

Extended Property Paths: Writing More SPARQL Queries in a Succinct Way*

Valeria Fionda¹, Giuseppe Pirrò² and Mariano P. Consens³

¹ DeMaCS, University of Calabria, Italy

² WeST, University of Koblenz-Landau, Germany

³ MIE, University of Toronto, Canada

Abstract

We introduce Extended Property Paths (EPPs), a significant enhancement of SPARQL property paths. EPPs allow to capture in a succinct way a larger class of *navigational queries* than property paths. We present the syntax and formal semantics of EPPs and introduce two different evaluation strategies. The first is based on an algorithm implemented in a custom query processor. The second strategy leverages a translation algorithm of EPPs into SPARQL queries that can be executed on existing SPARQL processors. We compare the two evaluation strategies on real data to highlight their pros and cons.

1 Introduction

Motivated by the spread of graph-like data, the area of graph databases has received renewed attention. SPARQL, the W3C standard query language for data in RDF, has recently been enhanced with *property paths* (Harris and Seaborne 2013) to support graph navigation capabilities. The regular-expression-like syntax of property paths enables to write SPARQL navigational queries in a more *succinct* way and extend the matching of triple patterns to arbitrary length paths.

However, property paths offer limited navigation capabilities; queries like “*Find my exclusive friends*”, that is, my friends that are not friend of any of my friends, cannot be expressed. To enhance the expressive power of property paths and enable to write more navigational queries in a succinct way, we introduce Extended Property Paths (EPPs). In designing EPPs we identified a core of new features and investigated how to make available such features; indeed, one can devise a custom query processor like in (Alkhateeb, Baget, and Euzenat 2009; Fionda and Pirrò 2013) and/or leverage existing (SPARQL) processors. The goal of this paper is to discuss the design and implementation of EPPs.

Related Work. Graph query languages have been deeply studied (Wood 2012). It emerged that for certain classes

of languages, like Conjunctive Regular Path Queries (CRPQs) (Barceló et al. 2012; Alkhateeb, Baget, and Euzenat 2009), the evaluation problem can become too expensive, making the class unattractive for practical purposes. Hence, restricted classes of languages like acyclic CRPQs and Nested Regular Expressions (NREs) (Pérez, Arenas, and Gutierrez 2010) have been proposed. NREs are at the core of navigational languages for RDF such as nSPARQL (Pérez, Arenas, and Gutierrez 2010) and NautiLOD (Fionda, Gutierrez, and Pirrò 2012). Languages like `TriAL` (Libkin, Reutter, and Vrgoč 2013), `TriQ` (Arenas, Gottlob, and Pieris 2014) and NEMODEQ (Rudolph and Krötzsch 2013) are grounded on (extensions of) Datalog. Other expressive languages like SPARQL`LeR` (Kochut and Janik 2007) and *extended* CRPQs (Barceló et al. 2012) focus on discovering paths.

In terms of expressiveness, Barceló et al. (Barceló, Pérez, and Reutter 2012) showed that NREs *cannot express* queries involving, for instance, path conjunction. SPARQL property paths are even less expressive than NREs because of the lack of nesting and tests also within (arbitrary length) paths. `TriAL` can capture queries not captured by NREs and nSPARQL. However, NREs/nSPARQL queries can be evaluated in linear time. Arenas et al. (Arenas, Gottlob, and Pieris 2014) studied the relationships between `TriQ`, `TriAL` and NEMODEQ in terms of Datalog[±] programs (Cali et al. 2010). Fletcher et al. (Fletcher et al. 2011) compared the expressiveness of navigational languages by considering different sets of features. We compare EPPs with closely-related languages in Section 3.3.

EPPs extend the expressive power of SPARQL property paths and NREs-based languages with new features such as *path conjunction*, *path negation* and *powerful types of tests*. These features are partially available in the W3C standard language XPath 2.0 used to query tree-like data (Berglund et al. 2010). From a concrete point of view, with EPPs we faced two main challenges: (i) how to make EPPs a conservative extension of SPARQL property paths; (ii) how to make EPPs readily available in existing SPARQL processors.

1.1 EPPs by Example

We now give an overview of EPPs via a concrete example. The syntax and semantics of EPPs will be discussed in detail

*V. Fionda’s work was supported by the European Commission, the European Social Fund and the Calabria region. G. Pirrò was partially supported by the EU Framework Programme for Research and Innovation under grant agreement no. 611242 (Sense4us).

Copyright © 2015, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.


```

epp ::= ‘^’ epp | epp ‘+’ | epp ‘?’ | epp ‘*’ | epp ‘/’ epp | epp ‘|’ epp | ‘(’ epp ‘)’ |
      [pos]1 test [pos]2 | epp ‘&&’ epp | epp ‘~’ epp
test ::= ‘!’ test | test ‘&&’ test | test ‘||’ test | ‘(’ test ‘)’ | base
base ::= uri | ‘TP’(‘pos’, epp ‘)’ | ‘T’(‘EExp’)
EBuiltInCall ::= BuiltInCall | pos
pos ::= ‘.s’ | ‘.p’ | ‘.o’

```

Table 1: Syntax of EPPs expressions. ¹If omitted is `.s`; ²If omitted is `.o`.

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

Table 2: Formal semantics of EPPs expressions.

a combination of forward/reverse predicates in a (negated) set are expressible in the EPPs syntax via the productions `test` \rightarrow `base` \rightarrow `u`. Note that EPPs use the symbol ‘||’ while PPs use ‘|’.

New features introduced by Extended Property Paths

The following *additional* features are introduced: path conjunction (`epp1&epp2`), path negation (`epp1~epp2`) and different types of tests (`test`) *within* a path, also by specifying the *starting* and *ending* positions (`pos`). EPPs enable to test from each of the subject, predicate and object positions in RDF triples, mapped in the syntax to `.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 a triple.

Tests (`test`) can be of different types and can be combined by using the logical operators AND (&&), OR (||) and NOT (!). A test can be a simple check for the existence of a URI or a *nested* EPP, i.e., `TP(pos, epp)`, which corresponds to the evaluation of `epp` starting from a position `pos` (of the last triple traversed) and whose evaluation returns true if, and only if, there exists at least one node that can be reached via `epp`. A base test (production `base`) can be of type `T`, which is a SPARQL boolean expression; here, `EExp` (not reported here for sake of space) extends the production [110] in the SPARQL grammar² where `BuiltInCall`³ is substituted with `EBuiltInCall`, which enables to use in EPPs tests available in SPARQL as built-in conditions also augmented with positions (`pos`). Built-in condi-

tions are constructed using elements of the set $\mathbf{U} \cup \mathbf{L}$ and constants, logical connectives (\neg, \wedge, \vee), (in)equality symbol(s) ($=, <, >, \leq, \geq$), unary (e.g., `isURI`) and binary (e.g., `STRSTARTS`) functions.

2.2 Extended Property Paths Semantics

The semantics (shown in Table 2) for the interpretation of an EPP expression `epp` on a graph \mathcal{G} uses two functions: (i) $\mathbf{E}[\text{epp}]^{\mathcal{G}}$ defined as a binary relation (u, v) such that u and v are nodes in \mathcal{G} and v is reachable from u via a path in \mathcal{G} satisfying `epp`; and (ii) $\mathbf{E}_T[\text{test}]_t^{\mathcal{G}}$ defined as a boolean function, which evaluates `true` if, and only if, the triple t satisfies the test `test`. The semantics also uses the position mapping function P_m defined as follows:

Definition 2 (Position mapping function). Let $t = \langle x, y, z \rangle$ be a triple pattern, $\{x, y, z\} \subseteq \mathcal{T} \cup \mathcal{V}$. The position mapping function $P_m(\text{pos}, t)$ is defined as: (i) $P_m(\text{.s}, t) = x$, (ii) $P_m(\text{.p}, t) = y$ (iii) $P_m(\text{.o}, t) = z$.

P_m selects one among the subject, predicate and object of a triple on the basis of the value of `pos`. Consider the triple $\langle u_1, p_1, u_2 \rangle$ and the test `T(.p=p1)`; this instantiates P_m as $P_m(\text{.p}, \langle u_1, p_1, u_2 \rangle) = p_1$, which checks $p_1 = p_1$ returning true while testing `T(.o=u3)` gives false.

3 Algorithms and Complexity

This section presents two strategies for the evaluation of EPPs expressions. The first via an ad-hoc algorithm; the second one via a translation into SPARQL queries that can be executed on existing processors.

²<http://www.w3.org/TR/sparql11-query/#rExpression>

³<http://www.w3.org/TR/sparql11-query/#rBuiltInCall>

```

Function EVALUATE( $n, \text{epp}, \mathcal{G}$ )
  Input: node  $n$ , expression  $\text{epp}$ , graph  $\mathcal{G}$ ; Output: node set  $Res$ .
1: if  $\text{epp} = (\text{epp}_1)^*$  then
2:   return CLOSURE( $n, \text{epp}_1, \mathcal{G}, \{\}, 0$ )
3: else if  $\text{epp} = (\text{epp}_1)^+$  then
4:   return CLOSURE( $n, \text{epp}_1, \mathcal{G}, \{\}, 1$ )
5: else
6:   return BASE( $n, \text{epp}, \mathcal{G}$ )

Function CLOSURE( $n, \text{epp}, \mathcal{G}, Res, l$ )
  Input: node  $n$ , EPPs expression  $\text{epp}$ , graph  $\mathcal{G}$ , node set  $Res$ , lower bound  $l$ ;
  Output: node set  $Res$ .
1: if  $l = 1$  then
2:    $S = \text{EVALUATE}(n, \text{epp}, \mathcal{G})$ 
3: else
4:    $S = \{n\}$ 
5: while  $S \neq \emptyset$  do
6:    $n = \text{extractNode}(S)$  /* extract one node */
7:   if  $n \notin Res$  then
8:      $Res = Res \cup \{n\}$ 
9:      $S = S \cup \text{EVALUATE}(n, \text{epp}, \mathcal{G})$ 
10: return  $Res$ 

Function BASE( $n, \text{epp}, \mathcal{G}$ )
  Input: node  $n$ , EPPs expression  $\text{epp}$ , graph  $\mathcal{G}$ ; Output: node set  $Res$ .
1: if  $\text{epp} = \text{epp}_1 | \text{epp}_2$  then
2:   return  $\text{EVALUATE}(n, \text{epp}_1, \mathcal{G}) \cup \text{EVALUATE}(n, \text{epp}_2, \mathcal{G})$ 
3: if  $\text{epp} = \text{epp}_1 / \text{epp}_2$  then
4:    $Res' := \text{EVALUATE}(n, \text{epp}_1, \mathcal{G})$ 
5:    $Res = \{\}$ 
6:   for all nodes  $n' \in Res'$  do
7:      $Res = Res \cup \text{EVALUATE}(n', \text{epp}_2, \mathcal{G})$ 
8:   return  $Res$ 
9: if  $\text{epp} = \text{epp}_1 \& \text{epp}_2$  then
10:  return  $\text{EVALUATE}(n, \text{epp}_1, \mathcal{G}) \cap \text{EVALUATE}(n, \text{epp}_2, \mathcal{G})$ 
11: if  $\text{epp} = \text{epp}_1 \sim \text{epp}_2$  then
12:  return  $\text{EVALUATE}(n, \text{epp}_1, \mathcal{G}) \setminus \text{EVALUATE}(n, \text{epp}_2, \mathcal{G})$ 
13: if  $\text{epp} = \text{epp}_1 ?$  then
14:  return  $\{n\} \cup \text{EVALUATE}(n, \text{epp}_1, \mathcal{G})$ 
15: if  $\text{epp} = \text{pos}_1 \text{test} \text{pos}_2$  then
16:   $Res = \{\}$ 
17:  for all triple  $t \in \mathcal{G}$  s.t.  $P_m(\text{pos}_1, t) = n$  do
18:    if  $\text{EVALTEST}(t, \text{test}, \mathcal{G})$  then
19:       $Res = Res \cup \{P_m(\text{pos}_1, t), P_m(\text{pos}_2, t)\}$ 
20:  return  $Res$ 

Function EVALTEST( $t, \text{test}, \mathcal{G}$ )
  Input: triple  $t$ , graph  $\mathcal{G}$ ; Output:  $\text{true}$  if  $t$  satisfy  $\text{test}$ .
1: if  $\text{test} = \text{test}_1 \& \& \text{test}_2$  then
2:   return  $\text{EVALTEST}(t, \text{test}_1, \mathcal{G}) \wedge \text{EVALTEST}(t, \text{test}_2, \mathcal{G})$ 
3: if  $\text{test} = \text{test}_1 || \text{test}_2$  then
4:   return  $\text{EVALTEST}(t, \text{test}_1, \mathcal{G}) \vee \text{EVALTEST}(t, \text{test}_2, \mathcal{G})$ 
5: if  $\text{test} = !\text{test}_1$  then
6:   return  $\neg \text{EVALTEST}(t, \text{test}_1, \mathcal{G})$ 
7: if  $\text{test} = u$  then
8:   return  $P_m(\cdot, t) = u$ 
9: if  $\text{test} = \text{TP}(\text{pos}, \text{epp})$  then
10:  return  $\text{EVALUATE}(P_m(\text{pos}, t), \text{epp}, \mathcal{G}) \neq \emptyset$ 
11: if  $\text{test} = \text{T}(\text{EExp})$  then
12:  return  $\text{EvalSPARQLBuilt-in}(\text{EExp}, t)$ 

```

Figure 3: EPPs evaluation algorithm.

3.1 Recursion-based Algorithm

The algorithm for the evaluation of an EPP expression starts by invoking EVALUATE (Fig. 2), which receives as input a graph \mathcal{G} , an expression epp and a node n . If epp is non recursive (i.e., it does not contain the closure operators ‘+’

and ‘*’) then it is given as input to the function BASE, which considers the various forms of syntactic expressions. For recursive expressions the algorithm uses the function CLOSURE. Finally, the boolean function EVALTEST handles the different types of test .

The result of the evaluation of an EPP expression epp from a node n is a set of pairs of nodes (n, n_r) where nodes n_r are reachable from n via paths satisfying epp .

We now discuss the complexity of the algorithm. We assume \mathcal{G} to be stored by its adjacency list. In particular, for each $t \in \mathcal{T}$, a Hashtable is maintained where the set of keys is the set of predicates p such that there exists a triple in \mathcal{G} having as subject t and as predicate p , and the set of values are lists of objects o reachable by traversing p -predicates from t . We assume that given t and a predicate p the set of nodes reachable can be accessed in time $O(1)$. An additional Hashtable is used for inverse navigation, that is, for navigation starting on the object and ending on the subject. Both structures use space $O(|\mathcal{G}|)$. Let $|\text{epp}|$ be the size of the EPP expression epp .

Theorem 3 Given the EPP expression epp , an RDF graph \mathcal{G} and a node $n \in \mathcal{G}$ the evaluation of $\text{E}[\text{epp}]^{\mathcal{G}}$ can be performed in time $O(|\mathcal{G}| \cdot |\text{epp}|) + c_{\text{EExp}}$.

Proof: [SKETCH] The function EVALUATE is recursively called on each sub-expression of the epp in input; if such sub-expressions are not recursive (i.e., do not contain ‘*’, ‘+’), EVALUATE is invoked at most $O(|\text{epp}|)$ times. The base cases (lines 15-19 of function BASE) require to consider at most all the edges for all the nodes; this can be done in time $O(|\mathcal{G}|)$. If epp is recursive, the function CLOSURE is executed at most $O(\text{nodes}(\mathcal{G}))$ times; the procedure EVALUATE is invoked for each node in the worst case. When evaluating a subexpression from a node we use *memoization* to store its result (i.e., the set of reachable nodes) thus avoiding to recompute the same expression from the same node multiple times. Memoization guarantees that the total time required by CLOSURE is $O(|\text{epp}| \cdot |\mathcal{G}|)$. As for *nested* expressions, memoization enables to mark nodes of the graph satisfying a given subexpression.

Path conjunction and negation, corresponding to intersection and difference of set of nodes respectively (line 10 and 12 of BASE), can be computed in time $O(|\mathcal{G}|)$ by using a (perfect) hash function as the graph is known beforehand. As for tests, their cost is constant for *logical operators* and simple URI checking. The complexity is parametric wrt the cost of other SPARQL-based built-in conditions EExp (c_{EExp}). Finally, observe that with memoization the space complexity is $O(|\text{epp}| \cdot \text{nodes}(\mathcal{G})^2)$. \square

3.2 Translation into SPARQL

The W3C spec. (Harris and Seaborne 2013) *informally* mentions the fact that non-recursive property paths can be evaluated via a *translation* into equivalent SPARQL algebraic forms. However, no formal proof of the correctness and completeness of such translation is provided. Recursive property paths are handled in the standard via auxiliary functions called ALP (Harris and Seaborne 2013). With EPPs

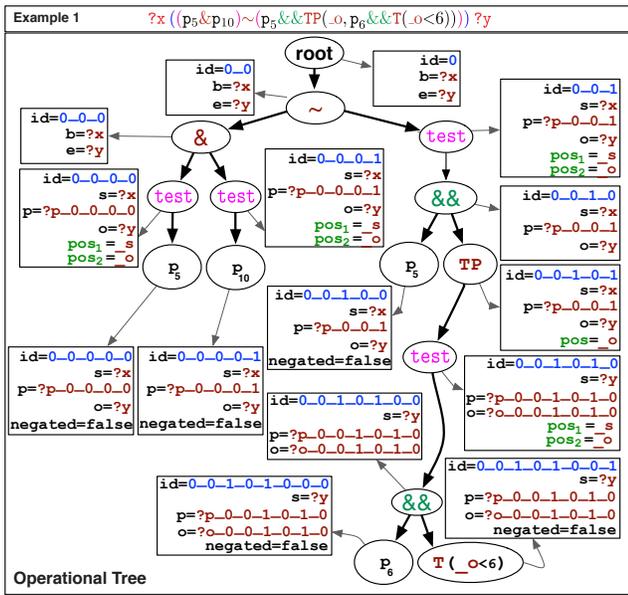


Figure 4: Operational tree for Ex. 1 after propagation.

(that are meant to extend property paths) we adopted a similar approach. In what follows, we give an overview of how *non-recursive* EPPs (NREPPs) are translated into SPARQL queries. For recursive EPPs we modified the ALP function defined in the W3C standard where subexpressions in recursive calls to ALP are evaluated via the SPARQL translation.

The NREPP to SPARQL translation

Let $\mathcal{P} = (\alpha, \text{epp}, \beta)$, $\alpha, \beta \in \mathcal{T} \cup \mathcal{V}$ be a NREPP pattern. The translation algorithm \mathcal{A}^t uses a structure called *operational tree* built from the *parse tree* of \mathcal{P} . The operational tree associated to \mathcal{P} is an ordered, labeled, rooted tree with node attributes reflecting the operational structure of \mathcal{P} . We now give a high level overview of the three phases of the translation algorithm \mathcal{A}^t :

R _m	$\Theta^{\mathcal{P}}(\text{root})$	$:= \text{'SELECT DISTINCT' root.b root.e ' (WHERE (' } \Theta^{\mathcal{P}}(\text{root.child}(1)) \text{')'}$
R ₀	$\Gamma(n)$	$:= n.s \ n.p \ n.o \text{'}$
R ₁	$\Theta^{\mathcal{P}}(n \sim)$	$:= \Theta^{\mathcal{P}}(n.\text{child}(1))$
R ₂	$\Theta^{\mathcal{P}}(n \setminus)$	$:= \Theta^{\mathcal{P}}(n.\text{child}(1)) \ \Theta^{\mathcal{P}}(n.\text{child}(2))$
R ₃	$\Theta^{\mathcal{P}}(n \cup)$	$:= \{ \Theta^{\mathcal{P}}(n.\text{child}(1)) \} \text{ UNION } \{ \Theta^{\mathcal{P}}(n.\text{child}(2)) \}$
R ₄	$\Theta^{\mathcal{P}}(n \&)$	$:= \Theta^{\mathcal{P}}(n.\text{child}(1)) \ \Theta^{\mathcal{P}}(n.\text{child}(2))$
R ₅	$\Theta^{\mathcal{P}}(n \sim)$	$:= \{ \Theta^{\mathcal{P}}(n.\text{child}(1)) \} \text{ MINUS } \{ \Theta^{\mathcal{P}}(n.\text{child}(2)) \}$
R ₆	$\Theta^{\mathcal{P}}(n^{\text{test}})$	$:= \Theta^{\mathcal{P}}(n.\text{child}(1))$
R ₇	$\Theta^{\mathcal{P}}(n^u)$	$:= \Gamma(n) \ \text{'FILTER' } (n.p[!]= u) \text{'}$ ^a
R ₈	$\Theta^{\mathcal{P}}(n^{\text{T(EEExp)}}$	$:= \Gamma(n) \ \text{'FILTER' } [!]= (\ \text{EEExp} \text{'}) \text{'}$ ^a
R ₉	$\Theta^{\mathcal{P}}(n^{\text{TP}})$	$:= \Gamma(n) \ \text{'FILTER' } [\text{'NOT' }] \ \text{'EXISTS' } \{ \Theta^{\mathcal{P}}(n.\text{child}(1)) \} \text{'}$ ^a
R ₁₀	$\Theta^{\mathcal{P}}(n \&\&)$	$:= \Theta^{\mathcal{P}}(n.\text{child}(1)) \ \Theta^{\mathcal{P}}(n.\text{child}(2))$
R ₁₁	$\Theta^{\mathcal{P}}(n \parallel)$	$:= \{ \Theta^{\mathcal{P}}(n.\text{child}(1)) \} \text{ UNION } \{ \Theta^{\mathcal{P}}(n.\text{child}(2)) \}$

Figure 5: Translation rules. ^aIf the node in input has the attribute `negated=true` use `!` (resp., `NOT`).

(i) Building of the operational tree

Consider a NREPP pattern \mathcal{P} . In the operational tree associated to \mathcal{P} , nodes are of two types: (i) *operational nodes* (labeled with $\sim, \&, /, |, \wedge$) and (ii) *test nodes* (labeled with `test`⁴, `||`, `&&`, `TP`, `u`, `T(EEExp)`). Each node n_i in the operational tree has a single parent and its children are an ordered set. Moreover, n_i has an attribute `id` whose unique value is computed as the *concatenation* of the parent's `id` and n_i 's position in the ordered set of children. The operational tree for Example 1 is shown in Fig. 4.

(ii) Propagation of variables and terms⁵

The key point is the propagation of variable names and RDF terms (kept in nodes' attributes). This is done by traversing the operational tree top-down and propagating attributes from each parent node to its children.

(iii) Application of the translation rules

This phase always starts by applying rule R_m in Fig. 5 on the `root` of the operational tree. This generates the outermost part of the final rewritten query into the SPARQL syntax. The translation proceeds by applying rules at each node of the operational tree visited according to a pre-order depth-first traversal.

For instance, in Fig. 4, after the `root`, the node with `id=0.0` and labeled with \sim is visited. This causes the triggering of rule R_5 , which enables to generate another (internal) chunk of the final SPARQL query. Node `0.0` is an operational node representing path negation \sim . As it can be noted, our translation procedure uses the SPARQL `MINUS` operator to reflect the semantics of EPPs dealing with path negation. The semantics of nodes representing EPPs `test` (e.g., `0.0-1.0-1`) is reflected into SPARQL via the `FILTER` operator.

Correctness of the translation

In order to prove the correctness of the translation of NREPPs into SPARQL we have defined a SPARQL based semantics (not reported here for sake of space) where conjunction ($\&$) and negation (\sim) are translated into join (\bowtie) and difference (\setminus) of (multi)sets of solution mappings (Pérez, Arenas, and Gutierrez 2009). Let $\llbracket \mathcal{P} \rrbracket_{\mathcal{G}}^{\text{epps}}$ denote the SPARQL-based semantics of EPPs, where \mathcal{P} is an NREPP pattern, and $\llbracket S \rrbracket_{\mathcal{G}}^{\text{W3C}}$ denote the SPARQL W3C semantics (Harris and Seaborne 2013), where S is a SPARQL query. The following theorem shows the correctness of the translation.

Theorem 4 The translation algorithm \mathcal{A}^t is correct and runs in polynomial time in the size of the expression to be translated. Moreover, for any RDF graph \mathcal{G} given \mathcal{P} it holds that $\llbracket \mathcal{P} \rrbracket_{\mathcal{G}}^{\text{epps}} = \llbracket S \rrbracket_{\mathcal{G}}^{\text{W3C}}$, where $S = \mathcal{A}^t(\mathcal{P})$ is the query produced by translating \mathcal{P} .

Proof: [SKETCH] \mathcal{A}^t is polynomial as it requires one scan of the operational tree (whose size is polynomial in the size of the EPP expression). As for the correctness, it is enough to prove that the propagation of variable names and RDF terms in the operational tree is correct. This proof can be

⁴In the syntax it corresponds to `pos1 test pos2`.

⁵Negation (!) for test nodes is propagated top-down when building the tree via the attribute `negated`.

done by structural induction on the depth of the operational tree. Finally, the semantic equivalence can be proved by associating to all types of EPP patterns the corresponding semantics as per SPARQL specification (Harris and Seaborne 2013). \square

Complexity of the translated queries. Consider a NREPP e_{pp} in a pattern \mathcal{P} and its translation into SPARQL $\mathcal{S}_{e_{pp}}$. Clearly, the complexity of evaluating $\mathcal{S}_{e_{pp}}$ depends on the fragment of SPARQL used in the translation. In particular (see Fig. 5) we make usage of SELECT, UNION, MINUS and FILTER and we *do not* use OPTIONAL. The complexity of this and other SPARQL fragments has been studied in (Pérez, Arenas, and Gutierrez 2009).

Summary. EPPs can be seen as a language per se and can be evaluated by implementing the algorithm in Fig. 3 in a custom processor, thus being independent from existing SPARQL processors. EPPs can also be seen as a conservative extension of SPARQL property paths and, thanks to the translation algorithm, can be evaluated on existing SPARQL processors. We have implemented both strategies; their pros and cons will be discussed in Section 4.

3.3 Comparison with Related Languages

As the goal of EPPs is to extend the expressive power of SPARQL property paths (PPs) and NREs-based languages, we compared such proposals with EPPs; Table 3 summarizes the (informal) comparison. We consider the following features of EPPs: path conjunction (&), path negation (~), nesting (TP), tests over node values (*test*) and usage of positions (*pos*). As it can be observed, PPs are the least expressive language; they do not support any of the new features of EPPs. This motivated our choice of *extending* PPs as described in Section 2.1. As for NREs, they clearly support nesting but neither other types of tests (e.g., node equality) nor path conjunction/negation as discussed in (Barceló, Pérez, and Reutter 2012).

Language	&	~	TP	<i>test</i>	<i>pos</i>
EPPs	Yes	Yes	Yes	Yes	Yes
PPs	No	No	No	No	No
NREs	No	No	No	No	No
nSPARQL	Yes	Yes	Yes	No	Yes

Table 3: Comparison of EPPs with related languages.

As for nSPARQL (based on NREs), it supports path conjunction and negation only via the SPARQL algebra; it also supports nesting and positions. However, it does not allow to test (in)equalities of nodes reached with a nested expression. EPPs support logical combination of tests representing nesting and tests representing (in)equalities as well as *safe-negation*⁶. Note that neither PPs nor NREs nor nSPARQL can express Example 1. nSPARQL supports positions by transforming an RDF graph into another graph where *navigational axes* (similar to those defined in XPath) are made

⁶The first element must be positive path.

explicit. On the other hand, EPPs do not require any transformation. In the syntax of EPPs, positions enable to perform *edge to node* traversals allowing to reach the node representing the predicate of a triple whence it is possible to traverse e.g. the property hierarchy. Finally, we want to emphasize that EPPs are syntactically compatible with PPs and are the only language, which thanks to the NREPP_{to}SPARQL translation, can be used in existing SPARQL processors.

4 Implementation and Evaluation

We have implemented both a custom query processor for EPPs and the NREPP_{to}SPARQL translation⁷.

Dataset and query set. We used a crawl of the FOAF social network (~500MBs) obtained from the BTC2012⁸ by traversing from the URI of T. Berners-Lee (TBL) `foaf:knows` predicates up to distance 4. We call the resulting graph \mathcal{G}^F , which has ~4M triples. We created 4 groups $G_i, i \in \{1, \dots, 4\}$ of EPP expressions each with 3 queries; this gives a queryset Q of 12 queries. The experiments have been performed on an Intel i5 machine with 8GBs RAM. Results are the average of 5 runs after removing lowest and highest values.

Experiment 1: Running time

For each $e_{pp} \in Q$ we generated the corresponding SPARQL query $\mathcal{S}_{e_{pp}}$ via the algorithm described in Section 3.2. We measured the execution time for each $e_{pp} \in Q$ with *our processor* and the execution time for $\mathcal{S}_{e_{pp}}$ in *Jena ARQ*⁹. Fig. 6 shows the running times. The number of results ranges from ~50 to ~8000.

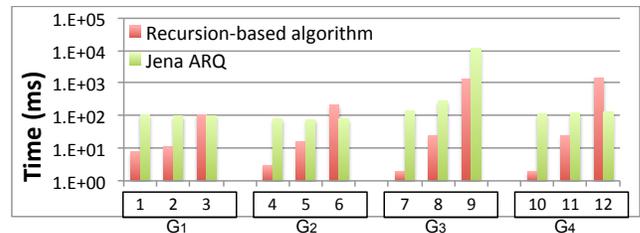


Figure 6: EPPs-custom-processor vs. Jena ARQ.

As it can be observed, for G_1 , which contains queries asking for friends of TBL up to distance 3, our evaluator performs better than ARQ at distance 1 and 2; at distance 3 times are comparable. The group G_2 additionally considers a test based on *nesting*. Again, our evaluator performs better at distance 1 and 2; at distance 3 it shows a higher running time. In G_3 , which considers *negation* (e.g., *exclusive friends* at various distances) our approach performs consistently better. Finally, in G_4 that includes conjunction (to ask for *mutual friends* at various distances) our approach performs better at distance 1 and 2 and obtains a higher running time at distance 3.

⁷Available at <http://extendedpps.wordpress.com>

⁸<http://km.aifb.kit.edu/projects/btc-2012>

⁹<http://jena.apache.org>

The complexity of the algorithm described in Section 3 is lower than that of Jena ARQ, which implements SPARQL. However, our evaluation suggests that for real-world data and natural queries (e.g., mutual friends) working with the translation of EPPs and using existing processors is still useful. Hence, the advantage of our approach is that navigational queries can be written in a succinct way via the EPPs syntax and execute after their translation via $\text{NREPPt} \circ \text{SPARQL}$. Anecdotally, while the EPP asking for mutual friends at distance 3 contains ~ 200 characters, the SPARQL query (obtained from the translation) contains ~ 700 characters. Moreover, when writing the SPARQL query one has also to deal with a large number of variables that need to be consistently joined.

Experiment 2: Translation overhead

We compared the elapsed times of our $\text{NREPPt} \circ \text{SPARQL}$ with the $\text{SPARQLt} \circ \text{ALGEBRA}$ translation of Jena ARQ. We used 28 queries generated in two steps. We started with 4 base expressions plus a fifth one combining them. Second, starting from them we generated increasingly longer expressions. We found that our $\text{NREPPt} \circ \text{SPARQL}$ translation performs comparably (in the order of ms) to the existing $\text{SPARQLt} \circ \text{ALGEBRA}$ translation.

5 Conclusions and Future Work

We introduced Extended Property Paths (EPPs), a significant extension of SPARQL property paths. We have provided an algorithm for their evaluation and also a translation procedure into SPARQL.

EPPs bring some benefits: (i) users can leverage their syntax to write more expressive *navigational* queries in a succinct way; (ii) EPPs are immediately available into existing SPARQL processors. Our study opens some research questions such as investigating properties of the translation (e.g., minimality) and studying optimization techniques. Also containment of EPPs is an open question.

References

- Alkhateeb, F.; Baget, J.-F.; and Euzenat, J. 2009. Extending SPARQL with Regular Expression Patterns (for querying RDF). *J. Web Semantics* 7(2):57–73.
- Arenas, M.; Gottlob, G.; and Pieris, A. 2014. Expressive Languages for Querying the Semantic Web. In *Proc. of the 33rd symposium on Principles of Database Systems*, 14–26. ACM.
- Barceló, P.; Libkin, L.; Lin, A. W.; and Wood, P. T. 2012. Expressive Languages for Path Queries over Graph-Structured Data. *ACM Transactions on Database Systems* 37(4):31.
- Barceló, P.; Pérez, J.; and Reutter, J. L. 2012. Relative Expressiveness of Nested Regular Expressions. In *Proc. of the 6th Alberto Mendelzon International Workshop on Foundations of Data Management*, 180–195.
- Berglund, A.; Boag, S.; Chamberlin, D.; Fernández, M.; Kay, M.; Robie, J.; and Siméon, J. 2010. XML Path Language (XPath) 2.0, W3C.
- Cali, A.; Gottlob, G.; Lukasiewicz, T.; Marnette, B.; and Pieris, A. 2010. Datalog+/-: A Family of Logical Knowledge Representation and Query Languages for New Applications. In *Proc. of the 25th Annual IEEE Symposium on Logic in Computer Science*, 228–242. IEEE.
- Fionda, V., and Pirrò, G. 2013. Querying Graphs with Preferences. In *Proc. of the 22nd ACM international conference on Conference on Information & Knowledge Management*, 929–938. ACM.
- Fionda, V.; Gutierrez, C.; and Pirrò, G. 2012. Semantic Navigation on the Web of Data: Specification of Routes, Web Fragments and Actions. In *Proc. of the 21st International Conference on World Wide Web*, 281–290. ACM.
- Fletcher, G. H.; Gyssens, M.; Leinders, D.; Van den Bussche, J.; Van Gucht, D.; Vansummeren, S.; and Wu, Y. 2011. Relative Expressive Power of Navigational Querying on Graphs. In *Proc. of the 14th International Conference on Database Theory*, 197–207. ACM.
- Harris, S., and Seaborne, A. 2013. SPARQL 1.1 Query Language, W3C.
- Kochut, K., and Janik, M. 2007. SPARQLer: Extended SPARQL for Semantic Association Discovery. In *Proc. of the 4th European Semantic Web Conf.*, 145–159.
- Libkin, L.; Reutter, J.; and Vrgoč, D. 2013. TriAL for RDF: Adapting Graph Query Languages for RDF Data. In *Proc. of the 32nd Symposium on Principles of Database Systems*, 201–212. ACM.
- Pérez, J.; Arenas, M.; and Gutierrez, C. 2009. Semantics and Complexity of SPARQL. *ACM Transactions on Database Systems* 34(3).
- Pérez, J.; Arenas, M.; and Gutierrez, C. 2010. nSPARQL: A Navigational Language for RDF. *J. Web Semantics* 8(4):255–270.
- Rudolph, S., and Krötzsch, M. 2013. Flag & check: Data Access with Monadically Defined Queries. In *Proc. of the 32nd Symposium on Principles of Database Systems*, 151–162. ACM.
- Wood, P. T. 2012. Query Languages for Graph Databases. *SIGMOD Rec.* 41(1).