $(\alpha, \mathbf{P}_1 \text{ TP}(\mathbf{P}, \mathbf{e}) \mathbf{P}_2, \beta) \in q$, then $(\alpha, y, \beta) \in G$ and $\theta(y, \mathbf{e}, r) \in G$ or $\theta(\alpha, \mathbf{e}, r) \in G$ or $\theta(\beta, \mathbf{e}, r) \in G$, where y and r are fresh IRIs.

Since G is an instance of q , the evaluation of q over G is not empty. Now, we construct an RDF transition system $K_G = (S, R, L)$ in the same way as it is done in Definition 9 of [7]. It is possible to verify K_G is a model of $E_\mu(q)$ by working inductively on the construction of $E_\mu(q)$. This is because atomic propositions encoding the constants and distinguished variables are true in K_G as they exist in G ; therefore, $E_\mu(q)$ is satisfiable in K_G . To elaborate, if l is a distinguished variable (i.e., either x or z) or a constant \mathbf{e} in q , then

- for l a distinguished variable, l is satisfiable in K_G since $\llbracket l \rrbracket^{K_G}$ is non empty,
- for $l = \mathbf{e}$ an EPP, its encoding $\text{RE}_\mu(\mathbf{e}, \top)$ is satisfiable in K_G if \mathbf{e} is satisfiable in K_G , this can be proved inductively on the structure of the encoding.

Thus, since K_G is an RDF transition system, we obtain that $\llbracket E_\mu(q) \rrbracket^{K_G} \neq \emptyset$.

(\Leftarrow) Assume that $\llbracket E_\mu(q) \rrbracket^{K_G} \neq \emptyset$. In order to test if the evaluation of q over an RDF graph obtained from K_G is non empty, we produce an RDF graph G from K_G as done in Lemma 4 of [7]. Thus, since G is a technical construction obtained from an RDF transition system associated to q , it holds that $\mathbf{E}[q]^G \neq \emptyset$.

Example 3. The μ -calculus encoding of the EPPs in Example 1 is given below. $\Gamma(\text{et1}, \text{et2})$ is a shorthand for $\Gamma(\text{:leaderParty} \ \&\& \ \text{TP}(\text{_, :formationYear}), \text{^{:leaderParty}})$.

$$\begin{aligned}
E_\mu(\mathbf{e1} \ \&\ \mathbf{e2}) &= \text{Ifp}(X, \langle \bar{s} \rangle \top \wedge \text{RE}_\mu(\mathbf{e1} \ \&\ \mathbf{e2}, \top)) \\
\text{RE}_\mu(\mathbf{e1} \ \&\ \mathbf{e2}, \top) &= \text{RE}_\mu(\mathbf{e1}, \top) \wedge \text{RE}_\mu(\mathbf{e2}, \top) \\
\text{RE}_\mu(\mathbf{e1}, \top) &= \text{RE}_\mu(\text{:country} / \text{^{:country}} \sim \text{:region} / \text{^{:region}}}, \top) \\
&= \text{RE}_\mu(\text{:country} / \text{^{:country}}, \top) \wedge \neg \text{RE}_\mu(\text{:region} / \text{^{:region}}, \top) \\
&= \text{blue}(\text{RE}_\mu(\text{:country}, \Gamma(\text{:country}, \text{^{:country}}) \text{RE}_\mu(\text{^{:country}}, \top)) \text{blue}) \wedge \\
&\quad \neg \text{blue}(\text{RE}_\mu(\text{:region}, \Gamma(\text{:region}, \text{^{:region}}) \text{RE}_\mu(\text{^{:region}}, \top)) \text{blue}) \\
&= \text{blue}(\langle p \rangle \text{:country} \wedge \langle o \rangle \langle \bar{o} \rangle (\langle p \rangle \text{:country} \wedge \top) \text{blue}) \wedge \neg \text{blue}(\langle p \rangle \text{:region} \wedge \langle o \rangle \langle \bar{o} \rangle (\langle p \rangle \text{:region} \wedge \langle \bar{s} \rangle \top) \text{blue}) \\
\text{RE}_\mu(\mathbf{e2}, \top) &= \text{RE}_\mu(\text{blue}(\text{:leaderParty} \ \&\& \ \text{TP}(\text{_, :formationYear}) \text{blue}) / \text{^{:leaderParty}}, \top) \\
&= \text{RE}_\mu(\text{:leaderParty} \ \&\& \ \text{TP}(\text{_, :formationYear}), \Gamma(\text{et1}, \text{et2}) \text{RE}_\mu(\text{^{:leaderParty}}, \top)) \\
&= \text{RE}_\mu(\text{:leaderParty}, \Gamma(\text{et1}, \text{et2}) \text{RE}_\mu(\text{^{:leaderParty}}, \top)) \wedge \\
&\quad \text{RE}_\mu(\text{TP}(\text{_, :formationYear}), \Gamma(\text{et1}, \text{et2}) \text{RE}_\mu(\text{^{:leaderParty}}, \top)) \\
&= \text{blue}(\langle p \rangle \text{:leaderParty} \wedge \langle o \rangle \langle \bar{o} \rangle (\langle p \rangle \text{:leaderParty} \wedge \top) \text{blue}) \wedge \\
&\quad \text{RE}_\mu(\text{TP}(\text{_, :formationYear}), \langle o \rangle \langle \bar{o} \rangle (\langle p \rangle \text{:leaderParty} \wedge \top)) \\
&= \text{blue}(\langle p \rangle \text{:leaderParty} \wedge \langle o \rangle \langle \bar{o} \rangle (\langle p \rangle \text{:leaderParty} \wedge \top) \text{blue}) \wedge \\
&\quad (\langle \bar{s} \rangle \top \wedge \text{T}_\mu(\text{TP}(\text{_, :formationYear})) \wedge \langle o \rangle \top \wedge \langle o \rangle \langle \bar{o} \rangle (\langle p \rangle \text{:leaderParty} \wedge \top)) \\
&= \text{blue}(\langle p \rangle \text{:leaderParty} \wedge \langle o \rangle \langle \bar{o} \rangle (\langle p \rangle \text{:leaderParty} \wedge \top) \text{blue}) \wedge \\
&\quad (\langle \bar{s} \rangle \top \wedge \langle o \rangle \text{:formationYear} \wedge \langle p \rangle \top \wedge \langle o \rangle \top \wedge \langle o \rangle \langle \bar{o} \rangle (\langle p \rangle \text{:leaderParty} \wedge \top))
\end{aligned}$$

So far, we have shown that our encoding of **cEPPs** queries into the μ -calculus is correct. Next, we formally present the reduction of the containment test of **cEPPs** into unsatisfiability test in the μ -calculus.

5 Containment Results

In this section we state our complexity results. We start with **cEPPs** and then move to SPARQL with **cEPPs**. We start with the containment of **cEPPs** which can be stated as follows:

Problem: cEPPs Containment

Input: **cEPPs** q_1 and q_2

Output: Is $q_1 \sqsubseteq q_2$?

Given two **cEPPs** queries q_1 and q_2 , $q_1 \sqsubseteq q_2$ if and only if for any RDF graph G , the answers of q_1 are included in the answers of q_2 , i.e., $\mathbf{E}[q_1]^G \subseteq \mathbf{E}[q_2]^G$.

In the previous section, we have presented various functions to produce \mathcal{L}_μ formulas corresponding to the encodings of **cEPPs**. Hence, the problem of containment can be reduced to formula unsatisfiability test in \mathcal{L}_μ as:

Problem: cEPPs Containment via \mathcal{L}_μ

Input: \mathcal{L}_μ encodings $E_\mu(q_1)$ and $E_\mu(q_2)$

Output: Is $E_\mu(q_1) \wedge \neg E_\mu(q_2)$ unsatisfiable?

Therefore, we obtain that $q_1 \sqsubseteq q_2 \Leftrightarrow E_\mu(q_1) \wedge \neg E_\mu(q_2)$ is unsatisfiable. We now show that our decision procedure is sound and complete.

Theorem 1 (Soundness). *Given two cEPPs queries q_1 and q_2 if $E_\mu(q_1) \wedge \neg E_\mu(q_2)$ is unsatisfiable, then $q_1 \sqsubseteq q_2$.*

Proof. We show the contrapositive. If $q_1 \not\sqsubseteq q_2$, then $E_\mu(q_1) \wedge \neg E_\mu(q_2)$ is satisfiable. One can verify that every RDF graph G in which there is at least one tuple satisfying q_1 but not q_2 can be turned into a RDF transition system model for $E_\mu(q_1) \wedge \neg E_\mu(q_2)$. To do so, consider an RDF graph G and assume that there exists a mapping $m \in \mathbf{E}[q_1]^G$ and $m \notin \mathbf{E}[q_2]^G$. Let us construct an RDF transition system K_G from G . In fact, we produce G which contains at least a canonical instantiation of q_1 , i.e., by replacing the variables in q_1 with their mappings in m . It has been shown that an RDF graph can be turned into an RDF transition system (cf. Lemma 1). At this point, it remains to verify that $\llbracket E_\mu(q_1) \rrbracket^{K_G} \neq \emptyset$ and $\llbracket E_\mu(q_2) \rrbracket^{K_G} = \emptyset$.

Let us construct the formulas $E_\mu(q_1)$ and $E_\mu(q_2)$ by first skolemizing the distinguished variables using their mappings in m . Consequently, from Lemma 2 one obtains, $\llbracket E_\mu(q_1) \rrbracket^{K_G} \neq \emptyset$. However, $\llbracket E_\mu(q_2) \rrbracket^{K_G} = \emptyset$, this is because the atomic propositions in the formula corresponding to the constants and variables are not satisfied in K . This implies that $\llbracket \neg E_\mu(q_2) \rrbracket^{K_G} \neq \emptyset$. This is justified by the fact that if a formula φ is satisfiable in an RDF transition system, then $\llbracket \varphi \rrbracket^{K_G} = S$ thus $\llbracket \neg \varphi \rrbracket^{K_G} = \emptyset$. So far we have: $\llbracket E_\mu(q_1) \rrbracket^{K_G} \neq \emptyset$ and $\llbracket \neg E_\mu(q_2) \rrbracket^{K_G} \neq \emptyset$. Without loss of generality, $\llbracket E_\mu(q_1) \wedge \neg E_\mu(q_2) \rrbracket^{K_G} \neq \emptyset$. Therefore, $E_\mu(q_1) \wedge \neg E_\mu(q_2)$ is satisfiable.

Theorem 2 (Completeness). *Given two cEPPs queries q_1 and q_2 , if $E_\mu(q_1) \wedge \neg E_\mu(q_2)$ is satisfiable, then $q_1 \not\sqsubseteq q_2$.*

Proof. $E_\mu(q_1) \wedge \neg E_\mu(q_2)$ is satisfiable $\Rightarrow \exists K_G. \llbracket E_\mu(q_1) \wedge \neg E_\mu(q_2) \rrbracket^{K_G} \neq \emptyset$. Consequently, K_G is an RDF transition system as shown in Proposition 1 of [7].

Using $K_G = (S' \cup S'', R, L)$, we construct an RDF graph G such that we can show that $\mathbf{E} \llbracket q_1 \rrbracket^G \not\subseteq \mathbf{E} \llbracket q_2 \rrbracket^G$ and hence $q_1 \not\sqsubseteq q_2$ holds:

- for each `iri` or constant u in q_1 and q_2 , we create an RDF triple (α, u, β) such that $\{(\alpha, \beta) \mid \exists n \in \llbracket u \rrbracket^{K_G} \wedge n' \in S'' \wedge (\alpha, n') \in R(s) \wedge (n', n) \in R(p) \wedge (n', \beta) \in R(o)\}$.

Thus, it remains to show that $\mathbf{E} \llbracket q_1 \rrbracket^G \not\subseteq \mathbf{E} \llbracket q_2 \rrbracket^G$. From our assumption, we obtain the following:

$$\begin{aligned} \llbracket E_\mu(q_1) \wedge \neg E_\mu(q_2) \rrbracket^{K_G} \neq \emptyset &\Rightarrow \llbracket E_\mu(q_1) \rrbracket^{K_G} \neq \emptyset \text{ and } \llbracket \neg E_\mu(q_2) \rrbracket^{K_G} \neq \emptyset. \\ &\Rightarrow \llbracket E_\mu(q_1) \rrbracket^{K_G} \neq \emptyset \text{ and } \llbracket E_\mu(q_2) \rrbracket^{K_G} = \emptyset. \end{aligned}$$

Note here that, if a formula φ is satisfiable in an RDF transition system K , then $\llbracket \varphi \rrbracket^{K_G} = S$ as shown in [7]. Consequently, using Lemma 2 and G , we obtain $\mathbf{E} \llbracket q_1 \rrbracket^G \neq \emptyset$ and $\mathbf{E} \llbracket q_2 \rrbracket^G = \emptyset$ because G contains all those triples that satisfy q_1 and not q_2 . Therefore, we get $\mathbf{E} \llbracket q_1 \rrbracket^G \not\subseteq \mathbf{E} \llbracket q_2 \rrbracket^G$ and thus $q_1 \not\sqsubseteq q_2$.

From Theorems 1 and 2, we obtain the following result.

Proposition 1 (Complexity of *cEPPs* Containment). *The containment of *cEPPs* queries can be determined in a double exponential amount time.*

Note that the size of the encoding n is exponential. Hence, we obtain a 2EXPTIME upper bound for containment. If we remove path conjunction from the *cEPPs* we obtain an EXPTIME upper bound; this is due to the complexity of satisfiability in the μ -calculus which is $2^{\mathcal{O}(n^2 \log n)}$ where n is the size of the encoding (which is exponential for *cEPPs*). To provide a lower bound to the problem of containment in *cEPPs*, we utilize the close connection between *cEPPs* and PDL (Propositional Dynamic Logic) with path intersection and converse [17]. Like \mathcal{L}_μ , PDL formulas are interpreted over transition systems (aka. Kripke structures). *cEPPs* without tests are exactly the same as that of PDL navigational programs (which are regular expressions). PDL formulas can be interpreted over RDF transitions systems by using a fixed set of programs $\{s, p, o, \bar{s}, \bar{p}, \bar{o}\}$ as we have done for \mathcal{L}_μ . Likewise, *cEPPs* expressions can be evaluated over RDF transition systems by a simple rewriting. Given a *cEPPs* expression e that can be expressed over a standard RDF graph, it can be turned into an expression e' that can be expressed over an RDF transition system by a bijective rewriting function ϑ . ϑ is defined inductively:

$$\begin{aligned} \vartheta(iri) &= p/\tau(\underline{o} = iri)/\wedge p/o & \vartheta(\wedge e) &= p/\tau(\underline{o} = iri)/\wedge p/\wedge s \\ \vartheta(e_1/e_2) &= \vartheta(e_1)/pos(e_1, e_2)/\vartheta(e_2) & \vartheta((e)^*) &= (\vartheta(e))^* \\ \vartheta((e)^+) &= (\vartheta(e))^+ & \vartheta(e_1 \mid e_2) &= \vartheta(e_1) \mid \vartheta(e_2) \\ \vartheta(e_1 \& e_2) &= \vartheta(e_1) \& \vartheta(e_2) & \vartheta(p_1 \text{ test } p_2) &= \vartheta_\tau(p_1)/\vartheta_\tau(\text{test})/\vartheta_\tau(p_2) \\ \vartheta_\tau(iri) &= p/\tau(\underline{o} = iri)/\wedge p & \vartheta_\tau(\text{test}_1 \&\& \text{test}_2) &= \vartheta_\tau(\text{test}_1) \&\& \vartheta_\tau(\text{test}_2) \\ \vartheta_\tau(\text{test}_1 \parallel \text{test}_2) &= \vartheta_\tau(\text{test}_1) \parallel \vartheta_\tau(\text{test}_2) & \vartheta_\tau(\text{TP}(p, e)) &= \vartheta_\tau(p)/\vartheta(e) \\ \vartheta_\tau(!\text{test}) &= !\vartheta_\tau(\text{test}) & \vartheta_\tau(\underline{s}) &= s \quad \vartheta_\tau(\underline{p}) = p \quad \vartheta_\tau(\underline{o}) = o \\ pos(e_1, e_2) &= s \quad pos(e_1, \wedge e_2) = \wedge o & pos(\wedge e_1, e_2) &= s \quad pos(\wedge e_1, \wedge e_2) = \wedge o \end{aligned}$$

where s , p and o are navigational programs of an RDF transition system. In order to show that **cEPPs** tests are just a syntactic variant of PDL node formulas (aka XPath node expressions), we prove the following.

Lemma 3. *Satisfiability of **cEPPs** queries on RDF transition systems can be reduced in polynomial time to satisfiability of **cEPPs** queries on RDF graphs.*

Proof. We can turn an RDF transition system into a standard RDF graph by using a bijective mapping function similar to the one in [7]. It is also possible to transform a **cEPPs** expression e over an RDF transition system K_G into a **cEPPs** expression e' that can be expressed over a standard RDF graph G via the function ϑ . Thus, it remains to show that, e is satisfiable in K_G iff e' is satisfiable in G . This can be done inductively on the construction of the expressions.

Theorem 3. *Testing containment of **cEPPs** queries is 2EXPTIME-hard.*

The above result is obtained by examining the connection between **cEPPs** and Propositional Dynamic Logic with path intersection (ICPDL). **cEPPs** can be as seen as a notational variant of ICPDL with the fixed set of atomic programs $\{s, p, o, \bar{s}, \bar{p}, \bar{o}\}$. For example, the **cEPPs** expression $\vartheta(a/\text{TP}(\underline{s}, b^+)/c) = s/a/\dots$ corresponds to the ICPDL formula $a \circ \langle s \rangle \langle b^+ \rangle \circ \langle c \rangle$.

The double exponential upper bound is unavoidable as it has already been shown that path conjunction makes the containment problem very difficult for XPath; it has been proved that the complexity of the containment of CoreXPath(\ast , \cap) expressions is 2EXPTIME-complete [30]. Consider now the language of **cEPPs** without path intersection (**cEPPs_{nn}**). This language is closely related to PDL and CoreXPath [30].

Theorem 4. *The containment problem of **cEPPs_{nn}** is EXPTIME-complete.*

The following result follows from the containment problem for GXPath_{reg} (Graph XPath) [21] and the satisfiability problem for PDL with negation on paths [18].

Theorem 5. *The containment problem of **EPPs** is undecidable.*

5.1 Containment of SPARQL with EPPs

In this section we briefly discuss the containment problem for SPARQL queries *with EPPs* (from here onwards, we refer to it as just SPARQL). We restrict ourselves to the union of conjunctive fragment of SPARQL, i.e., *AND-UNION fragment of SPARQL, without projection*.

Definition 6. *A SPARQL query q is defined inductively as: $q ::= (\bigvee \cup I, e, \bigvee \cup I \cup L) \mid q \text{ AND } q' \mid q \text{ UNION } q'$.*

The problem of the containment of SPARQL queries can be reduced into unsatisfiability test in \mathcal{L}_μ as given below:

Problem: SPARQL Containment via \mathcal{L}_μ Input: \mathcal{L}_μ encodings $E(q_1)$ and $E(q_2)$ Output: Is $E(q_1) \wedge \neg E(q_2)$ unsatisfiable?

In the following, we extend the EPPs encoding function E_μ to translate the containment test of q_1 in q_2 into the μ -calculus.

Definition 7 (SPARQL containment). *The encoding of a SPARQL query q is obtained inductively as follows:*

$$\begin{aligned} E_\mu((\alpha, \mathbf{e}, \beta)) &= \text{Ifp}(X, \langle \bar{s} \rangle \alpha \wedge \langle p \rangle \text{RE}_\mu(\mathbf{e}, \beta)) \\ E_\mu(q \text{ AND } q') &= E_\mu(q) \wedge E_\mu(q') \\ E_\mu(q \text{ UNION } q') &= E_\mu(q) \vee E_\mu(q') \end{aligned}$$

The variables on the left hand side query are encoded into nominals and the IRIs in the path expressions are encoded into atomic propositions; to encode variables that appear on the right hand side query, we follow two steps: (1) if a variable appears uniquely in the query, it is encoded into \top , and (2) if a variable appears multiple times, it is encoded by using the IRIs in the path expression (the triple pattern in which it appears). This technique has been already used in [9]. For instance, if the query is $q = (x, \mathbf{e}_1, y) \text{ AND } (y, \mathbf{e}_2, z)$, then the encoding of x (resp. z) is \top and $y \rightarrow \langle \bar{o} \rangle \langle p \rangle \mathbf{e}_1$ or $y \rightarrow \langle s \rangle \langle p \rangle \mathbf{e}_2$. Thus, the encoding becomes:

$$E_\mu(qp) = E_\mu((\top, \mathbf{e}_1, \langle \bar{o} \rangle \langle p \rangle \mathbf{e}_1) \text{ AND } (\langle \bar{o} \rangle \langle p \rangle \mathbf{e}_1, \mathbf{e}_2, \top)) \vee$$

$$E_\mu((\top, \mathbf{e}_1, \langle s \rangle \langle p \rangle \mathbf{e}_2) \text{ AND } (\langle s \rangle \langle p \rangle \mathbf{e}_2, \mathbf{e}_2, \top))$$

The size of the encoding for the containment problem is polynomial in the number of variables that appear more than once. To be more precise, for a given query, the size of the encoding is: $\Pi_{x \in \text{multiVar}} |x|$ where multiVar is the set of variables occurring more than once, and $|x|$ is the number of times variable x appears in the query. Note that the size of the encoding is linear if all the variables on the right-hand side query appear uniquely.

Theorem 6 (Soundness and Completeness). *Given two SPARQL queries q_1 and q_2 , $E_\mu(q_1) \wedge \neg E_\mu(q_2)$ is unsatisfiable if and only if $q_1 \sqsubseteq q_2$.*

From the above theorem, we obtain the following result:

Proposition 2. *Given two SPARQL queries q_1 and q_2 , the containment test can be solved in a double exponential amount of time in the size of the encoding $|E_\mu(q_1)| + |E_\mu(q_2)|$.*

This result is not surprising, as we have shown in Proposition 1 that the complexity of the containment of EPPs is 2EXPTIME in the worst case. Furthermore, if projection was part of the SPARQL fragment we consider here, then there is a further jump in the complexity of containment, i.e., the complexity increases by one exponential. Thereby, we obtain 3EXPTIME upper bound for the containment of SPARQL queries with projection. This is due to an exponential blow up in the size of the encoding as one needs to take care of multiple occurring non-distinguished variables.

6 Conclusions and Future Work

We have discussed a preliminary study of the problem of query containment for an expressive class of navigational queries captured by the EPPs language. Our study leverages μ -calculus to encode EPPs and then use this encoding to get a 2EXPTIME complexity upper bound. This bound remains the same when considering EPPs in SPARQL without projection. However, if one considers projection, there is an exponential increase in the complexity which results in a 3EXPTIME bound. While the results obtained are of theoretical interest, from a practical point of view an implementation is available (<http://sparql-qc-bench.inrialpes.fr/>) which can be extended for EPPs. Furthermore, an additional benefit of using μ -calculus is that by exploiting *logical combinators* the size of the encoding can be reduced by upto exponentiation. Thus, the complexity bounds that we obtained are not prohibitive in terms of a practical implementation.

A natural line of future research is to provide a tighter complexity bound for the problem of containment of EPPs. Moreover, the investigation of how the inclusion of constraints in EPPs affects the complexity of query containment is in our research agenda.

References

1. Barcelo, P., Libkin, L., Lin, A.W., Wood, P.T.: Expressive languages for path queries over graph-structured data. *ACM Trans. Database Syst. (TODS)* **37**, 31 (2012)
2. Calvanese, D., De Giacomo, G., Lenzerini, M., Vardi, M.Y.: Containment of conjunctive regular path queries with inverse. In: *Proceedings of 7th International Conference on the Principles of Knowledge Representation and Reasoning (KR 2000)*, pp. 176–185 (2000)
3. Calvanese, D., Giacomo, G.D., Lenzerini, M.: Conjunctive query containment and answering under description logic constraints. *ACM Trans. Comput. Logic (TOCL)* **9**(3), 22 (2008)
4. Calvanese, D., Ortiz, M., Simkus, M.: Containment of regular path queries under description logic constraints. In: *IJCAI Proceedings-International Joint Conference on Artificial Intelligence*, vol. 22, p. 805 (2011)
5. Chandra, A.K., Merlin, P.M.: Optimal implementation of conjunctive queries in relational data bases. In: *Proceedings of the Ninth Annual ACM Symposium on Theory of Computing, STOC 1977*, pp. 77–90. ACM, New York (1977)
6. Chaudhuri, S., Vardi, M.Y.: Optimization of real conjunctive queries. In: *Proceedings of the Twelfth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems, PODS 1993*, pp. 59–70. ACM, New York (1993)
7. Chekol, M.W.: Static analysis of semantic web queries. Ph.D. thesis, Université de Grenoble (2012)
8. Chekol, M.W., Euzenat, J., Genevès, P., Layaïda, N.: PPARQL query containment. In: *Proceedings of DBPL (2011)*
9. Chekol, M.W., Euzenat, J., Genevès, P., Layaïda, N.: SPARQL query containment under RDFS entailment regime. In: Gramlich, B., Miller, D., Sattler, U. (eds.) *IJCAR 2012. LNCS*, vol. 7364, pp. 134–148. Springer, Heidelberg (2012)

10. Chekol, M.W., Euzenat, J., Genevès, P., Layaïda, N.: SPARQL query containment under SHI axioms. In: *Proceedings of AAAI*, pp. 10–16 (2012)
11. Fionda, V., Pirrò, G.: Querying graphs with preferences. In: *22nd ACM International Conference on Information and Knowledge Management, CIKM 2013, San Francisco, CA, USA, October 27– November 1*, pp. 929–938 (2013)
12. Fionda, V., Pirrò, G., Consens, M.P.: Extended property paths: writing more SPARQL queries in a succinct way. In: *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence, 25–30 January, Austin, Texas, USA*, pp. 102–108 (2015)
13. Fionda, V., Pirrò, G., Gutierrez, C.: NautiLOD: a formal language for the web of data graph. *Trans. Web* **9**(1), 5: 1–5: 43 (2015)
14. Genevès, P., Layaïda, N.: A system for the static analysis of XPath. *ACM Trans. Inf. Syst.* **24**(4), 475–502 (2006)
15. Genevès, P., Layaïda, N., Schmitt, A.: Efficient static analysis of XML paths and types. In: *Proceedings of PLDI*, pp. 342–351. ACM (2007)
16. Genevès, P., Schmitt, A. : Expressive logical combinators for free. In: *IJCAI*, pp. 311–317 (2015)
17. Göller, S., Lohrey, M., Lutz, C.: PDL with intersection and converse: satisfiability and infinite-state model checking. *J. Symb. Logic* **74**(01), 279–314 (2009)
18. Harel, D., Kozen, D., Tiuryn, J.: *Dynamic Logic*. MIT Press, Cambridge (2000)
19. Ioannidis, Y.E., Ramakrishnan, R.: Containment of conjunctive queries: beyond relations as sets. *ACM Trans. Database Syst.* **20**(3), 288–324 (1995)
20. Kostylev, E.V., Reutter, J.L., Romero, M., Vrgoč, D.: SPARQL with property paths. In: Arenas, M., et al. (eds.) *ISWC 2015. LNCS*, vol. 9366, pp. 3–18. Springer, Heidelberg (2015)
21. Kostylev, E.V., Reutter, J.L., Vrgoc, D.: Containment of data graph queries. In: *ICDT*, pp. 131–142 (2014)
22. Kostylev, E.V., Reutter, J.L., Vrgoč, D.: Static analysis of navigational XPath over graph databases. *Inf. Process. Lett.* **116**(7), 467–474 (2016)
23. Kozen, D.: Results on the propositional μ -calculus. *Theoret. Comput. Sci.* **27**(3), 333–354 (1983)
24. Krötzsch, M., Rudolph, S.: Conjunctive queries for EL with composition of roles. In: *Proceedings of the 2007 International Workshop on Description Logics (DL 2007)*, Brixen-Bressanone, near Bozen-Bolzano, 8–10 Italy 2007, June 2007
25. Pérez, J., Arenas, M., Gutierrez, C.: Semantics and complexity of SPARQL. *ACM Trans. Database Syst. (TODS)* **34**(3), 16 (2009)
26. Pérez, J., Arenas, M., Gutierrez, C.: nSPARQL: a navigational language for RDF. *Web Semant. Sci. Serv. Agents WWW* **8**(4), 255–270 (2010)
27. Pichler, R., Skritek, S.: Containment and Equivalence of Well-designed SPARQL, pp. 39–50 (2014)
28. Reutter, J.L.: Containment of nested regular expressions. *CoRR*, abs/1304.2637 (2013)
29. Sagiv, Y., Yannakakis, M.: Equivalences among relational expressions with the union and difference operators. *J. ACM* **27**(4), 633–655 (1980)
30. ten Cate, B., Lutz, C.: The complexity of query containment in expressive fragments of XPath 2.0. *J. ACM (JACM)* **56**(6), 31 (2009)
31. Vardi, M.Y., Murano, A., Lutz, C., Bonatti, P.A.: The complexity of enriched μ -calculi. *Log. Methods Comput. Sci.* **4**, 1–27 (2008)