

# Query Nesting, Assignment, and Aggregation in SPARQL 1.1

MARK KAMINSKI, EGOR V. KOSTYLEV, and BERNARDO CUENCA GRAU, University of Oxford

---

Answering aggregate queries is a key requirement of emerging applications of Semantic Technologies, such as data warehousing, business intelligence and sensor networks. In order to fulfil the requirements of such applications, the standardisation of SPARQL 1.1 led to the introduction of a wide range of constructs that enable value computation, aggregation, and query nesting. In this paper we provide an in-depth formal analysis of the semantics and expressive power of these new constructs as defined in the SPARQL 1.1 specification, and hence lay the necessary foundations for the development of robust, scalable and extensible query engines supporting complex numerical and analytics tasks.

CCS Concepts: • **Information systems** → **Query languages for non-relational engines; Resource Description Framework (RDF)**;

Additional Key Words and Phrases: SPARQL

## ACM Reference format:

Mark Kaminski, Egor V. Kostylev, and Bernardo Cuenca Grau. 0000. Query Nesting, Assignment, and Aggregation in SPARQL 1.1. *ACM Trans. Datab. Syst.* 0, 0, Article 00 (0000), 54 pages.  
<https://doi.org/0000001.0000001>

---

## 1 INTRODUCTION

The Resource Description Framework (RDF) is the World Wide Web Consortium (W3C) standard language for representing information in the Semantic Web [49]. RDF has been adopted in many applications as a convenient format for storing semistructured data. Moreover, the advent of major initiatives such as Linked Open Data [12] have led to a constant growth in the availability of RDF data online. RDF is based on a directed graph data model, where both nodes and edges are labelled. An RDF graph is a set of triples of the form  $(s, p, o)$ , which can be interpreted as edges labelled by  $p$  (the *predicate*) from nodes labelled by  $s$  (the *subject*) to nodes labelled by  $o$  (the *object*). The elements of a triple are typically *Internationalised Resource Identifiers (IRIs)*—global names that uniquely identify resources on the Web.

SPARQL, which became a W3C recommendation in 2008, is the standard query language for RDF [52]. Since its standardisation, the study of SPARQL as a database query language has attracted significant attention [43, 50]. A key distinguishing feature of SPARQL over previous RDF query languages is that it comes with a well-defined algebraic semantics, which has been the subject of intensive research and has laid the foundations for subsequent implementations [3, 41, 43, 50, 51, 53, 56]. The basic element of a SPARQL algebra query is a *basic graph pattern*, which is essentially

---

This work was supported by the Royal Society and the EPSRC projects MaSI<sup>3</sup>, Score!, DBOnto, and ED<sup>3</sup>.

Authors' address: Department of Computer Science, University of Oxford, Wolfson Building, Parks Road, OXFORD, OX1 3QD, UK.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 0000 Association for Computing Machinery.

0362-5915/0000/0-ART00 \$15.00

<https://doi.org/0000001.0000001>

an RDF graph where triples may also contain variables; the evaluation of a graph pattern over an RDF graph then consists of the mappings that homomorphically embed the pattern into the graph. Graph patterns can be combined with each other as well as with expressions using a rich set of algebraic operators.

Despite being a powerful query language, the first version of the SPARQL standard provided no support for queries which, rather than simply retrieving data, involve some form of computation or summarisation. This significantly limited the applicability of the language. Indeed, answering such queries is a key requirement in data warehousing and business intelligence, where data is aggregated across dimensions looking for patterns [1, 10, 22–24, 32, 58], as well as in emerging applications involving sensor networks and streaming RDF data [6, 11, 15].

The standardisation of a new version of the SPARQL query language, called SPARQL 1.1 [29], addressed these limitations by introducing a wide range of constructs supporting aggregation and value computation in line with those available in SQL:

- a collection of *aggregate functions* for value computation, such as Min, Max, Avg, Sum, and Count;
- the *grouping* constructs GROUP BY and HAVING, which restrict the application of aggregate functions to groups of solutions satisfying certain conditions;
- the *variable assignment* constructs BIND, VALUES, and AS which are used to assign the value of a complex (e.g., arithmetic) expression to a variable; and
- a *query nesting* mechanism for embedding queries within graph patterns as well as within expressions.

Similarly to its predecessor SPARQL 1.0, the semantics of SPARQL 1.1 is specified by means of an (extended) normative algebra. Many of the new features introduced in SPARQL 1.1, such as property paths [7, 40, 48, 57], query federation [13, 14], or entailment regimes [2, 8, 38, 39] have already received significant attention in the literature. In contrast, the theoretical properties of the algebraic operators that enable value computation, aggregation, and query nesting remain largely unexplored [4]. This is in stark contrast to the case of relational databases, where the formal properties of aggregation have been studied in depth [17–21, 30, 36, 44, 45, 54].

Our aim in this paper is to provide a systematic study of the semantics and expressive power of the SPARQL 1.1 algebra with variable assignment, aggregates and query nesting. Understanding the capabilities of the new constructs and their inter-dependencies is a key requirement for the development of query engines supporting complex numerical and analytics tasks while providing correctness, robustness, scalability and extensibility guarantees.

We take the well-known SPARQL algebra [50] as a starting point. Most existing works on SPARQL, including [50], assume that graph patterns are interpreted as sets of solution mappings rather than multisets (or bags) as in the normative specification. In this paper, we follow the normative specification and adopt multiset semantics from the outset, which we formalise in Section 2.

In Sections 3 and 4 we study the query nesting mechanisms in SPARQL 1.1. We first consider, in Section 3, the nesting of SELECT and SELECT DISTINCT query blocks. In algebraic terms, this amounts to allowing the unrestricted use of operators *Project* and *Distinct* rather than limiting them to the outermost query level as in SPARQL 1.0. We show that there is no gain of expressive power by allowing unrestricted use of just one of these operators. In contrast, if *both* operators are allowed without restrictions, we show how to construct queries that cannot be equivalently expressed without nesting. The additional expressive power is due to the interplay between the set semantics enforced by the *Distinct* operator and the bag semantics of *Project*—a phenomenon that was first observed in the relational case by Cohen [18], and was later conjectured by Angles and Gutierrez to yield additional expressive power in SPARQL [4]. As argued in Section 3.2, however,

the evidence given in [4] in support of their conjecture is unsatisfactory. Our results settle this question and provide a detailed account of which combinations of constructs lead to expressivity gains and which ones are redundant.

In addition to subqueries as patterns, SPARQL 1.1 also provides a mechanism for embedding graph patterns within expressions in filter conditions using the `exists` construct. We investigate this form of query nesting in Section 4. We first argue that the normative semantics of queries with `exists` expressions is not always well-defined. Thus, we propose a new semantics, which is defined for all queries. Our proposed semantics is a conservative extension of the normative one, in the sense that they coincide whenever the normative semantics is well-defined. We finally show that the `exists` construct can be simulated, thus not resulting in any additional expressive power.

In Section 5 we turn our attention to variable assignment. It is enabled in the normative algebra by the *Extend* operator, which extends solution mappings with a fresh variable assigned to the value of an expression. We show that *Extend* increases expressive power, as it introduces arithmetic into the language [54].

In Section 6 we analyse the SPARQL 1.1 aggregate algebra. It turns out that this algebra is rather non-standard when compared to its relational counterpart. It provides a great deal of power and flexibility, and we show that it can express all forms of query nesting and variable assignment previously described. We then define a normal form, which we exploit to define a substantial simplification of the normative aggregate algebra that eliminates most of its unconventional aspects. The resulting algebra is much closer to its relational counterpart and can be exploited to provide a more transparent algebraic translation of the SPARQL 1.1 syntax.

In Section 7 we study property paths, which bring a limited form of recursion into the language, and discuss their interplay with aggregation and arithmetic. We first demonstrate that the combination of aggregation and property paths brings additional expressive power into the language, as it allows us to construct queries that cannot be captured if either of these features is disallowed. We then show that the normative multiset semantics of property paths can be equivalently recast in terms of sets.

In Section 8 we provide a semantics for analytic OLAP queries (such as cube and window-based queries) in terms of our simplified aggregate algebra. We also discuss OLAP queries in the specific case where RDF data conforms to the W3C Data Cube specification [23].

Finally, in Section 9 we revisit the simplifying assumptions adopted in our formal presentation with respect to the standard, and discuss their implications.

The paper is accompanied with an online appendix, accessible in the ACM Digital Library, which contains the more technical aspects of certain proofs.

## 2 THE SPARQL ALGEBRA

We assume familiarity with the normative SPARQL syntax, which we will use for writing example queries. In this section we recapitulate the core fragment of the SPARQL 1.1 algebra capturing SPARQL 1.0, which we refer to as *Sparql*; in subsequent sections we gradually introduce and study other algebraic constructs. We also recapitulate here basic notions on RDF, query equivalence, and expressive power. In contrast to most papers in the literature, we follow the W3C standard by using the ternary left join (`OPTIONAL`) operator and by defining the semantics of the algebra in terms of bags rather than sets.

### 2.1 RDF Graphs

Let  $\mathbf{I}$ ,  $\mathbf{L}$ , and  $\mathbf{B}$  be countably infinite pairwise disjoint sets of *IRIs*, *literals*, and *blank nodes*, respectively, where literals can be numbers, strings, or Boolean values *true* and *false*. The set  $\mathbf{T}$  of (*RDF terms*) is

$I \cup L \cup B$ . An (RDF) triple is an element  $(s, p, o)$  of  $(I \cup B) \times I \times T$ , with  $s$  called the *subject*,  $p$  the *predicate*, and  $o$  the *object*. An (RDF) graph is a finite set of RDF triples. In this paper, we assume that all graphs are non-empty. Although several of our results hold only under that assumption, we do not see this as a limitation in practice: graphs can be trivially tested for emptiness and querying empty graphs is of no use in most applications.

## 2.2 SPARQL Algebra Syntax

We adopt the basic algebra from the normative specification [52], but omit certain features which are immaterial to our results, such that ternary expression operator *if* and auxiliary difference pattern operator *Diff*; these omissions are discussed in Section 9.

We distinguish three types of syntactic building blocks—expressions, patterns, and queries, which we define next for our basic algebra Sparql. They are built over terms  $T$  and an infinite set  $X = \{?x, ?y, \dots\}$  of *variables*, disjoint from  $T$ . In this paper we will often use the following notion: a *renaming* of a set of variables  $X$  is an injective substitution from  $X$  into fresh variables; then, given a syntactic object  $s$  such as an expression, pattern, query, and a renaming  $\theta$  of  $X$ ,  $s\theta$  is the object obtained from  $s$  by a simultaneous replacement of all  $?x \in X$  by  $\theta(?x)$ .

*Expressions* in Sparql are inductively defined as follows:

- all variables in  $X$  and all terms in  $I \cup L$  are expressions;
- if  $?x \in X$ , then  $\text{bound}(?x)$  is an expression;
- if  $t \in T \cup X$ , then  $\text{isIRI}(t)$ ,  $\text{isLiteral}(t)$ , and  $\text{isBlank}(t)$  are expressions;
- if  $E_1$  and  $E_2$  are expressions, then so are
  - $(E_1 + E_2)$ ,  $(E_1 - E_2)$ ,  $(E_1 * E_2)$ ,  $(E_1/E_2)$ ;
  - $(E_1 \doteq E_2)$ ,  $(E_1 < E_2)$ ,  $(E_1 > E_2)$ ;
  - $(\neg E)$ ,  $(E_1 \wedge E_2)$ , and  $(E_1 \vee E_2)$ .

We use  $E_1 \rightarrow E_2$  and  $E_1 \leftrightarrow E_2$  as the usual abbreviations for  $\neg E_1 \vee E_2$  and  $(E_1 \wedge E_2) \vee (\neg E_1 \wedge \neg E_2)$ , respectively. Furthermore, given a set of variables  $X$  and a renaming  $\theta$  of  $X$ , we denote with  $\text{eq}(X, \theta)$  the expression

$$\bigwedge_{?x \in X} ((\text{bound}(?x) \leftrightarrow \text{bound}(?x\theta)) \wedge (\text{bound}(?x) \rightarrow ?x \doteq ?x\theta)).$$

This expression checks that values of variables  $X\theta$  duplicate the values of  $X$ .

The set  $\text{var}(E)$  of variables *in the scope of* an expression  $E$  is defined as follows:

- if  $E$  is one of  $?x$ ,  $\text{bound}(?x)$ ,  $\text{isIRI}(?x)$ ,  $\text{isLiteral}(?x)$ , and  $\text{isBlank}(?x)$  then  $\text{var}(E) = \{?x\}$ ;
- if  $E$  is a term from  $I \cup L$  or one of  $\text{isIRI}(t)$ ,  $\text{isLiteral}(t)$ , and  $\text{isBlank}(t)$  for  $t \in T$  then  $\text{var}(E) = \emptyset$ ;
- if  $E = E_1 \circ E_2$ , for  $\circ \in \{+, -, *, /, \doteq, <, >, \wedge, \vee\}$  then  $\text{var}(E) = \text{var}(E_1) \cup \text{var}(E_2)$ ;
- if  $E = \neg E'$  then  $\text{var}(E) = \text{var}(E')$ .

Note that for a Sparql expression  $E$  the set  $\text{var}(E)$  just consists of all the variables occurring in  $E$ . We have given a more explicit definition because this will be different for some extensions of Sparql studied later on.

*Patterns* in Sparql are inductively defined as follows:

- a *basic graph pattern (BGP)* is a possibly empty set of *triple patterns*, that is, elements of the set

$$(I \cup L \cup X) \times (I \cup X) \times (I \cup L \cup X);$$

- $\text{Join}(P_1, P_2)$  and  $\text{Union}(P_1, P_2)$  are patterns if  $P_1$  and  $P_2$  are patterns;
- $\text{Filter}(E, P)$  is a pattern if  $P$  is a pattern and  $E$  is an expression; and
- $\text{LeftJoin}(E, P_1, P_2)$  is a pattern if  $P_1$  and  $P_2$  are patterns and  $E$  is an expression.

For brevity, we will often omit braces when writing BGPs consisting of single triple patterns (e.g., we write  $(?x, ?y, ?z)$  instead of  $\{(?x, ?y, ?z)\}$ ).

The set  $\text{var}(P)$  of variables *in the scope of* a pattern  $P$  is as follows:

- $\text{var}(B)$  for a BGP  $B$  is the set of all its variables;
- $\text{var}(Op(S_1, S_2)) = \text{var}(S_1) \cup \text{var}(S_2)$  if  $Op$  is one of *Join*, *Union* and *Filter*;
- $\text{var}(LeftJoin(E, P_1, P_2)) = \text{var}(E) \cup \text{var}(P_1) \cup \text{var}(P_2)$ .

Similarly to expressions,  $\text{var}(P)$  is just the set of all variables in  $P$  for Sparql, which will be different for fragments defined later on.

*Queries* in Sparql are expressions of the form  $Project(X, P)$  or  $Distinct(Project(X, P))$ , for  $P$  a pattern and  $X$  a set of variables (called *free variables*).

The variables  $\text{var}(Q)$  *in the scope of* a query  $Q$  are its free variables.

### 2.3 SPARQL Algebra Semantics

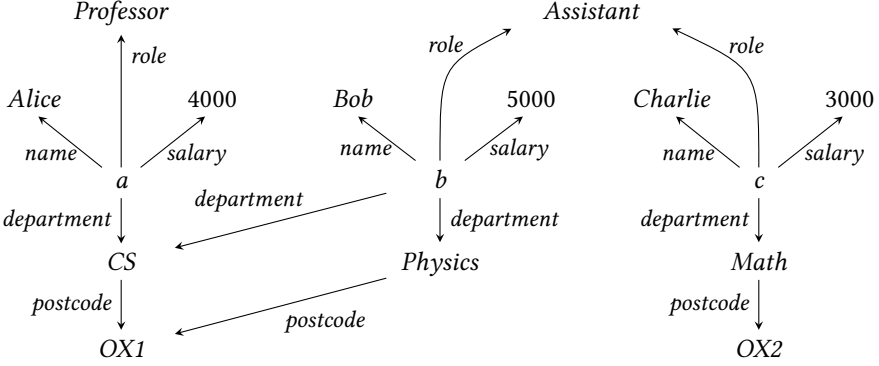
The semantics of Sparql is defined in terms of (*solution*) *mappings*, that is, partial functions  $\mu$  from variables  $X$  to terms  $T$ . The domain of  $\mu$ , denoted  $\text{dom}(\mu)$ , is the set of variables over which  $\mu$  is defined. Mappings  $\mu_1$  and  $\mu_2$  are *compatible*, written  $\mu_1 \sim \mu_2$ , if  $\mu_1(?x) = \mu_2(?x)$  for each  $?x$  in  $\text{dom}(\mu_1) \cap \text{dom}(\mu_2)$ . If  $\mu_1 \sim \mu_2$ , then  $\mu_1 \cup \mu_2$  is the mapping obtained by extending  $\mu_1$  according to  $\mu_2$  on all the variables in  $\text{dom}(\mu_2) \setminus \text{dom}(\mu_1)$ . We also write  $\mu_1 \subseteq \mu_2$  if  $\mu_1 \sim \mu_2$  and  $\text{dom}(\mu_1) \subseteq \text{dom}(\mu_2)$ . The *restriction*  $\mu|_X$  of a mapping  $\mu$  to a set of variables  $X$  is the mapping that coincides with  $\mu$  on  $X$  and is undefined on all other variables.

*Expressions.* The *evaluation*  $\llbracket E \rrbracket_G^\mu$  of an expression  $E$  over a graph  $G$  with respect to a mapping  $\mu$  is a *value* in  $T \cup \{\text{error}\}$ , as defined next (graph  $G$  does not affect the semantics of expressions in Sparql, but it will do so in relevant extensions of Sparql discussed later on; hence, we include here the subscript  $G$  for completeness):

- $\llbracket ?x \rrbracket_G^\mu$  is  $\mu(?x)$  if  $?x \in \text{dom}(\mu)$  and *error* otherwise;
- $\llbracket \ell \rrbracket_G^\mu$  is  $\ell$  for  $\ell \in I \cup L$ ;
- $\llbracket \text{bound}(?x) \rrbracket_G^\mu$  is *true* if  $?x \in \text{dom}(\mu)$ , and it is *false* otherwise;
- $\llbracket \text{isIRI}(t) \rrbracket_G^\mu$ ,  $\llbracket \text{isLiteral}(t) \rrbracket_G^\mu$ , or  $\llbracket \text{isBlank}(t) \rrbracket_G^\mu$  is *true* if  $t$  is an IRI, literal, or blank node, respectively, and *false* otherwise;
- $\llbracket E_1 \circ E_2 \rrbracket_G^\mu$ , for an arithmetic or comparison operator  $\circ \in \{+, -, *, /, \doteq, <, >\}$ , is either  $\llbracket E_1 \rrbracket_G^\mu \circ \llbracket E_2 \rrbracket_G^\mu$  if  $\llbracket E_1 \rrbracket_G^\mu$  and  $\llbracket E_2 \rrbracket_G^\mu$  are both different from *error* and of suitable types, or *error* otherwise;
- $\llbracket \neg E' \rrbracket_G^\mu$  is *true* if  $\llbracket E' \rrbracket_G^\mu$  is *false*, it is *false* if  $\llbracket E' \rrbracket_G^\mu$  is *true*, and it is *error* otherwise;
- $\llbracket E_1 \wedge E_2 \rrbracket_G^\mu$  is *true* if both  $\llbracket E_1 \rrbracket_G^\mu$  and  $\llbracket E_2 \rrbracket_G^\mu$  are *true*, it is *false* if  $\llbracket E_1 \rrbracket_G^\mu$  or  $\llbracket E_2 \rrbracket_G^\mu$  is *false*, and it is *error* otherwise;
- $\llbracket E_1 \vee E_2 \rrbracket_G^\mu$  is equal to  $\llbracket \neg(\neg E_1 \wedge \neg E_2) \rrbracket_G^\mu$ .

*Patterns.* The semantics of patterns is based on *multisets*  $\Omega = (S_\Omega, \text{card}_\Omega)$ , where  $S_\Omega$  is the *base set* of mappings, and the *multiplicity* function  $\text{card}_\Omega$  assigns a positive number to each element of  $S_\Omega$ . Abusing notation, we often write  $\mu \in \Omega$  instead of  $\mu \in S_\Omega$ , as well as  $\Omega = \emptyset$  instead of  $S_\Omega = \emptyset$ . We use the following operations on multisets.

- The *union*  $\Omega_1 \uplus \Omega_2$  of multisets  $\Omega_1$  and  $\Omega_2$  of mappings is the multiset with base set  $S_{\Omega_1} \cup S_{\Omega_2}$  such that the multiplicity of each mapping  $\mu$  is  $\text{card}_{\Omega_1}(\mu) + \text{card}_{\Omega_2}(\mu)$  if  $\mu$  is in both  $S_{\Omega_1}$  and  $S_{\Omega_2}$ , and  $\text{card}_{\Omega_i}(\mu)$  if  $\mu$  is only in one  $S_{\Omega_i}$ .
- The *restriction*  $\llbracket \mu \mid \mu \in \Omega, \text{Cond}(\mu) \rrbracket$  of a multiset  $\Omega$  of mappings given a Boolean condition  $\text{Cond}$  on a mapping  $\mu$  is the multiset with the base set consisting of all  $\mu \in \Omega$  with  $\text{Cond}(\mu)$  true such that the multiplicity of each such  $\mu$  coincides with the multiplicity of  $\mu$  in  $\Omega$ .

Fig. 1. Example RDF graph  $G_{ex}$ 

We consider two generalisations of multiset restriction:

- $\llbracket \mu' \mid \mu \in \Omega, \text{Cond}(\mu', \mu) \rrbracket$  for a multiset  $\Omega$  of mappings given a Boolean condition  $\text{Cond}$  on two mappings  $\mu'$  and  $\mu$  is the multiset with the base set consisting of all  $\mu'$  for which there exists  $\mu \in \Omega$  with  $\text{Cond}(\mu', \mu)$  true such that the multiplicity of each such  $\mu'$  is the sum of the multiplicities of all the contributing  $\mu$  in  $\Omega$  (i.e., all  $\mu$  for which  $\text{Cond}(\mu', \mu)$  is true);
- $\llbracket \mu' \mid \mu_1 \in \Omega_1, \mu_2 \in \Omega_2, \text{Cond}(\mu', \mu_1, \mu_2) \rrbracket$  for multisets  $\Omega_1$  and  $\Omega_2$  of mappings given a Boolean condition  $\text{Cond}$  on three mappings  $\mu', \mu_1$  and  $\mu_2$  is the multiset with the base set consisting of all  $\mu'$  for which there exist  $\mu_1$  in  $\Omega_1$  and  $\mu_2$  in  $\Omega_2$  with  $\text{Cond}$  true such that the multiplicity of each such  $\mu'$  is the following sum ranging over the pairs of contributing  $\mu_1$  and  $\mu_2$ :

$$\sum \text{card}_{\Omega_1}(\mu_1) \times \text{card}_{\Omega_2}(\mu_2).$$

The *evaluation*  $\llbracket P \rrbracket_G$  of a pattern  $P$  over a graph  $G$  is defined as follows:

- $\llbracket B \rrbracket_G$  for a BGP  $B$  is the multiset with  $S_{\llbracket B \rrbracket_G}$  consisting of all  $\mu$  with  $\text{dom}(\mu) = \text{var}(B)$  and  $\mu(B) \subseteq G$ , and with  $\text{card}_{\llbracket B \rrbracket_G}(\mu) = 1$  for each such  $\mu$ , where  $\mu(B)$  is the BGP obtained from  $B$  by replacing its variables according to  $\mu$ ;
- $\llbracket \text{Join}(P_1, P_2) \rrbracket_G = \llbracket \mu \mid \mu_1 \in \llbracket P_1 \rrbracket_G, \mu_2 \in \llbracket P_2 \rrbracket_G, \mu = \mu_1 \cup \mu_2 \rrbracket$ ;
- $\llbracket \text{Union}(P_1, P_2) \rrbracket_G = \llbracket P_1 \rrbracket_G \uplus \llbracket P_2 \rrbracket_G$ ;
- $\llbracket \text{Filter}(E, P') \rrbracket_G = \llbracket \mu \mid \mu \in \llbracket P' \rrbracket_G, \llbracket E \rrbracket_G^\mu = \text{true} \rrbracket$ ; and
- $\llbracket \text{LeftJoin}(E, P_1, P_2) \rrbracket_G = \llbracket \mu \mid \mu \in \llbracket \text{Join}(P_1, P_2) \rrbracket_G, \llbracket E \rrbracket_G^\mu = \text{true} \rrbracket \uplus \llbracket \mu \mid \mu \in \llbracket P_1 \rrbracket_G, \forall \mu_2 \in \llbracket P_2 \rrbracket_G. (\mu \approx \mu_2 \text{ or } \llbracket E \rrbracket_G^{\mu \cup \mu_2} \neq \text{true}) \rrbracket$ .<sup>1</sup>

Note that the condition  $\mu = \mu_1 \cup \mu_2$  in the case of *Join* implies that  $\mu_1 \sim \mu_2$ ; the latter is not explicitly written here and in similar places throughout the paper for brevity. Note also that the semantics of a non-empty BGP  $B = \{T_1, \dots, T_n\}$  can be alternatively defined as  $\llbracket \text{Join}(\dots \text{Join}(T_1, T_2), \dots, T_n) \rrbracket_G$ .

*Queries.* We conclude with the *evaluation* of queries, which are also evaluated as multisets of mappings in our formalisation:

- $\llbracket \text{Project}(X, P) \rrbracket_G$  is a multiset whose base set consists of the restrictions  $\mu|_X$  of all  $\mu$  in  $\llbracket P \rrbracket_G$ , and the multiplicity of each such  $\mu|_X$  is the sum of multiplicities of all corresponding  $\mu$  in  $\llbracket P \rrbracket_G$ ;

<sup>1</sup>The SPARQL 1.1 specification requires  $\llbracket E \rrbracket_G^{\mu \cup \mu_2} = \text{false}$  instead of  $\llbracket E \rrbracket_G^{\mu \cup \mu_2} \neq \text{true}$  in the second argument, but we follow the errata document [31] here. This choice does not affect any results of this paper; see a discussion about this issue in [38].

- $\llbracket \text{Distinct}(Q) \rrbracket_G$  is the multiset with the same base set as  $\llbracket Q \rrbracket_G$ , but with multiplicity 1 for all mappings.

*Example 2.1.* Consider RDF graph  $G_{ex}$  depicted in Figure 1, which we will use as a running example. Let us evaluate the following SPARQL query over  $G_{ex}$ .

(Q0) Find the departments having at least one employee with salary over 3500.

It can be written in SPARQL 1.1 as follows:

```
SELECT ?d
WHERE { ?x department ?d . { ?x salary ?y . FILTER (?y > 3500) } }.
```

Query (Q0) is translated to the Sparql algebra as follows:

$$\text{Project}(\{?d\}, \text{Join}((?x, \text{department}, ?d), \text{Filter}((?y > 3500), (?x, \text{salary}, ?y))))).$$

The *Filter* subpattern evaluates to mappings

$$\mu_1 = \{?x \mapsto a, ?y \mapsto 4000\}, \quad \mu_2 = \{?x \mapsto b, ?y \mapsto 5000\},$$

both with multiplicity 1. Alice works only for the CS department, whereas Bob works for both CS and *Physics*; thus, (Q0) evaluates to the multiset of mappings consisting of  $\mu_3 = \{?d \mapsto \text{CS}\}$  with multiplicity 2 and  $\mu_4 = \{?d \mapsto \text{Physics}\}$  with multiplicity 1.

## 2.4 Extensions of SPARQL

In this paper, we compare Sparql with several richer languages. Formally, a Sparql *extension* is a language that allows, besides all the constructs of Sparql, for some other forms of expressions, patterns, or queries; additionally, we impose the following natural requirements on the evaluation of any pattern or query  $S$  over a graph  $G$ :

- $\text{dom}(\mu) \subseteq \text{var}(S)$  for any  $\mu \in \llbracket S \rrbracket_G$ ,
- for any renaming  $\theta$  of  $\text{var}(S)$  to variables not mentioned in  $S$ ,  $\mu\theta \in \llbracket S\theta \rrbracket_G$  if and only if  $\mu \in \llbracket S \rrbracket_G$ , and the multiplicity of  $\mu\theta$  in  $\llbracket S\theta \rrbracket_G$  equals that of  $\mu$  in  $\llbracket S \rrbracket_G$ .

Note that, under this definition, not every language that includes all of the Sparql queries can be seen as its extension; however, all the languages considered in this paper satisfy these requirements. These restrictions guarantee that  $S\theta$  for a renaming  $\theta$  of  $\text{var}(S)$  always give a copy of the answers of  $S$  over fresh variables. We silently use this guarantee throughout many proofs in this paper.

## 2.5 Expressive Power of Query Languages

Extending Sparql can lead to an increase in expressive power; that is, some queries in the extended language may not be equivalently expressible in the original language. We next make this notion precise.

A *query language for RDF* is a pair  $(\mathbf{Q}, \llbracket \cdot \rrbracket)$  where  $\mathbf{Q}$  is a class of *queries* and  $\llbracket \cdot \rrbracket$  is an *evaluation function* that maps queries and RDF graphs to multisets of solution mappings (e.g., Sparql is a query language for RDF).

A query  $Q$  in a language  $\mathcal{L} = (\mathbf{Q}, \llbracket \cdot \rrbracket)$  is *equivalent* to a query  $Q'$  in a language  $\mathcal{L}' = (\mathbf{Q}', \llbracket \cdot \rrbracket')$  if  $\llbracket Q \rrbracket_G = \llbracket Q' \rrbracket'_G$  for every graph  $G$ . If such  $Q'$  exists, then  $Q$  is  *$\mathcal{L}'$ -expressible*. Language  $\mathcal{L}_1$  is *more expressive* than  $\mathcal{L}_2$  if every  $\mathcal{L}_1$  query is  $\mathcal{L}_2$ -expressible. For example, any Sparql extension is trivially more expressive than Sparql. Languages  $\mathcal{L}_1$  and  $\mathcal{L}_2$  have the *same expressive power* if each of them is more expressive than the other. Finally,  $\mathcal{L}_1$  is *strictly more expressive* than  $\mathcal{L}_2$  if it is more expressive than  $\mathcal{L}_2$ , but does not have the same expressive power as  $\mathcal{L}_2$ .

### 3 SUBQUERIES AS PATTERNS

In this and the next section we investigate the expressive power provided by nested queries, that is, those having other queries (called *subqueries*) embedded within. Subqueries can itself be nested queries; thus, queries can have a deeply nested structure.

The SPARQL 1.1 specification provides two mechanisms for query nesting. First, subqueries can play the role of patterns within the WHERE clause of another query. This is tantamount to allowing the arbitrary use of the operators *Project* and *Distinct* within patterns (in which case there is no real distinction between queries and patterns anymore), rather than allowing them only on the outermost level of queries. This basic form of nested queries is the topic of this section; in particular, after introducing a useful tool in Section 3.1, we define subqueries as patterns in Section 3.2 and study their expressive power in Section 3.3. In Section 4, we investigate the second mechanism for query nesting, where patterns can be embedded within expressions.

#### 3.1 Set Difference Operator

Before moving into further particulars, we first show that Sparql can express a “set difference” operator, which we will exploit to encode other constructs (note that this operator is different from the *Minus* and *Diff* operators discussed in Section 9 and in [3, 9, 37]).

*Definition 3.1.* For  $P_1$  and  $P_2$  patterns,  $SetMinus(P_1, P_2)$  is a pattern, whose semantics for a graph  $G$  is as follows:

$$\llbracket SetMinus(P_1, P_2) \rrbracket_G = \{ \mu \mid \mu \in \llbracket P_1 \rrbracket_G, \mu \notin \llbracket P_2 \rrbracket_G \}.$$

In contrast to the relational multiset difference operator, where the occurrences of an element in the right operand are subtracted from those in the left one, this operator yields each  $\mu$  in  $\llbracket P_1 \rrbracket_G$  but not in  $\llbracket P_2 \rrbracket_G$  with the same multiplicity as in  $\llbracket P_1 \rrbracket_G$ . Thus, mapping  $\mu$  is not returned whenever  $\mu \in \llbracket P_2 \rrbracket_G$ , regardless of its multiplicity in  $\llbracket P_1 \rrbracket_G$ .

*Example 3.2.* Consider the evaluation of the following pattern over graph  $G_{ex}$ :

$$SetMinus(Union(\{?x, department, CS\}, \{?x, salary, 5000\}), \{?x, department, Physics\}).$$

The inner *Union* pattern yields the multiset of mappings  $\{?x \mapsto a\}$  with multiplicity 1 and  $\{?x \mapsto b\}$  with multiplicity 2. Then, the evaluation of the overall *SetMinus* pattern yields  $\{?x \mapsto a\}$  with multiplicity 1. Note that the relational multiset difference operator would have additionally produced the mapping  $\{?x \mapsto b\}$  with multiplicity 1.

This pattern can be expressed in Sparql using *Filter*, *LeftJoin*, and bound, analogously to the well-known encoding of *Diff* (see e.g., [9]):

$$Filter(\neg bound(?u), LeftJoin(true, Union(\{?x, department, CS\}, \{?x, salary, 5000\}), \{?x, department, Physics\}, (?u, ?v, ?w))).$$

For more complex patterns, however, this encoding needs to be generalised to account for the possibility of some variables in the arguments of *LeftJoin* being unbound.

**LEMMA 3.3.** *For any extension Sparql<sub>★</sub> of Sparql, the pattern SetMinus( $P_1, P_2$ ) with Sparql<sub>★</sub> patterns  $P_1$  and  $P_2$  has an equivalent pattern in Sparql<sub>★</sub>.*

**PROOF.** We need to establish a mechanism that distinguishes, for any graph  $G$ , between the mappings in  $\llbracket P_1 \rrbracket_G$  that occur in  $\llbracket P_2 \rrbracket_G$  and those mappings that do not.

As the first step, we construct a pattern  $P'$  whose evaluation contains each mapping  $\mu$  in  $\llbracket P_1 \rrbracket_G$ ; if  $\mu$  is also in  $\llbracket P_2 \rrbracket_G$ , however, the corresponding mapping in  $\llbracket P' \rrbracket_G$  is extended with a copy of  $\mu$

(on fresh variables) as well as a binding for the fresh variables  $?x$ ,  $?y$ , and  $?z$ . Then we show that  $\text{SetMinus}(P_1, P_2)$  is equivalent to

$$P = \text{Filter}(\neg\text{bound}(?x), P').$$

Intuitively, the filter condition eliminates all mappings with fresh variables, thus producing the intended result. Before moving to the details, note that in this construction the multiplicities are important only for those mappings in  $\llbracket P' \rrbracket_G$  that, restricted to  $\text{var}(P_1)$ , do not appear in  $\llbracket P_2 \rrbracket_G$ .

We construct  $P'$  as follows, where  $X = \text{var}(P_1) \cup \text{var}(P_2)$  and  $\theta$  is a renaming of  $X$ :

$$P'_2 = \text{Join}(P_2, (?x, ?y, ?z)), \quad P' = \text{LeftJoin}(\text{eq}(X, \theta), P_1, P'_2\theta).$$

The following claim completes the proof of the lemma.

CLAIM 3.4. *Patterns  $\text{SetMinus}(P_1, P_2)$  and  $P$  are equivalent.*

The proof of this claim can be found in the appendix. □

To check if a mapping constitutes a match for a subpattern, the proof of Lemma 3.3 uses a join with a triple pattern  $(?x, ?y, ?z)$ , which matches any triple in a graph, coupled with a check of  $?x$  for boundedness. This idea is used in several proofs throughout the paper. As mentioned above, it is commonly used to express the auxiliary *Diff* operator via the core Sparql operators (see Section 9 for more details). Note, however, that the idea works only if empty graphs are not allowed [37].

### 3.2 Syntax and Semantics of Subqueries as Patterns

We start by discussing the basic nesting mechanism in SPARQL 1.1, which allows queries to be subpatterns of other patterns as illustrated by the following example.

*Example 3.5.* Consider the query (Q1) given next.

(Q1) *Find the names of people and the postcodes of the departments where they work.*

It can be written in SPARQL 1.1 as follows:

```
SELECT ?n ?p WHERE { ?x name ?n .
  { SELECT DISTINCT ?x ?p WHERE { ?x department ?d . ?d postcode ?p } } }.
```

Let us evaluate (Q1) over RDF graph  $G_{ex}$ . Variable  $?d$  is only visible within the subquery, whereas  $?x$  and  $?p$  are projected and hence visible in the outer query. Since a person can work in two departments with the same postcode (e.g., *Bob* works in *CS* and *Physics*, both located in *OX1*), the subquery uses *DISTINCT* to ensure that the result of the subquery does not contain duplicates. The evaluation of (Q1) over  $G_{ex}$  yields the multiset of mappings  $\mu_a = \{?n \mapsto \text{Alice}, ?p \mapsto \text{OX1}\}$ ,  $\mu_b = \{?n \mapsto \text{Bob}, ?p \mapsto \text{OX1}\}$ , and  $\mu_c = \{?n \mapsto \text{Charlie}, ?p \mapsto \text{OX2}\}$ , all with multiplicity 1. Omitting *DISTINCT* in the subquery would yield two copies of  $\mu_b$  in the evaluation.

To support queries such as (Q1), we extend our basic algebra Sparql by allowing *Project* and *Distinct* in patterns. After this extension, there is no longer a meaningful distinction between patterns and queries.

*Definition 3.6.* Language Sparql<sub>PD</sub> extends Sparql by allowing the query constructs *Project*( $X, P$ ) and *Distinct*( $P$ ) as patterns, called *subquery patterns*. The semantics of such patterns is the same as in the case of queries. The intermediate languages Sparql<sub>p</sub> and Sparql<sub>D</sub> allow only for *Project* and only for *Distinct* as patterns, respectively.

Language Sparql<sub>PD</sub> captures the query nesting functionality in SPARQL 1.1.

*Example 3.7.* Query (Q1) from Example 3.5 translates to  $\text{Sparql}_{PD}$  as follows:

$$\text{Project}(\{?n, ?p\}, \text{Join}((?x, \text{name}, ?n), \\ \text{Distinct}(\text{Project}(\{?x, ?p\}, \{(?x, \text{department}, ?d), (?d, \text{postcode}, ?p)\}))))).$$

### 3.3 Expressive Power of Subqueries as Patterns

At first sight, query nesting provides a great deal of power and flexibility to the language. It can lead to sophisticated interactions between set and bag semantics, which may be difficult (or, as we will soon see, impossible) to simulate within plain Sparql. Furthermore, subquery nesting can be arbitrarily deep, and it is reasonable to expect that each additional level of nesting may increase the expressive power of the language.

We next argue that  $\text{Sparql}_{PD}$  queries can be brought into a normal form where the nesting depth is bounded by two; thus, there is a natural bound on the level of nesting after which no further increase in expressive power can be achieved. This normal form is defined next, and one can check that query (Q1) satisfies its requirements.

*Definition 3.8.* A  $\text{Sparql}_{PD}$  query is in *s-normal form* if it has either the form  $\text{Distinct}(\text{Project}(X, P))$  with  $P$  subquery-free, or the form  $\text{Project}(X, P)$  where all subquery patterns in  $P$  are of the form  $\text{Distinct}(\text{Project}(X', P'))$  with  $P'$  subquery-free.

This normal form not only limits the nesting depth, but also restricts the ways in which *Project* and *Distinct* can be combined. If a query  $Q$  is in  $\text{Sparql}_P$  (or in  $\text{Sparql}_D$ ) and hence *Distinct* (respectively, *Project*) only occurs in the outermost level, then Definition 3.8 requires that pattern  $P$  is subquery-free and hence  $Q$  is a Sparql query.

We next show that each  $\text{Sparql}_{PD}$  query can be brought into s-normal form. The normalisation is based on two ideas that are illustrated in the following example.

*Example 3.9.* For the first idea, consider the  $\text{Sparql}_D$  query

$$\text{Distinct}(\text{Project}(\{?n\}, \text{Distinct}((?x, \text{name}, ?n)))).$$

Clearly, the inner application of *Distinct* is redundant as the outermost application of *Distinct* already ensures that the answer will contain no duplicates. Hence, the query can be simplified to  $\text{Distinct}(\text{Project}(\{?n\}, (?x, \text{name}, ?n)))$ .

For the second idea, consider the following modification of (Q1):

$$\text{Project}(\{?n, ?p\}, \text{Join}((?x, \text{name}, ?n), \\ \text{Project}(\{?x, ?p\}, \{(?x, \text{department}, ?n), (?n, \text{postcode}, ?p)\}))).$$

To remove the inner projection, it suffices to rename the variable  $?n$  in the *Project* subpattern so that it no longer clashes with any variable outside of the subpattern:

$$\text{Project}(\{?n, ?p\}, \{(?x, \text{name}, ?n), (?x, \text{department}, ?d), (?d, \text{postcode}, ?p)\}).$$

**THEOREM 3.10.** *Let  $\text{Sparql}_\star$  be one of  $\text{Sparql}_P$ ,  $\text{Sparql}_D$ , and  $\text{Sparql}_{PD}$ . Every query in  $\text{Sparql}_\star$  has an equivalent  $\text{Sparql}_\star$  query in s-normal form.*

**PROOF.** The first relevant observation is that all occurrences of *Distinct* in a  $\text{Sparql}_\star$  query  $Q$  that are in the scope of other *Distinct* subpatterns can be removed upfront without affecting the semantics. In the appendix we show the following claim.

**CLAIM 3.11.** *Any  $\text{Sparql}_{PD}$  query  $Q$  that has a subpattern  $\text{Distinct}(P)$  such that  $P$  has in turn a subpattern  $\text{Distinct}(P')$  is equivalent to the query obtained from  $Q$  by replacing  $\text{Distinct}(P')$  with  $P'$ .*

Second, *Project* can be seamlessly “pushed” upwards through all operators except *Distinct*. For instance,  $Join(Project(X, P_1), P_2)$  can be equivalently rewritten as  $Project(X \cup \text{var}(P_2), Join(P_1\theta, P_2))$  with  $\theta$  a renaming of  $\text{var}(P_1) \setminus X$ . Moreover, subsequent occurrences of *Project* can be merged, since the patterns  $Project(X_1 \cap X_2, P)$  and  $Project(X_1, Project(X_2, P))$  are equivalent.

To complete the proof, we note that every subpattern of  $Q$  of the form  $Distinct(P)$  can be extended with a “dummy” projection to  $Distinct(Project(\text{var}(P), P))$  without changing the semantics.  $\square$

From Theorem 3.10 and the definition of  $s$ -normal form it immediately follows that extending Sparql by allowing only *Project* or only *Distinct* as pattern constructs does not increase the expressive power of the language.

**COROLLARY 3.12.** *Languages Sparql<sub>P</sub> and Sparql<sub>D</sub> have the same expressive power as Sparql.*

Note that under set semantics, *Distinct* is redundant and hence subqueries as patterns do not add expressivity to Sparql. An obvious question is whether or not this is still the case under bag semantics.

Cohen [18] was the first to observe that query nesting in SQL can cause a complex interplay between bag and set semantics. In her paper, Cohen proposes a logical semantics for SQL queries where the aforementioned interplay between bags and sets occurs, and studies the query equivalence problem under such semantics. Cohen’s observations were used by Angles and Gutierrez [4] to conjecture that query nesting adds expressive power to SPARQL. Cohen’s paper, however, does not study questions related to expressive power and her results do not imply the conjecture by Angles and Gutierrez. In particular, Angles and Gutierrez claim that the query given in Example 3.13 below has no equivalent in Sparql. As we will see in Example 4.22, however, this query can be equivalently expressed in Sparql.

*Example 3.13 (Angles and Gutierrez [4]).* Consider the following query.

(Q2) *Find the names of people working for either the CS or the Math department.*

It can be written in SPARQL 1.1 as follows:

```
SELECT ?n WHERE { ?x name ?n .
  { SELECT DISTINCT *
    WHERE {?x department CS} UNION {?x department Math} } }.
```

Query (Q2) is translated to Sparql<sub>PD</sub> algebra as follows:

$$Project(\{?n\}, Join((?x, name, ?n), Distinct(Project(\{?x\}, P_u))))),$$

where  $P_u$  is  $Union((?x, department, CS), (?x, department, Math))$ .

We now settle the question whether query nesting in Sparql<sub>PD</sub> can be fully eliminated by showing that there exist Sparql<sub>PD</sub> queries that cannot be expressed in Sparql (and hence, by Corollary 3.12, in Sparql<sub>P</sub> or Sparql<sub>D</sub>).

**THEOREM 3.14.** *Language Sparql<sub>PD</sub> is strictly more expressive than Sparql.*

**PROOF.** We claim that the following Sparql<sub>PD</sub> query  $Q$  in  $s$ -normal form is not expressible in Sparql:

$$Project(\{?x, ?y\}, Join((?x, p, ?z), Distinct(Project(\{?y\}, (?y, q, ?u))))).$$

Assume for the sake of contradiction that there is a query  $Q'$  that is equivalent to  $Q$ . Consider graphs  $G_{m,n}$ , for  $m, n \geq 1$ , of the following form, where IRIs  $a, b_i, c$ , and  $d_j$  are chosen such that they do not occur in  $Q'$ :

$$\{ (a, p, b_1), \dots, (a, p, b_m), (c, q, d_1), \dots, (c, q, d_n) \}.$$

Query  $Q$  evaluates on every  $G_{m,n}$  to the multiset consisting of  $m$  copies of the mapping  $\mu = \{?x \mapsto a, ?y \mapsto c\}$ .

We claim that, in contrast, every Sparql query equivalent to  $Q$  modulo multiplicities evaluates on  $G_{m,n}$  to a multiset  $\Omega'$  such that either  $\text{card}_{\Omega'}(\mu) = 1$  or  $\text{card}_{\Omega'}(\mu) = m \cdot n$ . This contradicts our assumption that  $Q'$  is equivalent to  $Q$ . Intuitively, the contradiction arises because a Sparql query cannot distinguish between the different  $b_i$  and  $d_j$  and hence if a query of the form  $\text{Project}(X, P)$  returns  $\mu$  once, it must return its all  $m \cdot n$  copies. Of course, we can use *Distinct* in the query's outermost level, but then we obtain a single copy of  $\mu$  instead of the required  $m$  copies.

In what follows, we make this intuition precise. First, note that  $Q'$  cannot be of the form  $\text{Distinct}(\text{Project}(X, P))$  since such queries cannot yield mappings with multiplicity greater than one. Therefore,  $Q' = \text{Project}(X, P)$  for a set of variables  $X$  and pattern  $P$ . Let us consider the graph  $G_{1,2}$  and recall that  $Q$  evaluates to a single copy of  $\mu$  on  $G_{1,2}$ . Consequently, by assumption of equivalence between  $Q$  and  $Q'$ , we have that  $\llbracket Q' \rrbracket_{G_{1,2}}$  also consists of a single copy of  $\mu$ . Thus, by the structure of  $Q'$ , there exists  $\mu' \in \llbracket P \rrbracket_{G_{1,2}}$  such that  $\mu = \mu'|_X$ . We argue that there is a triple pattern  $T$  in  $P$  such that  $\mu'(T)$  is a triple in  $G_{1,2}$  mentioning IRI  $c$ ; this can be proved by an easy induction on the structure of  $P$  using the facts that  $\mu' \in \llbracket P \rrbracket_{G_{1,2}}$  and that  $\mu'$  mentions  $c$  (as it extends  $\mu$ ). Recall that, by construction of  $G_{m,n}$ , query  $Q'$  does not mention  $d_1$ ; hence,  $T$  has a variable  $?u$  at the object position and  $\mu'$  sends this variable to either  $d_1$  or  $d_2$ . The following claim shows the existence of a mapping in  $\llbracket P \rrbracket_{G_{1,2}}$  different from  $\mu'$ . Its proof proceeds by induction on the pattern structure and is given in the appendix.

**CLAIM 3.15.** *For every Sparql pattern  $P$  not mentioning  $d_1, d_2$  and every  $\mu' \in \llbracket P \rrbracket_{G_{1,2}}$ , the mapping  $\mu''$  such that  $\text{dom}(\mu'') = \text{dom}(\mu')$  and, for each  $?u \in \text{dom}(\mu')$ ,*

$$\mu''(?u) = \begin{cases} d_2 & \text{if } \mu'(?u) = d_1, \\ d_1 & \text{if } \mu'(?u) = d_2, \\ \mu'(?u) & \text{otherwise,} \end{cases}$$

*also belongs to  $\llbracket P \rrbracket_{G_{1,2}}$ .*

Claim 3.15 implies that  $\llbracket Q' \rrbracket_{G_{1,2}}$  contains at least two mappings, which contradicts our assumption that  $Q'$  is equivalent to  $Q$ . This concludes the proof of the theorem.  $\square$

#### 4 SUBQUERIES WITHIN EXPRESSIONS

Expressions in Sparql can only be constructed inductively from other expressions; that is, it is not possible for other constructs such as patterns or queries to occur within expressions. In SPARQL 1.1, however, the construct exists can be used to nest patterns within (possibly complex) expressions as illustrated next.

*Example 4.1.* Consider the following query (Q3).

(Q3) *Find the names of people not working in CS.*

It can be written in SPARQL 1.1 as follows:

```
SELECT ?n WHERE { ?x name ?n . FILTER NOT EXISTS {?x department CS} }.
```

Intuitively, query (Q3) should evaluate over our example graph  $G_{ex}$  to the single mapping  $\{?n \mapsto \text{Charlie}\}$ .

To support queries such as (Q3), the SPARQL 1.1 normative algebra introduces a dedicated construct exists. This construct can be syntactically incorporated into any extension of Sparql as given next.

*Definition 4.2.* For any Sparql extension Sparql<sub>★</sub>, the language Sparql<sub>★</sub><sup>∃</sup> further extends Sparql<sub>★</sub> by permitting expressions of the form exists( $P$ ) for a pattern  $P$ . Moreover, we define  $\text{var}(\text{exists}(P)) = \text{var}(P)$ .

Next we illustrate the translation of subqueries in expressions to Sparql<sup>∃</sup>.

*Example 4.3.* Query (Q3) from Example 4.1 translates to Sparql<sup>∃</sup> as follows:

$$\text{Project}(\{?n\}, \text{Filter}(\neg \text{exists}((?x, \text{department}, \text{CS}), (?x, \text{name}, ?n)))).$$

We discuss the formal semantics of exists expressions in the next two subsections.

#### 4.1 Normative Semantics of exists

The normative semantics of SPARQL 1.1 treats exists as a Boolean expression. Formally, it works with patterns normalised as follows: a pattern  $P$  is *rectified* if for each of its subpatterns  $\text{Project}(X, P')$ , no variable in  $\text{var}(P') \setminus X$  appears in  $P$  outside this subpattern. Intuitively, variables that occur in  $P'$  but do not belong to  $X$  are local for the pattern  $\text{Project}(X, P')$ . Every pattern  $P$  can be converted to an equivalent rectified pattern  $\bar{P}$  by renaming all local variables to fresh variables.

*Definition 4.4.* The (*normative*) *evaluation* of exists expressions over a graph  $G$  with respect to a mapping  $\mu$  is defined as follows, where  $\mu(\bar{P})$  is the result of replacing all variables  $?x \in \text{dom}(\mu)$  in  $\bar{P}$  by  $\mu(?x)$ :

$$\llbracket \text{exists}(P) \rrbracket_G^\mu = \begin{cases} \text{true} & \text{if } \llbracket \mu(\bar{P}) \rrbracket_G \neq \emptyset, \\ \text{false} & \text{otherwise.} \end{cases}$$

In other words, the evaluation of exists( $P$ ) for  $\mu$  and  $G$  is realised in four steps. First, the pattern  $P$  is rectified to  $\bar{P}$ .<sup>2</sup> Second, the mapping  $\mu$  is applied to  $\bar{P}$ , resulting in a new pattern  $\mu(\bar{P})$  where variables in  $\bar{P}$  have been replaced according to  $\mu$ . Third, the new pattern  $\mu(\bar{P})$  is evaluated against graph  $G$  to yield a multiset of mappings  $\llbracket \mu(\bar{P}) \rrbracket_G$ . Finally,  $\llbracket \text{exists}(P) \rrbracket_G^\mu$  is deemed *true* (respectively *false*) if  $\llbracket \mu(\bar{P}) \rrbracket_G$  is non-empty (respectively, empty).

*Example 4.5.* Consider query (Q3) from Examples 4.1 and 4.3. The triple pattern  $(?x, \text{name}, ?n)$  evaluates over  $G_{ex}$  to mappings

$$\mu_1 = \{?x \mapsto a, ?n \mapsto \text{Alice}\}, \mu_2 = \{?x \mapsto b, ?n \mapsto \text{Bob}\}, \text{ and } \mu_3 = \{?x \mapsto c, ?n \mapsto \text{Charlie}\},$$

all with multiplicity 1. The expression in the filter condition evaluates to *true* only for  $\mu_3$ . Specifically, when applying  $\mu_3$  to  $(?x, \text{department}, \text{CS})$  we obtain the ground triple pattern  $(c, \text{department}, \text{CS})$ , which evaluates to the empty multiset of mappings over  $G_{ex}$ . As a result, we have that  $\llbracket \text{exists}((?x, \text{department}, \text{CS})) \rrbracket_{G_{ex}}^{\mu_3}$  is *false* and its negation is therefore *true*.

Now, suppose the subpattern  $(?x, \text{department}, \text{CS})$  in query (Q3) is replaced by a subquery  $P = \text{Project}(\emptyset, (?x, \text{department}, \text{CS}))$ , resulting in the following Sparql<sup>∃</sup> query (Q3′):

$$\text{Project}(\{?n\}, \text{Filter}(\neg \text{exists}(\text{Project}(\emptyset, (?x, \text{department}, \text{CS})), (?x, \text{name}, ?n)))).$$

Since  $\text{var}(P) = \emptyset$  but variable  $?x$  occurs outside  $P$  in (Q3′), the rectified version of (Q3′) replaces  $?x$  in  $P$  by a fresh variable  $?y$ :

$$\text{Project}(\{?n\}, \text{Filter}(\neg \text{exists}(\text{Project}(\emptyset, (?y, \text{department}, \text{CS})), (?x, \text{name}, ?n)))).$$

Consequently, since  $\llbracket \text{exists}((?y, \text{department}, \text{CS})) \rrbracket_{G_{ex}}^{\mu_i}$  is *true* for all  $\mu_i$ , query (Q3′) evaluated over  $G_{ex}$  returns the empty multiset.

<sup>2</sup>The rectification step is missing in the SPARQL 1.1 specification itself, but should be done according to the errata document [31, query 8]. This reflects the intuition that renaming of local variables should not affect the semantics.

There are at least two entries in the SPARQL 1.1 errata document that ultimately relate to the semantics of exists (see errata queries 8 and 10 in [31]). The main issue behind both reports boils down to the fact that SPARQL 1.1 patterns are not closed under substitution; that is,  $\mu(\bar{P})$  may not be a syntactically valid pattern even if  $P$  is. In such cases, the semantics of exists is undefined.

*Example 4.6.* Consider the following SPARQL 1.1 query:

```
SELECT ?x
WHERE { ?x ?x ?x . FILTER EXISTS { ?y ?y ?y . FILTER BOUND(?x) } }.
```

The corresponding algebraic translation is as follows:

$$\text{Project}(\{?x\}, \text{Filter}(\text{exists}(\text{Filter}(\text{bound}(?x), (?y, ?y, ?y))), (?x, ?x, ?x))).$$

Given the graph  $\{(a, a, a)\}$ , one would intuitively expect the query to return the multiset containing a single mapping  $\mu = \{?x \mapsto a\}$  with multiplicity 1. The normative semantics is, however, ill-defined in this case: if we apply  $\mu$  to the pattern inside the exists expression, we obtain  $\text{Filter}(\text{bound}(a), (?y, ?y, ?y))$ , which is not a valid pattern since bound only accepts variables as arguments.

A similar problem can appear when the pattern inside an exists expression involves subqueries. Consider the SPARQL 1.1 query

```
SELECT ?x WHERE { ?x ?x ?x .
FILTER EXISTS { ?y ?y ?y . { SELECT ?x WHERE { (?x a ?y) } } } }
```

and its algebraic translation to  $\text{Sparql}_p^{\exists}$

$$\text{Project}(\{?x\}, \text{Filter}(\text{exists}(\text{Join}((?y, ?y, ?y), \text{Project}(\{?x\}, (?x, a, ?y)))), (?x, ?x, ?x))).$$

The normative evaluation of the query on the graph is also ill-defined as it requires evaluating  $\text{Project}(\{a\}, (a, a, ?y))$ , which is not a syntactically valid pattern.

Neither the SPARQL 1.1 specification nor the errata documents provide a fix for this issue and hence the normative semantics of exists is currently ill-defined.

In the following section, we propose a novel semantics for exists that provides a principled and elegant solution to the aforementioned issues.

## 4.2 Environment Semantics of exists

The semantics proposed in this section is a conservative extension of the normative one, in the sense that both semantics coincide whenever the normative semantics is well-defined. Unlike the normative specification, our semantics does not require any manipulation of syntactic objects. It is inspired by standard practices for defining denotational semantics of programming languages. The main idea is to make the semantics of patterns depend not only on the input graph, but also on a mapping, called *environment*. A pattern in an exists expression is then evaluated in the environment given by the mapping that the expression is tested on; in the absence of exists expressions, the environment remains empty.

*Definition 4.7.* The *environment evaluation*  $[E]_G^\mu$  of a  $\text{Sparql}_{pD}^{\exists}$  expression  $E$  over a graph  $G$  with respect to a mapping  $\mu$  is defined exactly the same as  $\llbracket E \rrbracket_G^\mu$  in Section 2.3 except that  $[E']_G^\mu$  is used instead of  $\llbracket E' \rrbracket_G^\mu$  for each participating subexpression  $E'$  of  $E$  (e.g.,  $[\text{bound}(?x)]_G^\mu = \llbracket \text{bound}(?x) \rrbracket_G^\mu$  and  $[E_1 \wedge E_2]_G^\mu$  is *true* if both  $[E_1]_G^\mu$  and  $[E_2]_G^\mu$  are *true*, it is *false* if  $[E_1]_G^\mu$  or  $[E_2]_G^\mu$  is *false*, and it is *error* otherwise), and for the case  $E = \text{exists}(P)$  it is defined as follows:

$$[\text{exists}(P)]_G^\mu = \begin{cases} \text{true} & \text{if } [P]_G^\mu \neq \emptyset, \\ \text{false} & \text{otherwise.} \end{cases}$$

The *environment evaluation*  $[P]_G^v$  of a Sparql<sub>PD</sub><sup>∃</sup> pattern  $P$  over a graph  $G$  in an environment  $v$  is defined as follows:

- $[B]_G^v$  for a BGP  $B$  is the multiset that consists of all  $\mu$ , with multiplicity 1, such that
  - $\mu \sim v$ ,
  - $\text{dom}(\mu) = \text{var}(B) \cup \text{dom}(v)$ , and
  - $\mu(B) \subseteq G$ ;
- $[\text{Project}(X, P')]_G^v = \{\mu \mid \mu' \in [P']_G^{v|_X}, \mu = \mu'|_X \cup v|_{\text{dom}(v) \setminus X}\}$ ;
- in all the remaining cases, the definition of  $[P]_G^v$  is the same as the definition of  $\llbracket P \rrbracket_G$  in Section 2.3 except that  $[P']_G^v$  and  $[E]_G^\mu$  are used instead of  $\llbracket P' \rrbracket_G$  and  $\llbracket E \rrbracket_G^\mu$ , respectively, for every participating subpattern  $P'$  and subexpression  $E$  of  $P$  (e.g.,  $[\text{Filter}(E, P')]_G^v = \{\mu \mid \mu \in [P']_G^v, [E]_G^\mu = \text{true}\}$ ).

Finally, the *environment evaluation*  $[Q]_G$  of a Sparql<sub>PD</sub><sup>∃</sup> query  $Q$  over a graph  $G$  is the evaluation  $[Q]_G^0$  of  $Q$ , seen as pattern, in the empty environment  $\emptyset$ .

These definitions require some discussion. First, note the difference to the normative semantics of  $\text{exists}(P)$ : instead of constructing a pattern  $\mu(\bar{P})$  and evaluating it over  $G$ , we now consider  $[P]_G^\mu$  where  $\mu$  is acting as the environment. Note also that for a BGP  $B$  we can equivalently define the semantics as

$$[B]_G^v = \{\mu \mid \mu' \in \llbracket B \rrbracket_G, \mu = \mu' \cup v\},$$

that is,  $[B]_G^v$  consists of the mappings in the evaluation of  $B$  without any environment, but extended by  $v$ ; if the extension is not possible, the mapping is just dropped.

Next, the definition of  $[\text{Project}(X, P)]_G^v$  reflects the intuition (mentioned in the errata document) that the variables that are projected out are local, and environments should not affect their values.

Finally, note that  $[P]_G^0 = \llbracket P \rrbracket_G$  for any Sparql<sub>PD</sub> pattern  $P$  and graph  $G$ , that is, for patterns without  $\text{exists}$  the normative and the environment semantics coincide.

*Example 4.8.* Consider the evaluation of query (Q3) from Examples 4.1 and 4.3 on  $G_{ex}$  using our proposed semantics. As in Example 4.5,  $[(?x, \text{name}, ?n)]_{G_{ex}}^0$  evaluates to mappings  $\mu_1, \mu_2$  and  $\mu_3$ . The expression in the filter condition evaluates to *true* only on

$$\mu_3 = \{?x \mapsto c, ?n \mapsto \text{Charlie}\}.$$

Specifically,  $[(?x, \text{department}, \text{CS})]_{G_{ex}}^{\mu_3}$  is the empty multiset, since the only mapping with domain  $\{?x, ?n\}$  compatible with  $\mu_3$  is precisely  $\mu_3$ , and the application of  $\mu_3$  to  $(?x, \text{department}, \text{CS})$  does not belong to  $G_{ex}$ .

*Example 4.9.* Consider the problematic query in Example 4.6 and its evaluation over the graph  $G = \{(a, a, a)\}$ . To this end, let us compute the evaluation  $[\text{Filter}(\text{exists}(P), (?x, ?x, ?x))]_G^0$  for  $P = \text{Filter}(\text{bound}(?x), (?y, ?y, ?y))$ . The only mapping in  $[(?x, ?x, ?x)]_G^0$  is  $\mu = \{?x \mapsto a\}$  (whose multiplicity is 1). Next we show that  $\mu$  is an answer to the whole query as well. By our semantics of  $\text{exists}$ ,  $\mu$  “survives” the filter condition only if  $[P]_G^\mu$  is non-empty. This holds because  $[P]_G^\mu$  consists of a single mapping  $\mu' = \{?x \mapsto a, ?y \mapsto a\}$ , which extends  $\mu$  and has  $[\text{bound}(?x)]_G^{\mu'} = \text{true}$ .

We conclude this section with a natural property of the environment semantics of Sparql<sub>PD</sub><sup>∃</sup> and, as we will see later, of every other language considered in this paper.

**PROPOSITION 4.10.** *Let  $P$  be a Sparql<sub>PD</sub><sup>∃</sup> pattern,  $G$  a graph, and  $v$  an environment. Then  $v \subseteq \mu$  for every  $\mu \in [P]_G^v$ .*

The inductive proof of this proposition is given in the appendix.

*Example 4.11.* Let  $P = (?x, name, ?n)$  and  $v = \{?x \mapsto c\}$ . Then  $[P]_{G_{ex}}^v$  consists of a single mapping  $\{?x \mapsto c, ?n \mapsto \text{Charlie}\}$ , which is a superset of  $v$ . The normative semantics, on the other hand, does not satisfy Proposition 4.10:  $\llbracket v(P) \rrbracket_{G_{ex}}$  consists of a single mapping  $\{?n \mapsto \text{Charlie}\}$ , which does not bind  $?x$ . In fact, for every pattern  $P$ , mapping  $v$ , graph  $G$  and mapping  $\mu \in \llbracket v(P) \rrbracket_G$ , we have  $\text{dom}(\mu) \cap \text{dom}(v) = \emptyset$ .

### 4.3 Comparison between Normative and Environment Semantics

Unlike the normative semantics, our proposed semantics is well-defined for all  $\text{Sparql}_{PD}^{\exists}$  patterns. We next show that our semantics is a conservative extension of the normative one in that the two coincide on all patterns and graphs for which the normative semantics is *well-defined*, that is, for which no ill-formed patterns appear in the process of evaluation.

The following lemma formally relates the two semantics, taking into account the differences observed in Example 4.11. For the lemma, it is convenient to generalise the notion of *rectification* to expressions:  $\bar{E}$  for an expression  $E$  is defined exactly the same as  $\bar{P}$  for a pattern  $P$ .

**LEMMA 4.12.** *Let  $G$  be a graph,  $E$  a  $\text{Sparql}_{PD}^{\exists}$  expression with rectification  $\bar{E}$ , and  $P$  a  $\text{Sparql}_{PD}^{\exists}$  pattern with rectification  $\bar{P}$ .*

1. *For every two mappings  $\mu, \nu$  such that the variables introduced by the rectification  $\bar{E}$  of  $E$  are not in  $\text{dom}(\mu) \cup \text{dom}(\nu)$  and  $\llbracket v(\bar{E}) \rrbracket_G^\mu$  is well-defined, if  $\mu \sim \nu$  then  $\llbracket v(\bar{E}) \rrbracket_G^\mu = [E]_G^{\mu \cup \nu}$ .*
2. *For every two mappings  $\mu, \nu$  such that the variables introduced by the rectification  $\bar{P}$  of  $P$  are not in  $\text{dom}(\nu)$  and  $\llbracket v(\bar{P}) \rrbracket_G$  is well-defined,*
  - (a) *if  $\mu \in \llbracket v(\bar{P}) \rrbracket_G$  then  $\mu \cup \nu \in [P]_G^v$  and  $\text{card}_{\llbracket v(\bar{P}) \rrbracket_G}(\mu) = \text{card}_{[P]_G^v}(\mu \cup \nu)$ ;*
  - (b) *if  $\mu \in [P]_G^v$  then  $\mu|_{\text{dom}(\mu) \setminus \text{dom}(\nu)} \in \llbracket v(\bar{P}) \rrbracket_G$  and  $\text{card}_{[P]_G^v}(\mu) = \text{card}_{\llbracket v(\bar{P}) \rrbracket_G}(\mu|_{\text{dom}(\mu) \setminus \text{dom}(\nu)})$ .*

The proof of this lemma, which proceeds by simultaneous induction on the structure of  $E$  and  $P$ , is given in the appendix. Note that, for the induction to go through, the lemma makes additional assumptions on the domain of environments, which are trivially satisfied by the empty environment.

We are now ready to prove that the two semantics are equivalent.

**THEOREM 4.13.** *For any  $\text{Sparql}_{PD}^{\exists}$  pattern  $P$  and graph  $G$  such that  $\llbracket P \rrbracket_G$  is well-defined it holds that  $\llbracket P \rrbracket_G = [P]_G^0$ .*

**PROOF.** By definition,  $\llbracket \bar{P} \rrbracket_G = \llbracket P \rrbracket_G$ . Statement 2(a) in Lemma 4.12 applied to the empty environment implies that every mapping in  $\llbracket P \rrbracket_G$  belongs to  $[P]_G^0$  with the same multiplicity. Statement 2(b) in the same lemma implies that every mapping in  $[P]_G^0$  belongs to  $\llbracket P \rrbracket_G$  also with the same multiplicity.  $\square$

In the rest of the paper we use (often silently) the environment semantics (i.e.,  $[\cdot]$ ) instead of the normative one (i.e.,  $\llbracket \cdot \rrbracket$ ). Note, however, that the notion of pattern equivalence is now parametrised by an environment.

**Definition 4.14.** Given a  $\text{Sparql}$  extension  $\text{Sparql}_\star$  and an environment  $\nu$ ,  $\text{Sparql}_\star$  patterns  $P_1$  and  $P_2$  are  $\nu$ -equivalent if  $[P_1]_G^v = [P_2]_G^v$  for any graph  $G$ .

The presence of environments may have non-trivial consequences. For example, the proof of Lemma 3.3 on the expressibility of *SetMinus* works only for the empty environment: the justifying rewriting involves the renaming  $P'_2\theta$ , whose evaluation  $[P'_2\theta]_G^v$  may be different, for some nonempty  $\nu$ , from  $[P'_2]_G^v$  renamed according to  $\theta$ , because  $\nu$  may interact differently with the mappings appearing during the evaluation of  $P'_2$  and  $P'_2\theta$ . Therefore, we shall be careful when exploiting

statements such as Lemma 3.3 later on and make sure that we never apply them in the scope of an exists expression.

Note also that the natural requirements we impose on extensions of  $\text{Sparql}^{\exists}$  apply only to the empty environment: for any pattern or query  $S$  and any graph  $G$

- $\text{dom}(\mu) \subseteq \text{var}(S)$  for any  $\mu \in [S]_G^0$ , and
- for any renaming  $\theta$  of  $\text{var}(S)$ ,  $\mu\theta \in [S\theta]_G^0$  if and only if  $\mu \in [S]_G^0$ .

#### 4.4 Expressive Power of exists

The exists construct seems rather powerful as it makes the languages of patterns and expressions mutually recursive. Furthermore, expressions can occur not only as parameters of *Filter*, but also in *LeftJoin* patterns. As we show next, however, exists provides no additional expressive power, since it can be simulated by other constructs.

We proceed according to the following three steps.

- (1) We show that exists can be eliminated from *LeftJoin* patterns. This is achieved by “pushing” complex expressions from *LeftJoin* into *Filter* patterns.
- (2) We prove that any *Filter* pattern can be expressed in terms of exists-free *Filter* patterns and patterns of the form  $\text{Filter}(\text{exists}(P_1), P_2)$  or  $\text{Filter}(\neg\text{exists}(P_1), P_2)$ .
- (3) We show that patterns of the form  $\text{Filter}(\text{exists}(P_1), P_2)$  and  $\text{Filter}(\neg\text{exists}(P_1), P_2)$  can be rewritten in terms of exists-free patterns.

The following example demonstrates the ideas behind these steps.

*Example 4.15.* Consider the  $\text{Sparql}^{\exists}$  pattern  $P_0$  defined as follows:

$$\text{LeftJoin}(\text{exists}((?x, \text{role}, \text{Professor})) \vee (?y > 4000), (?x, \text{name}, ?n), (?x, \text{salary}, ?y)).$$

Intuitively, the pattern returns the names of people as well as their salaries, where salaries are only reported for professors and people earning more than 4000. Evaluated over  $G_{ex}$ , the pattern returns three mappings:  $\mu_1 = \{?x \mapsto a, ?n \mapsto \text{Alice}, ?y \mapsto 4000\}$ ,  $\mu_2 = \{?x \mapsto b, ?n \mapsto \text{Bob}, ?y \mapsto 5000\}$ , and  $\mu_3 = \{?x \mapsto c, ?n \mapsto \text{Charlie}\}$ .

We begin by demonstrating how to move the join condition of  $P_0$  into *Filter* patterns. Let  $P'_1$  be the following pattern, capturing people whose salary is to be reported:

$$\text{Filter}(\text{exists}((?x, \text{role}, \text{Professor})) \vee (?y > 4000), \{(?x, \text{name}, ?n), (?x, \text{salary}, ?y)\}).$$

Then  $P_0$  can be equivalently restated as the following pattern  $P_1$ :

$$\text{Union}(P'_1, \text{Filter}(\neg\text{bound}(?u), \text{LeftJoin}(\text{true}, (?x, \text{name}, ?n), \text{Join}(P'_1, (?u, ?v, ?w))))).$$

Next, we rewrite  $P_1$  to an equivalent pattern  $P_2$  in which all occurrences of exists are of the form  $\text{Filter}(\neg\text{exists}(P'), P'')$ . For  $P_1$ , it suffices to rewrite  $P'_1$  to  $P'_2 = \text{Union}(P'_{21}, P'_{22})$  where

$$P'_{21} = \text{Filter}(\text{exists}((?x, \text{role}, \text{Professor})), \{(?x, \text{name}, ?n), (?x, \text{salary}, ?y)\}),$$

$$P'_{22} = \text{Filter}(\neg\text{exists}((?x, \text{role}, \text{Professor})), \text{Filter}(?y > 4000, \{(?x, \text{name}, ?n), (?x, \text{salary}, ?y)\})).$$

Pattern  $P_2$  is then obtained from  $P_1$  by replacing all occurrences of  $P'_1$  with  $P'_2$ .

Finally, we eliminate exists from  $P_2$ . Consider first the negative occurrence in  $P'_{22}$ . It can be rewritten using *Filter*, *bound*, and *LeftJoin*, yielding the pattern

$$\begin{aligned} &\text{Filter}(\neg\text{bound}(?u'), \text{Filter}(?y > 4000, \\ &\quad \text{LeftJoin}(\text{true}, \{(?x, \text{name}, ?n), (?x, \text{salary}, ?y)\}, \\ &\quad \{(?x, \text{name}, ?n), (?x, \text{salary}, ?y), (?x, \text{role}, \text{Professor}), (?u', ?v', ?w')\}))). \end{aligned}$$

The positive occurrence of exists in  $P'_{21}$  can, in turn, be rewritten using *SetMinus* in terms of its negation yielding the pattern

$$\text{SetMinus}(\{(?x, \text{name}, ?n), (?x, \text{salary}, ?y)\}, \\ \text{Filter}(\neg \text{exists}((?x, \text{role}, \text{Professor})), \{(?x, \text{name}, ?n), (?x, \text{salary}, ?y)\})).$$

This negative occurrence can then be eliminated analogously to the one in  $P'_{22}$ . By replacing  $P'_{21}$  and  $P'_{22}$  by their exists-free counterparts, we obtain an exists-free pattern equivalent to  $P_0$ .

We show correctness of the three steps only with respect to the empty environment. Thus, in the presence of nested exists expressions, these steps must be applied iteratively in a top-down manner. To this end, we should also guarantee that no step increases the nesting depth of exists (and the third step decreases it by one). We refer to subpatterns of patterns not in the scope of any exists expressions as *top-level*.

The first step is justified by the following lemma.

LEMMA 4.16. *Let  $\text{Sparql}_\star$  be any Sparql extension. Every  $\text{Sparql}_\star^{\exists}$  pattern is  $\emptyset$ -equivalent to a  $\text{Sparql}_\star^{\exists}$  pattern with the same nesting depth of exists all of whose top-level *LeftJoin* subpatterns are of the form  $\text{LeftJoin}(\text{eq}(X, \theta), P'_1, P'_2)$ .*

PROOF. Consider an arbitrary  $\text{Sparql}_\star^{\exists}$  pattern of the form  $\text{LeftJoin}(E, P_1, P_2)$ . We first construct  $P$  such that, for any graph  $G$ ,  $[P]_G^0$  captures (with possibly incorrect multiplicities) all the mappings  $\mu_1$  in  $[P_1]_G^0$  having a compatible  $\mu_2$  in  $[P_2]_G^0$  such that  $[E]_G^{\mu_1 \cup \mu_2}$  is *true*, extended with a “certificate” in the form of a possible compatible extension. Such a  $P$  can be defined as follows, where  $X = \text{var}(P_1) \cup \text{var}(P_2)$ , and  $\theta_1$  is a renaming of  $X$ :

$$\text{Filter}(E\theta_1, \text{Join}(\text{Filter}(\text{eq}(X, \theta_1), \text{Join}(P_1, P_1\theta_1)), P_2\theta_1)).$$

More formally, by construction, we have that  $\mu_1 \in [P_1]_G^0$ ,  $\mu_2 \in [P_2]_G^0$ ,  $\mu_1 \sim \mu_2$ , and  $[E]_G^{\mu_1 \cup \mu_2} = \text{true}$  if and only if  $\mu_1 \cup (\mu_1 \cup \mu_2)\theta_1 \in [P]_G^0$ .

We then construct  $P'$  such that  $[P']_G^0$  captures, with correct multiplicities, all mappings  $\mu_1 \in [P_1]_G^0$  that do not have any corresponding  $\mu_2$ . For this, we use the following construction, which involves fresh variables  $?x, ?y, ?z$  and another renaming  $\theta_2$  of  $X$ :

$$P' = \text{Filter}(\neg \text{bound}(?x), \text{LeftJoin}(\text{eq}(X, \theta_2), P_1, P^x\theta_2)),$$

where  $P^x = \text{Join}(P, (?x, ?y, ?z))$ . By construction,  $\mu_1 \in [P']_G^0$  if and only if  $\mu_1 \in [P_1]_G^0$  and there is no  $\mu_2$  such that  $\mu_1 \cup (\mu_1 \cup \mu_2)\theta_1 \in [P]_G^0$  (recall that  $G$  is non-empty, so  $(?x, ?y, ?z)$  always matches); moreover, whenever  $\mu_1 \in [P']_G^0$  it holds that  $\text{card}_{[P']_G^0}(\mu_1) = \text{card}_{[P_1]_G^0}(\mu_1)$ . In other words,  $P'$  captures the second operand of the bag union in the definition of *LeftJoin*.

Consequently, the pattern  $\text{LeftJoin}(E, P_1, P_2)$  is equivalent to  $\text{Union}(\text{Filter}(E, \text{Join}(P_1, P_2)), P')$ .  $\square$

In the second step, we show that all patterns involving exists can be reduced to patterns where exists appears in expressions only at the top or under a single negation.

LEMMA 4.17. *Let  $\text{Sparql}_\star$  be any Sparql extension. Every  $\text{Sparql}_\star^{\exists}$  pattern is  $\emptyset$ -equivalent to a  $\text{Sparql}_\star^{\exists}$  pattern with the same nesting depth of exists each of whose top-level *Filter* subpatterns involving exists has the form  $\text{Filter}(\text{exists}(P_2), P_1)$  or  $\text{Filter}(\neg \text{exists}(P_2), P_1)$ .*

PROOF. Consider a  $\text{Sparql}_\star^{\exists}$  pattern  $\text{Filter}(E, P)$ , where  $E$  contains a top-level occurrence of an expression  $\text{exists}(P')$  (i.e., an occurrence that is not inside any other exists expression). Since

exists( $P'$ ) must evaluate to either *true* or *false*, we can capture each possibility by replacing exists( $P'$ ) in  $E$  by the respective truth value to obtain the equivalent pattern

$$\text{Union}(\text{Filter}(\text{exists}(P'), \text{Filter}(E[\text{true}], P)), \text{Filter}(\neg \text{exists}(P'), \text{Filter}(E[\text{false}], P))),$$

where  $E[E']$  is  $E$  with exists( $P'$ ) replaced by  $E'$ .  $\square$

For the final step, we first show that patterns of the form  $\text{Filter}(\neg \text{exists}(P_2), P_1)$  can be expressed in Sparql by means of a structural transformation  $\text{Sub}_\theta(P_2, P_1)$ , for  $\theta$  a renaming of  $\text{var}(P_1)$ . Intuitively,  $\text{Sub}_\theta(P_2, P_1)$  joins each subpattern of  $P_2$  with a fresh copy of  $P_1$  in such a way that every mapping in the evaluation of the transformed pattern over a graph  $G$  has the form  $\mu_2 \cup \mu_1 \theta \cup \dots$ , for  $\mu_2 \in [P_2]_G^{\mu_1}$  and  $\mu_1 \in [P_1]_G^0$ . The construction of  $\text{Sub}_\theta(P_2, P_1)$  is complicated by the fact that  $[P_1]_G^0$  may contain distinct but pairwise compatible mappings. Thus, to ensure that all subpatterns of  $P_2$  are indeed evaluated in the same environment, the transformation employs techniques previously seen in the rewriting of the operator *SetMinus* in Lemma 3.3 and in the first step of exists elimination in Lemma 4.16.

Once we have a rewriting for  $\text{Filter}(\neg \text{exists}(P_2), P_1)$ , patterns of the form  $\text{Filter}(\text{exists}(P_2), P_1)$  are reduced to  $\text{Filter}(\neg \text{exists}(P_2), P_1)$  using *SetMinus*.

LEMMA 4.18. *Let Sparql $_\star$  be one of Sparql, Sparql $_p$ , Sparql $_D$  and Sparql $_{pD}$ . Any Sparql $_\star^{\exists}$  pattern with all top-level subpatterns with exists of the form  $\text{Filter}(\neg \text{exists}(P_2), P_1)$  or  $\text{Filter}(\text{exists}(P_2), P_1)$  is  $\emptyset$ -equivalent to a Sparql $_\star^{\exists}$  pattern with a smaller nesting depth of exists.*

PROOF. We begin with  $\text{Filter}(\neg \text{exists}(P_2), P_1)$ . We will use an abbreviation

$$\text{Copy}_\theta(P') = \text{Filter}(\text{eq}(\text{var}(P'), \theta), \text{Join}(P', P'\theta)),$$

for a pattern  $P'$  and renaming  $\theta$  of  $\text{var}(P')$ . Intuitively, pattern  $\text{Copy}_\theta(P')$  duplicates, for each mapping in the answer to  $P'$ , the values of variables of  $P'$  and stores them in fresh copies  $(\text{var}(P'))\theta$  (note that the multiplicities of mappings with duplicated values may be different from the original, but this is inessential, as we will see later). We also use the following abbreviations, for patterns  $P^1$  and  $P^2$ , a set of variables  $Y$ , renamings  $\theta^1$  and  $\theta^2$  of  $Y$ , and an expression  $E'$ :

$$\begin{aligned} \text{EqJoin}_{\theta^1, \theta^2}^Y(P^1, P^2) &= \text{Filter}(\text{eq}(Y\theta^1, (\theta^1)^{-1} \circ \theta^2), \text{Join}(P^1, P^2)), \\ \text{EqLeftJoin}_{\theta^1, \theta^2}^Y(E, P^1, P^2) &= \text{LeftJoin}(E \wedge \text{eq}(Y\theta^1, (\theta^1)^{-1} \circ \theta^2), P^1, P^2), \end{aligned}$$

where  $(\theta^1)^{-1}$  is the inverse of  $\theta^1$ , that is the renaming of the range of  $\theta^1$  such that  $?x((\theta^1)^{-1}) = ?y$  for each  $?y \in Y$  with  $?y\theta^1 = ?x$  (this renaming is well-defined because  $\theta^1$  is injective by definition), and  $(\theta^1)^{-1} \circ \theta^2$  is the composition of  $(\theta^1)^{-1}$  and  $\theta^2$ , that is, the renaming of the range of  $\theta^1$  such that  $?x((\theta^1)^{-1} \circ \theta^2) = (?x(\theta^1)^{-1})\theta^2$  for each  $?x$  in the domain of  $(\theta^1)^{-1}$ . Intuitively, pattern  $\text{EqJoin}_{\theta^1, \theta^2}^Y(P^1, P^2)$  yields all mappings in  $\text{Join}(P^1, P^2)$  that coincide on the variables in  $Y\theta^1$  and  $Y\theta^2$ , respectively. Pattern  $\text{EqLeftJoin}_{\theta^1, \theta^2}^Y(E, P^1, P^2)$  is an analogue of  $\text{EqJoin}_{\theta^1, \theta^2}^Y(P^1, P^2)$  for  $\text{LeftJoin}(E, P^1, P^2)$ .

Given patterns  $P$  and  $P_1$  and a renaming  $\theta$  of  $\text{var}(P_1)$ , we define  $\text{Sub}_\theta(P, P_1)$  by induction on the structure of  $P$  as follows, where  $\theta^1$  and  $\theta^2$  are fresh renamings of  $\text{var}(P_1)$ , and  $\theta_Y$  is a fresh

renaming of  $Y$  for every recursive application of  $\text{Sub}_\theta$ :

$$\begin{aligned}
\text{Sub}_\theta(B, P_1) &= \text{Join}(B, \text{Copy}_\theta(P_1)), \text{ for a BGP } B, \\
\text{Sub}_\theta(\text{Join}(P^1, P^2), P_1) &= \text{EqJoin}_{\theta, \theta^1}^{\text{var}(P_1)}(\text{EqJoin}_{\theta^1, \theta^2}^{\text{var}(P_1)}(\text{Sub}_{\theta^1}(P^1, P_1), \text{Sub}_{\theta^2}(P^2, P_1)), P_1\theta), \\
\text{Sub}_\theta(\text{Union}(P^1, P^2), P_1) &= \text{Union}(\text{Sub}_\theta(P^1, P_1), \text{Sub}_\theta(P^2, P_1)), \\
\text{Sub}_\theta(\text{Filter}(E, P'), P_1) &= \text{Filter}(E, \text{Sub}_\theta(P', P_1)), \\
\text{Sub}_\theta(\text{LeftJoin}(E, P^1, P^2), P_1) &= \\
&\quad \text{EqJoin}_{\theta, \theta^1}^{\text{var}(P_1)}(\text{EqLeftJoin}_{\theta^1, \theta^2}^{\text{var}(P_1)}(E, \text{Sub}_{\theta^1}(P^1, P_1), \text{Sub}_{\theta^2}(P^2, P_1)), P_1\theta), \\
\text{Sub}_\theta(\text{Project}(X, P'), P_1) &= \text{Project}(X \cup \text{var}(P_1) \cup (\text{var}(P_1))\theta, \text{Sub}_\theta(P'\theta_{\text{var}(P_1)\setminus X}, P_1)), \\
\text{Sub}_\theta(\text{Distinct}(P'), P_1) &= \text{Distinct}(\text{Sub}_\theta(P', P_1)).
\end{aligned}$$

As formalised in the following claim, in the empty environment  $\text{Sub}_\theta(P, P_1)$  always gives the same answers as  $P$ , except that these answers have some additional variables defined; for example, the answers to  $P_1$  are stored in  $(\text{var}(P_1))\theta$ . Also, importantly, multiplicities are inessential in this construction, because everything we need is to preserve emptiness and non-emptiness of the set of answers.

**CLAIM 4.19.** *For any Sparql $_\star$  patterns  $P$  and  $P_1$ , graph  $G$ , and mapping  $\mu_1 \in [P_1]_G^0$ , a mapping  $\mu$  belongs to  $[P]_G^{\mu_1}$  if and only if there is a mapping  $\mu_*$  in  $[\text{Sub}_\theta(P, P_1)]_G^0$  for a renaming  $\theta$  of  $\text{var}(P_1)$  such that  $\mu_*|_{\text{var}(P_1) \cup \text{var}(P)} = \mu$  and  $\mu_*|_{(\text{var}(P_1))\theta} = \mu_1\theta$ .*

The claim is shown in the appendix by induction on the structure of  $P$ .

The key corollary of Claim 4.19 is that  $[P]_G^{\mu_1}$  is empty for some  $\mu_1 \in [P_1]_G^0$  if and only if  $[\text{Sub}_\theta(P, P_1)]_G^0$  does not have any mapping with a copy of  $\mu_1$  in  $(\text{var}(P_1))\theta$ . Therefore, a pattern  $\text{Filter}(\neg\text{exists}(P_2), P_1)$ , which tests the evaluation of  $P_2$  for emptiness in answers to  $P_1$  as environments, is  $\emptyset$ -equivalent to the following pattern, where  $?x, ?y$  and  $?z$  are fresh variables, and  $\theta$  is a renaming of  $\text{var}(P_1)$ :

$$\text{Filter}(\neg\text{bound}(?x), \text{LeftJoin}(\text{eq}(\text{var}(P_1), \theta), P_1, \text{Join}(\text{Sub}_\theta(P_2, P_1), (?x, ?y, ?z)))).$$

Finally, patterns of the form  $\text{Filter}(\text{exists}(P_2), P_1)$  are dealt with by rewriting them to  $\emptyset$ -equivalent patterns  $\text{SetMinus}(P_1, \text{Filter}(\neg\text{exists}(P_2), P_1))$ .  $\square$

Using Lemmas 4.16–4.18 we can now establish that the  $\text{exists}$  construct does not add any expressive power to any language between Sparql and Sparql $_{PD}$ .

**THEOREM 4.20.** *Let Sparql $_\star$  be one of the languages Sparql, Sparql $_P$ , Sparql $_D$ , and Sparql $_{PD}$ . Then Sparql $_\star^{\exists}$  has the same expressive power as Sparql $_\star$ .*

From this theorem and Corollary 3.12 we can conclude that neither  $\text{Project}$  nor  $\text{Distinct}$  in isolation add any expressive power in the presence of  $\text{exists}$ .

**COROLLARY 4.21.** *Languages Sparql $_P^{\exists}$ , Sparql $_D^{\exists}$ , and Sparql $^{\exists}$  have the same expressive power as Sparql.*

We conclude this section with the translation of the query in [4] to Sparql.

**Example 4.22.** Recall query (Q2) from Example 3.13. It is straightforward to check that (Q2) can be expressed using the  $\text{exists}$  operator as follows:

$$\text{Project}(\{?n\}, \text{Filter}(\text{exists}(P_u), (?x, \text{name}, ?n))).$$

But then, applying Lemma 4.18 and some simplifications, we establish that (Q2) is equivalent to the Sparql query

$$\text{Project}(\{?n\}, \text{SetMinus}((?x, \text{name}, ?n), P)),$$

where  $P$  is

$$\text{Filter}(\neg \text{bound}(?u), \text{LeftJoin}(\text{true}, (?x, \text{name}, ?n), \text{Join}(P_u, (?u, ?v, ?w)))).$$

## 5 VARIABLE ASSIGNMENT TO EXPRESSIONS

In addition to retrieving data, many applications require the ability to perform numeric computations. SQL provides a wide range of constructs to this effect: on the one hand, it allows for Boolean and arithmetic expressions for computing new data values, which can subsequently be assigned to variables; on the other hand, it is equipped with powerful constructs for grouping and aggregation. Formalising these features requires a significant extension to the relational algebra, which involves *grouping* and *generalised projection* operators, as well as *aggregate functions* (e.g., see [25, Chapter 5] for an excellent introduction).

The original SPARQL 1.0 recommendation, however, did not provide any such features. Although arithmetic expressions were available, computed values could only be used as part of filter conditions; thus, the means for assigning such values to variables and subsequently return them in query answers was missing. Similarly, SPARQL 1.0 did not provide any support for grouping and aggregation, which limited its applicability in many practical scenarios.

The standardisation of SPARQL 1.1 addressed these limitations. As in SQL, introducing these features into the language required an extended algebra, which, however, turned out rather unconventional when compared to the algebra underpinning SQL.

Our aim in this and the next sections is to provide an in-depth formal analysis of the SPARQL 1.1 assignment and aggregation algebra, which (to the best of our knowledge) has not been studied in the literature. In particular, in this section we study the *Extend* operator, which captures the SPARQL 1.1 constructs BIND and VALUES, and provides the algebraic means for assigning variables to expressions. In Section 5.1 we formalise the assignment operator, and in Section 5.2 we look at the expressive power of the languages with this operator.

### 5.1 Syntax and Semantics of Assignment

We give a semantics of *Extend* that correctly deals with expressions involving exists and analyse the expressive power of the variants of Sparql discussed so far when enriched with *Extend*.

*Example 5.1.* Consider the following query.

(Q4) *Return people's names with their salaries after 20% tax together with flags indicating whether they work in the CS department.*

This query can be written in SPARQL 1.1 as follows:

```
SELECT ?n (0.8 * ?s AS ?t) ?c
WHERE {?x name ?n . ?x salary ?s . ?x department ?d BIND (?d=CS AS ?c)}.
```

Over graph  $G_{ex}$ , query (Q4) yields mappings such as  $\{?n \mapsto \text{Alice}, ?t \mapsto 3200, ?c \mapsto \text{true}\}$ , indicating Alice's net salary and the fact that she works in the CS department.

To support such queries, the SPARQL 1.1 algebra provides the *Extend* operator.

*Definition 5.2.* Given any extension  $\text{Sparql}_\star$  of Sparql, the language  $\text{Sparql}_{\star E}$  further extends  $\text{Sparql}_\star$  by permitting patterns of the form  $\text{Extend}(?x, E, P)$ , where  $?x$  is a variable not in  $\text{var}(P)$ ,  $E$  is an expression, and  $P$  is a pattern. Let also  $\text{var}(\text{Extend}(?x, E, P)) = \{?x\} \cup \text{var}(P)$ .

The *evaluation* of a pattern  $Extend(?x, E, P)$  over a graph  $G$  in an environment  $\nu$  is defined as follows:

$$\begin{aligned} [Extend(?x, E, P)]_G^\nu = \\ \llbracket \mu \mid \mu' \in [P]_G^\nu, [E]_G^{\mu'} \neq error, \mu = \mu' \cup \{?x \mapsto [E]_G^{\mu'}\} \rrbracket \uplus \llbracket \mu \mid \mu \in [P]_G^\nu, [E]_G^\mu = error \rrbracket. \end{aligned}$$

Intuitively,  $Extend$  assigns to variable  $?x$  the evaluation of  $E$  in each solution mapping to  $P$ , provided the evaluation does not yield *error*. Solution mappings in which  $E$  evaluates to *error* are also included, but with  $?x$  unbound.

*Example 5.3.* Query (Q4) is translated into the normative algebra as follows:

$$\begin{aligned} Project(\{?n, ?t, ?c\}, Extend(?t, 0.8 * ?s, \\ Extend(?c, ?d = CS, \{(?x, name, ?n), (?x, salary, ?s), (?x, department, ?d)\}))). \end{aligned}$$

Having defined the semantics, we next discuss the interactions between  $Extend$  patterns and exists expressions. To this end, we first illustrate that the issues with the normative semantics described in Section 4.1 are exacerbated in the presence of variable assignment. To make our argument precise, we need to formally extend the normative semantics to the new operator:  $\llbracket Extend(?x, E, P) \rrbracket_G$  is defined in exactly the same way as the environment semantics  $[Extend(?x, E, P)]_G^\nu$  in Definition 5.2, except that  $\llbracket P \rrbracket_G$  and  $\llbracket E \rrbracket_G^\mu$  are used instead of  $[P]_G^\nu$  and  $[E]_G^\mu$ , respectively.

*Example 5.4.* Consider the SPARQL 1.1 query

```
SELECT ?x WHERE { ?x a b .
  FILTER EXISTS { ?x c d . { ?y ?y ?y BIND (?y AS ?x) } } }
```

and its algebraic translation to  $\text{Sparql}_E^{\exists}$

$$Project(\{?x\}, Filter(exists(Join((?x, c, d), Extend(?x, ?y, (?y, ?y, ?y)))), (?x, a, b))).$$

Consider also the graph  $G$  consisting of the triples

$$(e, a, b), (e, c, d), (f, f, f).$$

As in Example 4.6, the normative semantics of the query on  $G$  is ill-defined as it requires evaluating  $Extend(e, ?y, (?y, ?y, ?y))$ , which is not a syntactically valid pattern.

A proposed fix to the normative semantics, discussed in the SPARQL 1.1 errata document [31] and associated forums, is to refrain from replacing the first argument of  $Extend$  when substituting an environment. We argue, however, that this fix leads to unintuitive behaviour. Indeed, in this example, the exists expression should intuitively return *false*, because  $\{?x \mapsto e\}$  is the only mapping satisfying the triple patterns  $(?x, a, b)$  and  $(?x, c, d)$ , whereas  $\{?x \mapsto f, ?y \mapsto f\}$  is the only mapping satisfying  $Extend(?x, ?y, (?y, ?y, ?y))$ , that is, there are no mappings in the evaluation of the first and second arguments of  $Join$  that have the same value of  $?x$ . However, the proposed fix would apply the mapping  $\{?x \mapsto e\}$  to the pattern inside the exists expression to yield the pattern  $Join((e, c, d), Extend(?x, ?y, (?y, ?y, ?y)))$ , which evaluates to the mapping  $\{?x \mapsto f, ?y \mapsto f\}$ ; as a result, the exists expression would evaluate to *true*.

In contrast to this fix, the environment semantics evaluates the exists expression in the query to *false*, as expected. Indeed, the environment  $\nu$ , in which exists is evaluated, is  $\{?x \mapsto e\}$ , so  $\llbracket Extend(?x, ?y, (?y, ?y, ?y)) \rrbracket_G^\nu = \emptyset$ , because  $\llbracket (?y, ?y, ?y) \rrbracket_G^\nu$  consists of a single mapping  $\{?x \mapsto e, ?y \mapsto f\}$ , which is incompatible with the extension  $\{?x \mapsto f\}$ , and therefore

$$\llbracket Join((?x, c, d), Extend(?x, ?y, (?y, ?y, ?y))) \rrbracket_G^\nu = \emptyset.$$

To summarise, the environment semantics of *exists* is always well-defined and provides intuitive answers.

We conclude this section with the observation that Theorem 4.13 generalises to the languages with *Extend*, that is, if the normative semantics is well-defined then it gives the same answers as the environment semantics.

**PROPOSITION 5.5.** *For any  $\text{Sparql}_{PDE}^{\exists}$  pattern  $P$  and graph  $G$  such that  $\llbracket P \rrbracket_G$  is well-defined it holds that  $\llbracket P \rrbracket_G = [P]_G^0$ .*

The proof, given in the appendix, uses an extension of Lemma 4.12 to  $\text{Sparql}_{PDE}^{\exists}$ .

## 5.2 Expressive Power of Languages with Assignment

First we show that *exists* does not add expressive power if we enrich the languages considered so far with *Extend*. For this, we prove the following analogue of Theorem 4.20.

**THEOREM 5.6.** *Let  $\text{Sparql}_{\star}$  be one of  $\text{Sparql}_E$ ,  $\text{Sparql}_{PE}$ ,  $\text{Sparql}_{DE}$  and  $\text{Sparql}_{PDE}$ . Language  $\text{Sparql}_{\star}^{\exists}$  has the same expressive power as  $\text{Sparql}_{\star}$ .*

**PROOF.** Lemmas 4.16 and 4.17 hold for language  $\text{Sparql}_{\star}$ , so it is enough to extend Lemma 4.18 to  $\text{Sparql}_{\star}^{\exists}$ . To this end, let

$$\text{Sub}_{\theta}(\text{Extend}(\?x, E, P'), P_1) = \text{Extend}(\?x, E, \text{Sub}_{\theta}(P', P_1)).$$

Claim 4.19 then generalises to  $\text{Sparql}_{\star}^{\exists}$  (with *Extend*) in a straightforward way: we just need to consider another case for the inductive step. The proof of this fact is in the appendix, at the end of the proof of Claim 4.19.  $\square$

Next we note that the proof of Theorem 3.10 goes through in the presence of *Extend* as well (in particular, Claim 3.11 generalises), meaning that every  $\text{Sparql}_{\star}^{\exists}$  query, for  $\text{Sparql}_{\star}$  one of  $\text{Sparql}_{PE}$ ,  $\text{Sparql}_{DE}$ , or  $\text{Sparql}_{PDE}$  (and their extensions with *exists*), has an equivalent  $\text{Sparql}_{\star}$  query in s-normal form. Therefore, we obtain the following analogue of Corollary 3.12.

**COROLLARY 5.7.** *Languages  $\text{Sparql}_{PE}$  and  $\text{Sparql}_{DE}$  have the same expressive power as  $\text{Sparql}_E$ .*

Finally, the proof of Theorem 3.14 is unaffected by the presence of *Extend*, so we also have the following inexpressibility result.

**PROPOSITION 5.8.** *Language  $\text{Sparql}_{PDE}$  is strictly more expressive than  $\text{Sparql}_E$ .*

Unsurprisingly, adding the *Extend* operator to any language considered so far increases its expressive power since *Extend* provides queries with a means of returning values that do not occur in the queried graph.

**PROPOSITION 5.9.** *Let  $\text{Sparql}_{\star}$  be either  $\text{Sparql}$  or  $\text{Sparql}_{PD}$ . Language  $\text{Sparql}_{\star E}$  is strictly more expressive than  $\text{Sparql}_{\star}$ .*

**PROOF.** Let the query  $Q$  be defined as follows:

$$\text{Project}(\{\?x\}, \text{Extend}(\?x, \text{bound}(\?y), (\?y, a, a))).$$

Over the graph  $G = \{(a, a, a)\}$ ,  $[Q]_G^0$  contains the mapping  $\{\?x \mapsto \text{true}\}$ . However, any  $\text{Sparql}_{PD}$  query  $Q'$  satisfies  $\mu(\?z) = a$  for each  $\mu \in [Q']_G^0$  and  $\?z \in \text{dom}(\mu)$ ; as a result,  $Q'$  cannot be equivalent to  $Q$ .  $\square$

We conclude this section with the following observation. Language  $\text{Sparql}_{PDE}^{\exists}$ , as any other language considered in this paper, possesses the property stated in Proposition 4.10:  $\nu \subseteq \mu$  for every mapping  $\mu \in [P]_G^V$ . This property is natural, but it might not hold for some other potential extensions of Sparql. For such extensions it seems to be reasonable to add the condition  $\{?x \mapsto [E]_G^{\mu'}\} \sim \nu$  to the first argument of the bag union in the definition of the semantics of *Extend* (while for languages satisfying the above property this condition holds automatically).

## 6 AGGREGATION

We now turn our attention to aggregation in SPARQL 1.1. In Section 6.1 we discuss the normative algebra for aggregation and present a formalisation that makes ambiguous aspects of the specification precise. In Section 6.2 we demonstrate the power of the aggregate algebra by showing that it is capable of expressing variable assignment as well as nested queries in their full generality. In Section 6.3 we present a normal form which leads to a substantial simplification of the SPARQL 1.1 aggregate algebra where most of its unconventional aspects are eliminated.

### 6.1 Normative Aggregate Algebra

SPARQL 1.1 and SQL provide similar functionality for aggregation: grouping is first used to define equivalence classes of solution mappings over which aggregate functions are subsequently applied.

*Example 6.1.* Consider the following query.

(Q5) *Return the salary average in CS as well as the average over all other departments combined, where the averages are only relevant if the minimal salary exceeds 3500.*

It can be written in SPARQL 1.1 as follows:

```
SELECT ?d (AVG(?s) AS ?n)
  WHERE { ?x department ?d . ?x salary ?s }
GROUP BY (?d = CS AS ?d)
HAVING (MIN(?s) > 3500).
```

Note that the grouping expression  $?d = \text{CS}$  is Boolean. The value of this expression is then again assigned to variable  $?d$ . These two occurrences of  $?d$  are semantically unrelated: the one in  $?d = \text{CS}$  is local to the aggregation and bound in the WHERE clause, while the one being assigned to is global and referred to in SELECT. Over  $G_{ex}$ , (Q5) evaluates to  $\{?d \mapsto \text{true}, ?n \mapsto 4500\}$ , since the salary of *Charlie* is less than 3500, while the average over the salaries of *Alice* and *Bob* is 4500.

The SPARQL 1.1 aggregate algebra has several unconventional features when compared with SQL:

- (F1) groups and aggregates are seen as first-class citizens of the algebra, which are defined independently using dedicated constructs *Group* and *Aggregate*;
- (F2) grouping is allowed on arbitrary lists of expressions, and not just on lists of variables (as in the GROUP BY expression in query (Q5)); and
- (F3) aggregation is also allowed on arbitrary lists of expressions, and not just on single expressions.

In what follows, we provide a rigorous formalisation of the normative aggregate algebra that addresses unspecified corner cases in the standard and makes ambiguous aspects of the specification precise. Our semantics also generalises the normative one in order to handle exists expressions correctly, and hence relies on the notion of environment introduced in Section 4.2.

We start our discussion by introducing groups as first-class citizens. A group induces a partitioning of a pattern's solution mappings into equivalence classes, each of which is determined by a key obtained from the evaluation of a list of expressions.

*Definition 6.2.* A *v-list* is a list of values in  $\mathbf{T} \cup \{\text{error}\}$ . The *evaluation*  $[E]_G^\mu$  of an expression list  $E = \langle E_1, \dots, E_k \rangle$  over a graph  $G$  with respect to a mapping  $\mu$  is the v-list  $\langle [E_1]_G^\mu, \dots, [E_k]_G^\mu \rangle$ .

*Definition 6.3.* A *group* is a construct  $\text{Group}(\mathbf{F}, P)$  with  $\mathbf{F}$  a list of expressions and  $P$  a pattern. The *evaluation*  $[\Gamma]_G$  of a group  $\Gamma = \text{Group}(\mathbf{F}, P)$  over a graph  $G$  is a partial function from v-lists to multisets of mappings that is defined for all v-lists  $\text{Key} = [\mathbf{F}]_G^\mu$  with  $\mu \in [P]_G^0$  as follows:

$$[\Gamma]_G(\text{Key}) = \{\mu' \mid \mu' \in [P]_G^0, [\mathbf{F}]_G^{\mu'} = \text{Key}\}.$$

Note that the evaluation of groups does not depend on environments. As we will see later on, groups can only appear in patterns that use the construct *AggregateJoin*, where the variables in the group are always local, that is, not visible outside the pattern. As a consequence, the environment in which the pattern is evaluated does not affect the semantics of the group.

As in SQL, *aggregate functions* in SPARQL 1.1 (e.g., AVG in query (Q5)) allow us to compute a single value for each group of solution mappings. In the relational case, they are functions from multisets of values to a single value [17]. Due to (F3), aggregate functions in SPARQL 1.1 deal with more complex structures involving multisets of v-lists. To handle them, SPARQL 1.1 introduces the function defined next.

*Definition 6.4.* Function Flatten maps each multiset  $\Lambda$  of v-lists to the multiset  $\Theta$  of values in  $\mathbf{T} \cup \{\text{error}\}$  whose base consists of values in  $\Lambda$  and whose multiplicity function is as follows, for each such value  $v$ :

$$\text{card}_\Theta(v) = \sum_{\lambda \in \Lambda} (\text{card}_\lambda(v) \times n_{v,\lambda}),$$

where  $n_{v,\lambda}$  is the number of appearances of  $v$  in  $\lambda$ .

SPARQL 1.1 provides aggregate functions analogous to those in SQL. Differences stem mostly from the treatment of lists and errors.

*Definition 6.5.* Let  $<$  be a total order on values that extends the usual orders on literals and such that  $\text{error} < b < u < \ell$  for any  $b \in \mathbf{B}$ ,  $u \in \mathbf{I}$ , and  $\ell \in \mathbf{L}$ .

A SPARQL 1.1 *aggregate function* is one of the following functions, mapping multisets of v-lists  $\Lambda$  to values in  $\mathbf{T} \cup \{\text{error}\}$ , where  $\Theta = \text{Flatten}(\Lambda)$ :

- $\text{Count}(\Lambda) = \sum_{v \in \Theta, v \neq \text{error}} \text{card}_\Theta(v)$ ;
- $\text{Sum}(\Lambda)$  is  $\sum_{v \in \Theta} (\text{card}_\Theta(v) \times v)$  if all the values in  $\Lambda$  are numbers, and *error* otherwise;
- $\text{Avg}(\Lambda)$  is 0 if  $\text{Count}(\Lambda) = 0$  and  $\text{Sum}(\Lambda) / \text{Count}(\Lambda)$  otherwise (in particular, it is *error* if  $\text{Sum}(\Lambda) = \text{error}$ );
- $\text{Min}(\Lambda)$  is the  $<$ -minimal value in  $\Theta$  if  $\Theta \neq \emptyset$  and *error* otherwise;
- $\text{Max}(\Lambda)$  is the  $<$ -maximal value in  $\Theta$  if  $\Theta \neq \emptyset$  and *error* otherwise.

Finally, CountD, SumD, and AvgD are defined as their counterparts Count, Sum, and Avg, but applied to the multiset of v-lists obtained from  $\Lambda$  by removing duplicates.<sup>3</sup>

Note that *error* does not contribute to Count, but may affect the results of other functions. We use Id as a synonym for Min whenever, by construction,  $\text{Flatten}(\Lambda)$  consists of a single value (with any multiplicity).

We define the aggregate construct, which computes a value for each group by means of aggregate functions.

<sup>3</sup>According to this definition  $\Theta$  may still contain repeated values, because duplicate elimination is applied before Flatten; this may seem counter-intuitive, but we follow the SPARQL 1.1 specification here.

*Definition 6.6.* An *aggregate* is a construct of the form  $Aggregate(E, f, \Gamma)$ , for  $E$  a list of expressions,  $f$  an aggregate function, and  $\Gamma = Group(F, P)$  a group.

The *evaluation*  $[A]_G$  of an aggregate  $A = Aggregate(E, f, \Gamma)$  over a graph  $G$  is the partial function from  $v$ -lists to values such that, for each  $Key$  in the domain of  $[A]_G$ ,

$$[A]_G(Key) = f(\llbracket \Lambda \mid \mu \in [\Gamma]_G(Key), \Lambda = [E]_G^\mu \rrbracket).$$

Note that, for brevity, we omit a special case in SPARQL 1.1, which additionally allows aggregates to use `Count` and `CountD` with a distinguished symbol `*` in place of the list  $E$ . However, such special aggregates can be expressed in terms of the ones in Definition 6.6, as will be discussed in Section 9.

Finally, the normative SPARQL 1.1 algebra provides the *AggregateJoin* construct, which combines a list of aggregates  $A_1, \dots, A_n$  to form a pattern  $P$ . The semantics mandates that  $[P]_G^v$  contain a mapping for each  $v$ -list in the domains of all  $[A_i]_G$ ; each such mapping assigns variables  $?x_i$  to record the values of aggregates  $A_i$  for the corresponding  $v$ -list.

*Definition 6.7.* Let  $Sparql_\star$  be an extension  $Sparql$ . The language  $Sparql_{\star A}$  extends  $Sparql_\star$  by allowing patterns of the form  $AggregateJoin_x(A)$ , with  $\mathbf{x} = \langle ?x_1, \dots, ?x_n \rangle$  a non-empty list of variables and  $\mathbf{A} = \langle A_1, \dots, A_n \rangle$  a list of aggregates of the same size.

The *evaluation* of a pattern  $AggregateJoin_x(A)$  over a graph  $G$  in an environment  $v$  is as follows, where  $\Lambda$  is the intersection of the domains of all  $A_i$ :

$$[AggregateJoin_x(A)]_G^v = \bigcup_{Key \in \Lambda} \llbracket \mu \mid \mu = \mu' \cup v, \mu' = \{?x_i \mapsto v \mid 1 \leq i \leq n, v = [A_i]_G(Key), v \neq error\} \rrbracket.$$

We also set  $var(AggregateJoin_x(A)) = \{?x_1, \dots, ?x_n\}$ .

*Example 6.8.* Query (Q5) in Example 6.1 translates into the normative algebra as given next, where we have a single group over departments, and aggregates  $A_2$  and  $A_3$  for computing the average and minimum; an additional aggregate  $A_1$  is required to store the keys of the groups and incorporate them into a pattern using *AggregateJoin*:

$$Project(\{?d, ?n\}, Extend(?n, ?v_2, Extend(?d, ?v_1, P_1))),$$

where pattern  $P_1$  is defined as given next:

$$\begin{aligned} P_1 &= Filter(?v_3 > 3500, AggregateJoin_{\langle ?v_1, ?v_2, ?v_3 \rangle}(\langle A_1, A_2, A_3 \rangle)), \\ A_1 &= Aggregate(\langle ?d = CS \rangle, Id, Group(\langle ?d = CS \rangle, P_2)), \\ A_2 &= Aggregate(\langle ?s \rangle, Avg, Group(\langle ?d = CS \rangle, P_2)), \\ A_3 &= Aggregate(\langle ?s \rangle, Min, Group(\langle ?d = CS \rangle, P_2)), \\ P_2 &= \{(?x, department, ?d), (?x, salary, ?s)\}. \end{aligned}$$

Note how variable  $?d$  is handled: first an expression with the inner occurrence of  $?d$  is aggregated to an auxiliary variable  $?v_1$ , and then the outer occurrence is assigned by *Extend* to the (trivial) expression  $?v_1$ . The two occurrences of  $?d$  are, essentially, different variables, and the inner one is local to the *AggregateJoin* (as are all other variables of *AggregateJoin*). This justifies the evaluation of groups in the empty environment.

In conclusion, we note that, similarly to the case of *Extend*, Theorem 4.13 generalises to the languages with *AggregateJoin*, that is, if the normative semantics is well-defined then it gives the same answers as the environment semantics. To make this statement formal, we need to define the normative semantics for aggregates, which requires a generalisation of the notion of rectification: a pattern (with aggregates) is *rectified* if, same as in the basic case, for each subpattern

$Project(X, P')$  no variable in  $\text{var}(P') \setminus X$  appears outside this subpattern and, additionally, for each subpattern  $AggregateJoin_x(\mathbf{A})$  no variable in  $\mathbf{A}$  appears outside this subpattern. Each pattern can be converted to a rectified pattern in the same way as in the case without aggregation. This treatment of aggregates is justified by the fact that all variables in aggregates are essentially local, and their precise names should not influence the semantics. Then,  $\llbracket AggregateJoin_x(\mathbf{A}) \rrbracket_G$  is defined in the same way as  $\llbracket AggregateJoin_x(\mathbf{A}) \rrbracket_G^v$  except that (i)  $\mathbf{A}$  is rectified, (ii) the condition  $\mu = \mu' \cup v$  in Definition 6.7 is replaced by  $\mu = \mu'$ , and (iii) subexpressions and subpatterns are evaluated in the obvious way according to the normative semantics rather than the environment semantics.

**PROPOSITION 6.9.** *For any Sparql<sub>PDEA</sub><sup>∃</sup> pattern  $P$  and graph  $G$  such that  $\llbracket P \rrbracket_G$  is well-defined it holds that  $\llbracket P \rrbracket_G = \llbracket P \rrbracket_G^0$ .*

The proof, given in the appendix, is based on an extension of Lemma 4.12 to Sparql<sub>PDEA</sub><sup>∃</sup>.

## 6.2 Expressive Power of Languages with Aggregates

Operators *Group*, *Aggregate* and *AggregateJoin* provide a great deal of power and flexibility to the query language. We next show that, when added to Sparql, these operators are sufficiently expressive to capture all forms of query nesting and variable assignment discussed so far.

*Example 6.10.* To see how *Extend* can be expressed in terms of aggregates, consider the pattern  $P_0 = Extend(?x, ?y_1 + ?y_2, (?y_1, p, ?y_2))$ . This pattern can be equivalently expressed by a Sparql<sub>A</sub> pattern  $P_1$  defined as follows:

$$\begin{aligned} A_{?y_i} &= Aggregate(\langle ?y_i \rangle, Id, Group(\langle ?y_1, ?y_2 \rangle, (?y_1, p, ?y_2))), \quad i = 1, 2, \\ A_{?x} &= Aggregate(\langle ?y_1 + ?y_2 \rangle, Id, Group(\langle ?y_1, ?y_2 \rangle, (?y_1, p, ?y_2))), \\ P_1 &= AggregateJoin_{\langle ?x, ?y_1, ?y_2 \rangle}(A_{?x}, A_{?y_1}, A_{?y_2}). \end{aligned}$$

To eliminate *Extend* over arbitrary patterns rather than simple triple patterns, this construction needs to be generalised to account for partial mappings as well as mappings with multiplicities greater than 1 (in this example it is not needed, because  $P_0$  is equivalent to  $Distinct(P_0)$  and all mappings in its evaluation have the same domain).

**THEOREM 6.11.** *Language Sparql<sub>PDEA</sub> has the same expressive power as Sparql<sub>A</sub>.*

**PROOF.** We start by expressing patterns  $Extend(?x, E, P)$  in Sparql<sub>PA</sub>. For an enumeration  $\mathbf{x} = \langle ?x_1, \dots, ?x_n \rangle$  of  $\text{var}(P)$ , let

$$\begin{aligned} P' &= AggregateJoin_{\langle ?x, ?x_1, \dots, ?x_n \rangle}(\langle A, A_1, \dots, A_n \rangle), \text{ where} \\ A_i &= Aggregate(\langle ?x_i \rangle, Id, Group(\mathbf{x}, P)), \quad 1 \leq i \leq n, \\ A &= Aggregate(\langle E \rangle, Id, Group(\mathbf{x}, P)). \end{aligned}$$

By construction, the evaluation of  $P'$  consists of the same mappings as the evaluation of the pattern  $Extend(?x, E, P)$ , but all with multiplicities being 1: in particular, for any distinct mapping  $\mu$  in  $\llbracket P \rrbracket_G^0$  there is a  $v$ -list in the domain of the evaluation of  $Group(\mathbf{x}, P)$ , namely the list of values of  $\mathbf{x}$  under  $\mu$ , where the unbound variables are represented as *error*; each  $A_i$  then evaluates to the  $i$ -th value in the list, while *AggregateJoin* combines these values back to the mapping extended with the evaluation of  $E$  in  $?x$  in such a way that variables corresponding to *error* values are left unbound.

Then the following pattern is fully equivalent to  $Extend(?x, E, P)$ , where  $\theta$  is a renaming of  $\text{var}(P)$ :

$$Project(\{?x\} \cup \text{var}(P), Filter(eq(\text{var}(P), \theta), Join(P, P'\theta))).$$

Patterns of the form  $Distinct(P)$  are expressed similarly: we can construct  $P'$  as above, but without aggregate  $A$  and variable  $?x$ , and take  $P'$  as an equivalent pattern.

Finally,  $Project$  can be pushed upwards through Sparql operators as in Theorem 3.10. Thus, it suffices to show that  $Project$  can be eliminated from any group  $\Gamma$  of the form  $Group(F, Project(X, P))$ . For this, replace  $Project(X, P)$  in  $\Gamma$  by  $P\theta'$ , with  $\theta'$  a renaming of  $\text{var}(P) \setminus X$ .  $\square$

Since all variables in groups are considered local, that is, groups are always evaluated in empty environments, we have the following corollary of Theorem 5.6, saying that exists does not add any expressive power in the presence of aggregation.

**COROLLARY 6.12.** *Language  $\text{Sparql}_{PDEA}^{\exists}$  has the same expressive power as  $\text{Sparql}_{PDEA}$ .*

**PROOF.** Lemmas 4.16 and 4.17 hold for language  $\text{Sparql}_{PDEA}$ , so it is enough to extend Lemma 4.18 to  $\text{Sparql}_{PDEA}$ . To this end, let

$$\text{Sub}_{\theta}(\text{AggregateJoin}_x(\mathbf{A}), P_1) = \text{Join}(\text{AggregateJoin}_x(\mathbf{A}), \text{Copy}_{\theta}(P_1)).$$

Claim 4.19 straightforwardly generalises to  $\text{Sparql}_{PDEA}$  by adding another base case.  $\square$

We conclude this section by showing that aggregation does provide additional expressive power to the language.

**PROPOSITION 6.13.** *Language  $\text{Sparql}_A$  is strictly more expressive than  $\text{Sparql}_{PDE}$ .*

**PROOF.** Let  $\mathcal{G}$  be the set of all graphs that contain no numbers.

On the one hand, for every  $\text{Sparql}_{PDE}$  query  $Q$  and graph  $G$  in  $\mathcal{G}$ , the maximal number occurring in any mapping  $\mu$  in  $[Q]_G^0$  is bounded by a function in the size  $|Q|$  of  $Q$  and the maximal number  $n$  occurring in  $Q$ . Indeed, whenever  $\mu(?x)$  is a number,  $Q$  contains a subpattern of the form  $Extend(?x, E, P)$ . Therefore, by a straightforward induction on the nesting depth  $d$  of  $Extend$  in  $Q$ , one can establish that all numbers in  $\mu$  are bounded by  $n^{|Q|^d}$ . In particular, this bound does not depend on  $G$ .

On the other hand, let  $Q$  be the  $\text{Sparql}_A$  query

$$Project(\{?c\}, \text{Extend}(?c, ?v, \text{AggregateJoin}_{(?v)}(\langle A \rangle))),$$

where  $A = \text{Aggregate}(\langle ?x \rangle, \text{Count}, \text{Group}(\langle \rangle, (?x, ?y, ?z)))$ ; this query is obtained from the SPARQL 1.1 query

```
SELECT (COUNT(?x) AS ?c) WHERE { ?x ?y ?z }.
```

Then, for every  $G \in \mathcal{G}$ , the multiset  $[Q]_G^0$  consists of a single mapping  $\mu$  such that  $\mu(?c)$  is the number of triples in  $G$ . Since  $\mu(?c)$  cannot be bounded independently of  $G$ , it follows that  $Q$  has no equivalent query in  $\text{Sparql}_{PDE}$ .  $\square$

### 6.3 Normalisation and Simplification

We next show that features (F2) and (F3) in the normative algebra do not add expressive power: every query can be rewritten into a normal form where grouping is only allowed over lists of variables rather than arbitrary expressions, and aggregation is done only over singleton lists. Moreover, our normal form dispenses with the functions  $\text{CountD}$ ,  $\text{SumD}$  and  $\text{AvgD}$ , and hence shows that it suffices to consider aggregate functions that do not involve duplicate elimination.

**Definition 6.14.** A  $\text{Sparql}_A$  pattern is in *a-normal form* if each group is of the form  $Group(\mathbf{y}, P)$  with  $\mathbf{y}$  a list of variables and each aggregate is of the form  $Aggregate(\langle E \rangle, f, \Gamma)$  with  $f$  different from  $\text{CountD}$ ,  $\text{SumD}$ , and  $\text{AvgD}$ .

*Example 6.15.* The algebraic translation of query (Q5), given in Example 6.8, is not in a-normal form since grouping is not performed over variables but a complex expression  $?d = CS$ . However, (Q5) can be easily brought to a-normal form by replacing every occurrence of  $Group(\langle ?d = CS \rangle, P_2)$  by the group

$$Group(\langle ?y \rangle, Extend(?y, ?d = CS, P_2)),$$

where  $?y$  is a fresh variable, and then eliminating  $Extend$  as done in the proof of Theorem 6.11. Note that elimination of  $Extend$  from otherwise a-normal formulas yields a-normal formulas.

Similarly, the following aggregate  $A_0$  is not a-normal since it uses  $AvgD$  and aggregates over a list with more than one component:

$$Aggregate(\langle ?x - ?z, ?z \rangle, AvgD, Group(\langle ?y \rangle, (?x, ?y, ?z))).$$

To bring  $A_0$  to a-normal form, we first rewrite it to the following equivalent aggregate  $A_1$ , which uses  $Avg$  instead of  $AvgD$  and fresh variables  $?u, ?v, ?w$ :

$$\begin{aligned} A_{?u} &= Aggregate(\langle ?x - ?z \rangle, Id, Group(\langle ?x - ?z, ?z, ?y \rangle, (?x, ?y, ?z))), \\ A_{?v} &= Aggregate(\langle ?z \rangle, Id, Group(\langle ?x - ?z, ?z, ?y \rangle, (?x, ?y, ?z))), \\ A_{?w} &= Aggregate(\langle ?y \rangle, Id, Group(\langle ?x - ?z, ?z, ?y \rangle, (?x, ?y, ?z))), \\ P'_1 &= AggregateJoin_{\langle ?u, ?v, ?w \rangle}(\langle A_{?u}, A_{?v}, A_{?w} \rangle), \\ A_1 &= Aggregate(\langle ?u, ?v \rangle, Avg, Group(\langle ?w \rangle, P'_1)). \end{aligned}$$

Note that pattern  $P'_1$  “simulates”  $(?x, ?y, ?z)$  in the sense that, for any graph  $G$  and v-list  $\lambda$ , we have  $\lambda \in [\langle ?x - ?z, ?z, ?y \rangle]_G^\mu$  for some  $\mu \in [(\langle ?x, ?y, ?z \rangle)]_G^0$  if and only if  $\lambda \in [\langle ?u, ?v, ?w \rangle]_G^{\mu'}$  for some  $\mu' \in [P'_1]_G^0$ . However, by construction,  $[P'_1]_G^0$  contains each mapping at most once, and each pair of distinct mappings in  $[P'_1]_G^0$  yields two distinct v-lists, which justifies using  $Avg$  in  $A_1$  in place of  $AvgD$  in  $A_0$ .

Note also that removing  $AvgD$  introduced grouping over a complex expression, which can in turn be removed as discussed above, resulting in an aggregate  $A_2$  of the form  $Aggregate(\langle ?u, ?v \rangle, Avg, \Gamma)$ , for  $\Gamma$  an a-normal  $Sparql_A$  rewriting of  $Group(\langle ?w \rangle, P'_1)$ .

Thus, it remains to remove the aggregation over  $\langle ?u, ?v \rangle$ , for which it suffices to note that  $A_2$  is equivalent to  $Aggregate(\langle (?u + ?v)/2 \rangle, Avg, \Gamma)$ , which is in a-normal form.

Next, we show that a-normalisation is always feasible.

**THEOREM 6.16.** *Every  $Sparql_A$  pattern has an  $\emptyset$ -equivalent  $Sparql_A$  pattern in a-normal form.*

**PROOF.** We first show that aggregate functions with duplicate elimination can be rewritten using their ordinary counterparts. Let  $fD \in \{CountD, SumD, AvgD\}$  and

$$A_1 = Aggregate(E, fD, Group(F, P_1))$$

with  $E = \langle E_1, \dots, E_m \rangle$  and  $F = \langle F_1, \dots, F_k \rangle$ . Then  $A_1$  is equivalent to the aggregate  $A'_1$  defined as follows, where  $\mathbf{x} = \langle ?x_1, \dots, ?x_m \rangle$  and  $\mathbf{y} = \langle ?y_1, \dots, ?y_k \rangle$  are lists of fresh variables, and  $\cdot$  denotes list concatenation:

$$\begin{aligned} A_{?x_i} &= Aggregate(\langle E_i \rangle, Id, Group(E \cdot F, P_1)), \quad 1 \leq i \leq m, \\ A_{?y_j} &= Aggregate(\langle F_j \rangle, Id, Group(E \cdot F, P_1)), \quad 1 \leq j \leq k, \\ P'_1 &= AggregateJoin_{\mathbf{x}, \mathbf{y}}(\langle A_{?x_1}, \dots, A_{?x_m}, A_{?y_1}, \dots, A_{?y_k} \rangle), \\ A'_1 &= Aggregate(\mathbf{x}, f, Group(\mathbf{y}, P'_1)). \end{aligned}$$

Second, we argue that grouping over lists of expressions can be reduced to grouping over lists of variables by exploiting *Extend*. For this, we note that an aggregate

$$A_2 = \text{Aggregate}(\mathbf{E}, f, \text{Group}(\mathbf{F}, P_2))$$

with  $\mathbf{F} = \langle F_1, \dots, F_k \rangle$  is equivalent to an aggregate  $A'_2$  defined as follows, where  $\mathbf{y} = \langle ?y_1, \dots, ?y_k \rangle$  is a list of fresh variables:

$$\begin{aligned} P'_2 &= \text{Extend}(?y_1, F_1, \dots \text{Extend}(?y_k, F_k, P_2) \dots), \\ A'_2 &= \text{Aggregate}(\mathbf{E}, f, \text{Group}(\mathbf{y}, P'_2)). \end{aligned}$$

By Theorem 6.11, *Extend* in  $P'_2$  is inessential as it is expressible using normalised grouping and aggregation constructs (in particular, no grouping over non-trivial expressions is introduced by the transformation in the proof of Theorem 6.11).

To summarise, every aggregate  $A$  can be rewritten to the form  $\text{Aggregate}(\mathbf{E}, f, \Gamma)$ , where  $\mathbf{E} = \langle E_1, \dots, E_m \rangle$ ,  $\Gamma = \text{Group}(\mathbf{y}, P)$ , and  $f$  is not CountD, SumD, or AvgD.

In the last step of the proof, we argue that every such  $A$  can be rewritten into an equivalent aggregate  $A'$  in a-normal form by reducing the list  $\mathbf{E}$  to a single expression; for example, computing Avg over a list  $\langle E_1, \dots, E_m \rangle$  is equivalent to aggregating over the expression  $(\sum_{i=1}^m E_i)/m$ . Unlike the previous two steps, this step is sensitive to the particular aggregate functions available in SPARQL 1.1, and hence we perform a case-by-case analysis.

Let  $f = \text{Count}$ . Let  $?x_1, \dots, ?x_m$  and  $?z_1, \dots, ?z_m$  be sets of fresh variables. We define

$$A' = \text{Aggregate}\left(\left\langle \sum_{i=1}^m ?x_i \right\rangle, \text{Sum}, \text{Group}(\mathbf{y}, P_m)\right),$$

where the patterns  $P_i$  are as follows:

$$\begin{aligned} P_0 &= \text{Extend}(?z_m, E_m, \dots \text{Extend}(?z_1, E_1, P) \dots), \\ P_{i+1} &= \text{Union}(\text{Filter}(\text{bound}(?z_{i+1}), \text{Extend}(?x_{i+1}, 1, P_i)), \\ &\quad \text{Filter}(\neg \text{bound}(?z_{i+1}), \text{Extend}(?x_{i+1}, 0, P_i))). \end{aligned}$$

Let  $f = \text{Sum}$ . We define  $A' = \text{Aggregate}\left(\left\langle \sum_{i=1}^m E_i \right\rangle, \text{Sum}, \Gamma\right)$ .

Let  $f = \text{Avg}$ . We define  $A' = \text{Aggregate}\left(\left\langle \left( \sum_{i=1}^m E_i \right) / m \right\rangle, \text{Avg}, \Gamma\right)$ .

Let  $f = \text{Min}$ . We define  $A' = \text{Aggregate}(\langle ?x \rangle, \text{Min}, \text{Group}(\mathbf{y}, P'))$ , where  $?x$  is a fresh variable and

$$\begin{aligned} P' &= \text{Union}(\text{AggregateJoin}_{\langle ?x \rangle}(\langle A_1 \rangle), \dots \\ &\quad \text{Union}(\text{AggregateJoin}_{\langle ?x \rangle}(\langle A_{m-1} \rangle), \text{AggregateJoin}_{\langle ?x \rangle}(\langle A_m \rangle)) \dots), \end{aligned}$$

and  $A_i = \text{Aggregate}(\langle E_i \rangle, \text{Min}, \text{Group}(\mathbf{y}, P'))$  for  $i = 1, \dots, m$ .

The case  $f = \text{Max}$  is analogous to the case for Min.

In all cases,  $[A]_G = [A']_G$  for every  $G$  by construction and the definitions of the aggregate functions.  $\square$

The normal form in Definition 6.14 already provides a significant simplification of the algebra. Indeed, it shows that features (F2) and (F3) are inconsequential. Our normal form also suggests that the definition of aggregate functions can be made more transparent than it is the case in the standard: not only the functions involving duplicate elimination can be dispensed with, but also the function Flatten is inessential since aggregation can always be performed over a single expression rather than a list.

We next show that feature (F1) is also immaterial; that is, we can collapse the *Group*, *Aggregate* and *AggregateJoin* constructs into a single pattern operator without affecting the expressive power of the language. This further simplification not only brings the normative aggregate algebra closer to its relational counterpart, but can also be exploited to make the mapping from the SPARQL 1.1 syntax into the algebra much more direct and transparent.

The following definition specifies the aforementioned combined operator.

*Definition 6.17.* Let  $\text{Sparql}_\star$  be an extension of  $\text{Sparql}$ . The language  $\text{Sparql}_{\star A_S}$  extends  $\text{Sparql}_\star$  by allowing patterns of the form  $\text{GroupAgg}(Z, ?x, f, E, P)$ , where  $Z$  is a set of variables, called *grouping variables*,  $?x$  is another variable, called *aggregation variable*,  $f$  is an aggregate function,  $E$  is an expression, and  $P$  is a pattern.

Given a pattern  $\text{GroupAgg}(Z, ?x, f, E, P)$ , a graph  $G$ , and a mapping  $\mu \in [P]_G^0$ , let

$$v_\mu = f(\{v \mid \mu' \in [P]_G^0, \mu'|_Z = \mu|_Z, v = [E]_{G'}^{\mu'}\}).$$

Then, the *evaluation*  $[\text{GroupAgg}(Z, ?x, f, E, P)]_G^v$  in an environment  $v$  over  $G$  is the multiset with the base set

$$\begin{aligned} &\{\mu' \mid \mu \in [P]_G^0, v_\mu \neq \text{error}, \mu' = \mu|_Z \cup \{?x \mapsto v_\mu\} \cup v\} \cup \\ &\{\mu' \mid \mu \in [P]_G^0, v_\mu = \text{error}, \mu' = \mu|_Z \cup v\}, \end{aligned}$$

and multiplicity 1 for each mapping in the base set. We also set

$$\text{var}(\text{GroupAgg}(Z, ?x, f, E, P)) = Z \cup \{?x\}.$$

The construct *GroupAgg* is close to the grouping operator in the relational algebra:  $Z$  represents the set of grouping variables,  $?x$  is the fresh variable storing the aggregation result,  $f$  is the aggregate function, and  $E$  is the expression (often a variable) we are aggregating over. The treatment of environments in Definition 6.17 is consistent with the normative algebra operators. Specifically,  $[\text{GroupAgg}(Z, ?x, f, E, P)]_G^v$  is defined by means of the evaluation  $[P]_G^0$  in the empty environment.

*Example 6.18.* Query (Q5) in Example 6.8 can be written in  $\text{Sparql}_{A_S}^3$  as follows (exists is just used for succinctness and can be dispensed with as usual):

*Project*( $\{?d, ?n\}$ , *Filter*(exists( $P_2$ ),  $P_1$ )), where

$$P_1 = \text{GroupAgg}(\{?d\}, ?n, \text{Avg}, ?s, P_3),$$

$$P_2 = \text{Filter}(?v > 3500, \text{GroupAgg}(\{?d\}, ?v, \text{Min}, ?s, P_3)),$$

$$P_3 = \text{GroupAgg}(\{?x, ?c, ?s\}, ?d, \text{Id}, ?c = \text{CS}, \{(?x, \text{department}, ?c), (?x, \text{salary}, ?s)\}).$$

The following theorem shows that the *GroupAgg* construct captures grouping and aggregation in SPARQL 1.1. While the reduction from  $\text{Sparql}_{A_S}$  to  $\text{Sparql}_A$  is straightforward, the other direction generalises Example 6.18. There, to express *AggregateJoin*, one needs to join *GroupAgg* subpatterns in a way that forces the respective mappings to coincide on a given set of variables (rather than just being compatible on the set). For this, we once again employ techniques developed in Sections 3.1 and 4.4.

**THEOREM 6.19.** *For any extension  $\text{Sparql}_{\star P}$  of  $\text{Sparql}_P$ , language  $\text{Sparql}_{\star P A_S}$  has the same expressive power as  $\text{Sparql}_{\star P A}$ .*

**PROOF.** We first show that each  $\text{Sparql}_{\star P A_S}$  pattern

$$P = \text{GroupAgg}(\{?z_1, \dots, ?z_k\}, ?x, f, E, P')$$

has a pattern in  $\text{Sparql}_{\star P A}$  that is  $v$ -equivalent to  $P$  for any environment  $v$ . For this, let

$$\Gamma = \text{Group}(\{?z_1, \dots, ?z_k\}, P').$$

We then use aggregates  $A_i = \text{Aggregate}(\langle ?z_i \rangle, \text{Id}, \Gamma)$ ,  $1 \leq i \leq k$ , to record the values of the grouping variables in each group; finally, we capture the value of  $E$  using  $A = \text{Aggregate}(\langle E \rangle, f, \Gamma)$ . Then, for any  $\nu$ ,  $P$  is  $\nu$ -equivalent to

$$\text{AggregateJoin}_{\langle ?z_1, \dots, ?z_k, ?x \rangle}(\langle A_1, \dots, A_k, A \rangle).$$

For the other direction, consider a  $\text{Sparql}_{\star PA}$  query  $Q$  with all  $\text{AggregateJoin}$  subpatterns in the  $a$ -normal form. We will replace each such subpattern  $P$  with a  $\text{Sparql}_{\star PAs}$  subpattern  $P''$  and then prove that the resulting  $\text{Sparql}_{\star PAs}$  query is equivalent to  $Q$ . To this end, let

$$P = \text{AggregateJoin}_{\langle ?x_1, \dots, ?x_n \rangle}(\langle A_1, \dots, A_n \rangle), \text{ where} \\ A_i = \text{Aggregate}(\langle E_i \rangle, f_i, \text{Group}(y_i, P_i)), \quad 1 \leq i \leq n.$$

Assume, without loss of generality, that all  $y_i$  are of the same length  $k$ , since otherwise the evaluation of  $P$  is empty. Let  $Z = \{?z_1, \dots, ?z_k\}$  be fresh variables and, for  $1 \leq i \leq n$ , let  $\theta_i$  be a renaming of variables  $y_i$  to the corresponding variables in  $Z$ . We simulate each  $A_i$  by the pattern

$$P'_i = \text{GroupAgg}(Z, ?x_i, f_i, E_i \theta_i, P_i \theta_i).$$

We then combine these patterns as follows, where each  $\theta'_i$ ,  $1 \leq i < n$ , is a renaming of  $Z$  to fresh  $Z_i$ :

$$P'' = \text{Project}(\langle ?x_1, \dots, ?x_n \rangle, P'), \text{ where} \\ P' = \text{Filter}(\text{eq}(Z, \theta'_1), \dots, \text{Filter}(\text{eq}(Z, \theta'_{n-1}), \text{Join}(P'_1 \theta'_1, \dots, \text{Join}(P'_{n-1} \theta'_{n-1}, P'_n) \dots)) \dots).$$

We argue that  $P''$  is  $\nu$ -equivalent to  $P$  for every environment  $\nu$  in which  $P$  can be evaluated as a subpattern of  $Q$ , that is,  $P$  can be safely replaced by  $P''$  in  $Q$  without changing the semantics of the overall query. Note that, for any graph  $G$ , each  $\mu \in [P]_G^\nu$  is of the form  $\{?x_1 \mapsto u_1, \dots, ?x_n \mapsto u_n\}$ , and it corresponds to a single key  $v^\mu = \langle v_1^\mu, \dots, v_k^\mu \rangle$  such that  $(v^\mu \mapsto u_i) \in [A_i]_G$  for all  $1 \leq i \leq n$ . Given such  $\mu$ , let  $\mu'$  be the mapping with  $\text{dom}(\mu') \subseteq Z$  such that  $\mu'(?z_j) = v_j^\mu$  whenever  $v_j^\mu \neq \text{error}$  and  $?z_j \notin \text{dom}(\mu')$  otherwise. Since variables  $Z_i$  and  $Z$  do not occur in  $Q$ , they cannot be bound in any environment  $\nu$  in which  $P$  is evaluated. Therefore, by construction,  $\mu \in [P]_G^\nu$  if and only if the mapping

$$\mu \cup \mu' \theta'_1 \cup \dots \cup \mu' \theta'_{n-1} \cup \mu'$$

is in  $[P']_G^\nu$  for any such  $\nu$ . Since

$$\{?x_1, \dots, ?x_n\} \cap (Z_1 \cup \dots \cup Z_{n-1} \cup Z) = \emptyset,$$

it follows that  $[P'']_G^\nu$  and  $[P]_G^\nu$  have the same mappings. Finally, the multiplicity of every mapping in each of these multisets is 1, by definition of  $[P]_G^\nu$  and by construction of  $[P'']_G^\nu$ , so  $P$  can be replaced by  $P''$  in  $Q$ .  $\square$

Note that in the theorem, the reduction from  $\text{Sparql}_{\star PAs}$  to  $\text{Sparql}_{\star PA}$  does not use *Project*. The reduction for the other direction introduces *Project*, which may occur at an intermediate level of the overall query. However, if  $\text{Sparql}_{\star PAs}$  allows to “push” projection up through all other pattern constructs, it follows that languages  $\text{Sparql}_{\star As}$  and  $\text{Sparql}_{\star A}$  are equivalent as well, because the subpattern  $\text{Project}(X, P)$  in  $\text{GroupAgg}(Z, ?x, f, E, \text{Project}(X, P))$  can be replaced by just  $P$  without changing the semantics. In particular, we have the following corollary.

**COROLLARY 6.20.** *Language  $\text{Sparql}_{As}$  has the same expressive power as  $\text{Sparql}_A$ .*

## 7 PROPERTY PATHS

Property paths were introduced in SPARQL 1.1 to allow for navigational queries, which give users a tool to discover how different resources in an RDF graph are connected. Property paths have been formally studied in recent papers (e.g., see [7, 40, 48]) and their complexity and expressive power is by now relatively well-understood.

Existing work on property paths is, however, restricted either to set semantics [40] or to bag semantics from the early draft of the specification [7, 48], which has changed since then. In this section, we study the normative bag semantics of property paths and describe the interplay between property paths and aggregation.

We start our discussion by introducing property paths syntactically in the algebra.

*Definition 7.1.* Path expressions are inductively defined as follows:

- $Link(u)$  is a path expression if  $u$  is an IRI,
- $Inv(PP)$  is a path expression if so is  $PP$ ,
- $Seq(PP_1, PP_2)$  and  $Alt(PP_1, PP_2)$  are path expressions if so are both  $PP_i$ ,
- $NPS(U)$  is a path expression if  $U$  is a set of IRIs,
- $ZeroOrOnePath(PP)$ ,  $ZeroOrMorePath(PP)$ , and  $OneOrMorePath(PP)$  are path expressions if so is  $PP$ .

Intuitively, a path expression specifies a possible route between two terms in a graph. In particular,  $Link(u)$  represents all pairs of terms that are directly connected via property  $u$ ;  $Inv(PP)$  inverses the pairs represented by  $PP$ ;  $Seq(PP_1, PP_2)$  and  $Alt(PP_1, PP_2)$  correspond to path concatenation and union, respectively;  $NPS(U)$  (standing for “negated property set”) represents pairs connected by a property not mentioned in  $U$ ;  $ZeroOrOnePath(PP)$  represents all pairs connected by  $PP$  and pairs with the same first and second elements; finally,  $ZeroOrMorePath(PP)$  and  $OneOrMorePath(PP)$  represent pairs of terms that are connected by a chain of  $PP$  edges, possibly empty or nonempty, respectively.

The *Path* construct provides a means to build patterns with path expressions.

*Definition 7.2.* Let  $Sparql_\star$  be an extension  $Sparql$ . The language  $Sparql_{\star, PP}$  extends  $Sparql_\star$  by allowing triple patterns of the form  $Path(t_1, PP, t_2)$  where  $t_1, t_2 \in \mathbf{I} \cup \mathbf{L} \cup \mathbf{X}$  and  $PP$  is a path expression that is not of the form  $Link(u)$ ,  $Inv(Link(u))$ , or  $Seq(PP_1, PP_2)$ .<sup>4</sup>

*Example 7.3.* Consider the following query.

(Q6) Compute all managers together with the total salary of the employees managed by them (directly or indirectly).

The query can be written in SPARQL 1.1 as follows:

```
SELECT ?m (SUM(?s) AS ?n)
  WHERE { ?m manages+ ?x . ?x salary ?s }
GROUP BY ?m.
```

Query (Q6) translates into  $Sparql_{As, PP}$  as

$$Project(\{?m, ?n\}, GroupAgg(\{?m\}, ?n, Sum, ?s, \{ Path(?m, OneOrMorePath(Link(manages)), ?x), (?x, salary, ?s) \})).$$

Note that this query combines the navigational features of property paths with grouping and aggregation: property paths are used to identify managers and the employees they manage (directly

<sup>4</sup>The aforementioned path expressions are excluded explicitly in the standard: their counterparts are allowed in SPARQL 1.1 syntax, but are rewritten away in the translation to the algebra. Their inclusion would not alter any results of this paper.

or indirectly), whereas the *GroupAgg* construct is used to group employees according to their manager and to sum over their salaries.

We now specify the normative bag semantics of path expressions and patterns. To handle some cases uniformly, we extend Definition 7.2 and allow  $PP$  in  $Path(t_1, PP, t_2)$  to be arbitrary path expressions, and  $t_1$  and  $t_2$  to be blank nodes.

*Definition 7.4.* The *evaluation*  $[Path(t_1, PP, t_2)]_G^v$  of a triple pattern  $Path(t_1, PP, t_2)$  with property paths over a graph  $G$  in an environment  $v$  is defined inductively as follows, assuming that mappings extend to terms as identity, and  $\text{var}(t_1, t_2)$  is the set of all variables among  $t_1$  and  $t_2$ :

- if  $PP = Link(u)$  for IRI  $u$  then it is  $[(t_1, u, t_2)]_G^v$ ;
- if  $PP = Inv(PP')$  then it is  $[Path(t_2, PP', t_1)]_G^v$ ;
- if  $PP = Seq(PP_1, PP_2)$  then it is  $[Project(\text{var}(t_1, t_2), Join(P_1, P_2))]_G^v$  for  $P_1 = Path(t_1, PP_1, ?x)$ ,  $P_2 = Path(?x, PP_2, t_2)$ , and  $?x$  a fresh variable;
- if  $PP = Alt(PP_1, PP_2)$  then it is  $[Union(P_1, P_2)]_G^v$  for  $P_1 = Path(t_1, PP_1, t_2)$  and  $P_2 = Path(t_1, PP_2, t_2)$ ;
- if  $PP = NPS(U)$  then it is the multiset with the base set

$$\{\mu \mid \exists \mu', \exists u \notin U : \mu = \mu' \cup v, \text{dom}(\mu') = \text{var}(t_1, t_2), (\mu'(t_1), u, \mu'(t_2)) \in G\}$$

and multiplicity 1 for each such  $\mu$ ;

- if  $PP = ZeroOrOnePath(PP')$  then it is the multiset consisting of all mappings  $\mu$  (with multiplicity 1) such that either  $\mu \in [Path(t_1, PP', t_2)]_G^v$  or  $\mu = \mu' \cup v$  where  $\text{dom}(\mu') = \text{var}(t_1, t_2)$  and  $\mu'(t_1) = \mu'(t_2)$ ;
- if  $PP = ZeroOrMorePath(PP')$  then it is the multiset consisting of all mappings  $\mu = \mu' \cup v$  (with multiplicity 1) such that  $\text{dom}(\mu') = \text{var}(t_1, t_2)$  and there is a sequence of terms  $s_0, \dots, s_n$ ,  $n \geq 0$ , with  $s_0 = \mu'(t_1)$ ,  $s_n = \mu'(t_2)$ , and  $[Path(s_i, PP', s_{i+1})]_G^v$  nonempty for all  $0 \leq i < n$ ;
- if  $PP = OneOrMorePath(PP')$  then it is defined exactly the same as for  $ZeroOrMorePath(PP')$ , except that  $n$  is restricted to be strictly greater 0.

The *evaluation* of a BGP with triple patterns  $T_1, \dots, T_n$ , some of which have path expressions, over  $G$  in  $v$  is then

$$[Join(\dots Join(T_1, T_2), \dots T_n)]_G^v.$$

*Example 7.5.* Query (Q6) in Example 7.3 computes mappings of the form  $\{?m \mapsto v, ?n \mapsto w\}$  where  $v$  is an employee and  $w$  is the total salary of all employees managed (directly or via lower-ranked managers) by  $v$ . For instance, if we extend our example graph  $G_{ex}$  with the triples  $(a, \text{manages}, b)$  and  $(a, \text{manages}, c)$  we would obtain  $\{?m \mapsto a, ?n \mapsto 8000\}$  as the only answer mapping.

Moving on to the expressive power of languages with property paths, we first note that all the relevant previous results, that is, Theorems 3.14, 4.13, 4.20, 5.6, 6.11 and 6.13, Corollaries 3.12, 5.7, 6.12 and 6.20, and Proposition 5.9 hold in the presence of property paths, with the proofs extending the existing ones in the obvious way.

On the other hand, using techniques in [45], it is possible to show that query (Q6) from Example 7.3 (as a variation of a reachability query) cannot be expressed in  $\text{Sparql}_{AS}$ , and hence property paths add expressive power to this language. Same can be shown for  $\text{Sparql}$ ,  $\text{Sparql}_{PD}$ ,  $\text{Sparql}_E$  and  $\text{Sparql}_{PDE}$ .

At first glance, the normative algebra allows for an intricate interaction between bag and set semantics in path expressions: some of the constructs keep multiplicities different from 1, while others performs implicit duplicate elimination. However, it is not difficult to see that this impression is wrong, and, as we formalise next, the algebra can be easily simplified in such a way that all paths are evaluated as sets.

*Definition 7.6.* Let *set path expressions* be defined exactly the same as path expressions, but using constructs  $sLink$ ,  $sInv$ , etc. instead of  $Link$ ,  $Inv$ , etc., respectively.

For any extension  $Sparql_{\star}$  of  $Sparql$ , let the language  $Sparql_{\star sPP}$  further extend  $Sparql_{\star}$  by permitting triple patterns of the form  $sPath(t_1, sPP, t_2)$  where  $t_1, t_2 \in I \cup L \cup X$  and  $sPP$  is a set path expression that is not of the form  $sLink(u)$ ,  $sInv(sPP')$ ,  $sSeq(sPP_1, sPP_2)$ , or  $sAlt(sPP_1, sPP_2)$ .

The *evaluation*  $[sPath(t_1, sPP, t_2)]_G^v$  over a graph  $G$  in an environment  $v$  is defined in exactly the same way as for a triple pattern with the corresponding path expression, except that it always works with set semantics, that is,  $card_{\Omega}(\mu) = 1$  for any  $\mu \in [sPath(t_1, sPP, t_2)]_G^v$ .

*Example 7.7.* Consider the  $Sparql_{pp}$  pattern

$$P = Path(?x, Alt(OneOrMorePath(Link(p)), Seq(Link(q), Link(r))), ?y).$$

By definition,  $P$  can be restated as the following  $Sparql_{ppp}$  pattern, for  $?z$  fresh:

$$Union(Path(?x, OneOrMorePath(Link(p)), ?y), Project(\{?x, ?y\}, \{(?x, q, ?z), (?z, r, ?y)\})).$$

Since the remaining path subpattern in this union is equivalent to its set path counterpart  $sPath(?x, sOneOrMorePath(sLink(p)), ?y)$ , it follows that  $P$  can be expressed in  $Sparql_{pspp}$ .

As suggested by the above example, in the presence of projection, path expressions can always be replaced by set path expressions.

**PROPOSITION 7.8.** *For any extension  $Sparql_{\star}$  of  $Sparql_p$ , language  $Sparql_{\star sPP}$  has the same expressive power as  $Sparql_{\star pp}$ .*

The proof of this proposition is in the appendix.

## 8 ANALYTIC AGGREGATE QUERIES

An increasing number of applications of semantic technologies require the analysis of data for effective decision making. In the data warehousing literature this activity is referred to as *online analytical processing (OLAP)* and it involves the execution of complex aggregate queries.

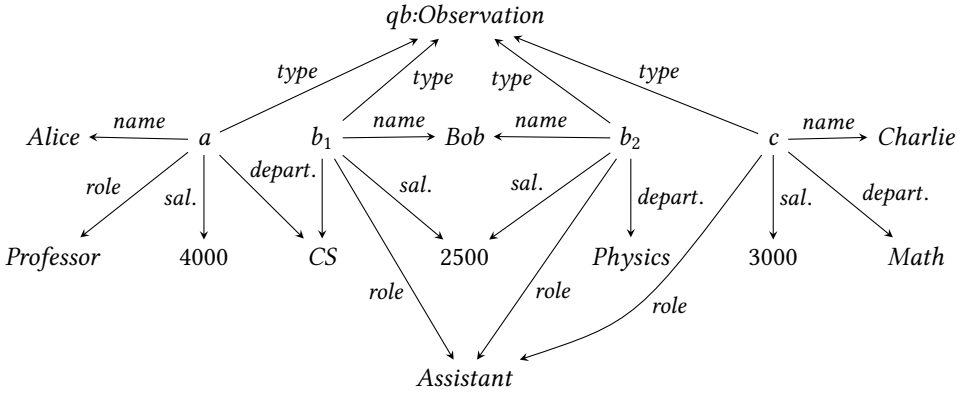
In this section, we briefly recap the multidimensional model underpinning data warehousing applications and existing W3C standards for publishing multidimensional data using RDF. We then give semantics for cube and window queries in SPARQL and finally provide a rewriting of such queries into our algebra  $Sparql_{As}$ .

### 8.1 Data Representation

The natural way of thinking about OLAP queries is in terms of a multidimensional data model, which defines a collection of numeric measures (also called observations) based on a set of dimensions. For instance, in an online retail application a typical measure would be the total amount of merchandise sold in a given currency, and typical dimensions include the kind of product sold, the location and time of the sale, and the method of payment. We can think of the sales information as being arranged in a hypercube, where the coordinates of each sale value are given by specific value IDs for the product, location, time, and method of payment dimensions.

The RDF Data Cube specification [23] is a standard for publishing such multidimensional data on the Semantic Web, where dimensions and measures are represented as RDF properties of a specific kind.

*Example 8.1.* Consider our running example, where we take as measure the employee salary, and define as dimensions the department, the employee name, and the departmental role. The RDF Data Cube vocabulary allows us to define the measure and dimensions using the IRIs  $qb:MeasureProperty$

Fig. 2. Example RDF graph  $G'_{ex}$ 

and  $qb:DimensionProperty$ , respectively. For instance, the salary measure and the role dimension are defined as follows:

$$\begin{aligned}
 (\text{salary}, \text{rdf:type}, qb:\text{MeasureProperty}), & & (\text{salary}, \text{rdfs:range}, \text{xsd:decimal}), \\
 (\text{role}, \text{rdf:type}, qb:\text{DimensionProperty}), & & (\text{role}, \text{rdfs:range}, \text{xsd:string}).
 \end{aligned}$$

Each individual observation is represented as an instance of type  $qb:Observation$ . In the basic case, values of each of the dimensions and measures are attached directly to the observation. The standard introduces a collection of integrity constraints that must be satisfied all datasets. For example, it requires that each observation has a value for each corresponding dimension, or that every dimension has a defined range.

*Example 8.2.* Our running example graph  $G_{ex}$  can be written using the RDF Data Cube Vocabulary as a graph  $G'_{ex}$  specified in Figure 2 (for brevity, triples involving  $rdfs:range$ ,  $qb:MeasureProperty$  and  $qb:DimensionProperty$  are omitted). Note that the data for *Bob* in  $G_{ex}$  has been split into two observations in  $G'_{ex}$ . The W3C specification requires a salary value for each tuple of name, department, and role values; as a result, we need to specify the salary earned by *Bob* for each of its individual departmental affiliations (specifically, 2500 per affiliation).

## 8.2 Cube Queries

In the relational world, a *cube query* involves aggregating over all possible subsets of dimensions at the same time. The query output involves the values of both the measure and the dimensions, and a special symbol is used to indicate that a particular dimension has been aggregated over [26]. In SQL, cube queries are supported by extending the GROUP BY construct with the CUBE keyword, which indicates that grouping must be performed on all subsets of the grouping attributes. In most cases, SQL cube queries are executed over databases that are structured according to a *star schema*—a faithful relational representation of multidimensional data. However, this is not required and the algebraic semantics of cube SQL queries is defined in general for any relational representation of the data.

Our algebra can be extended with a cube operator analogously to the relational case.

*Definition 8.3.* For  $Y$  a set of variables,  $?x$  a variable not in  $Y$ ,  $f$  an aggregate function,  $E$  an expression, and  $P$  a pattern,  $Cube(Y, ?x, f, E, P)$  is a pattern, whose evaluation  $[Cube(Y, ?x, f, E, P)]_G^Y$

over a graph  $G$  in an environment  $\nu$  is as follows, where *all* is a special value not in  $\mathbf{T}$ :

$$\bigsqcup_{Z \subseteq Y} \{ \mu' \mid \mu \in [GroupAgg(Z, ?x, f, E, P)]_G^{\nu}, \mu' = \mu \cup \{?y \mapsto all \mid ?y \in Y \setminus Z\} \}.$$

The special term *all* in the definition is introduced analogously to the relational formalisation (see [26]). Rather than a dedicated value *all*, we could have left the relevant variables unbound; this, however, could yield counter-intuitive results when further applying operators such as *Join*. The semantics of *Cube* suggests a simple translation to our algebra using *GroupAgg*, *Extend*, and *Union*, which is a straightforward adaptation to the SPARQL setting of the well-known translation of cube queries into standard relational algebra operators.

**PROPOSITION 8.4.** *The extension of  $Sparql_{PDEAs}^{\exists}$  that allows for *Cube* has the same expressive power as  $Sparql_{As}^{\exists}$ .*

**PROOF.** By Theorems 6.11 and 6.19 as well as Corollaries 6.12 and 6.20, it is enough to show that for any pattern  $P = Cube(Y, ?x, f, E, P')$  with  $P'$  in the extension of  $Sparql_{PDEAs}^{\exists}$  with *Cube* there exists a pattern in the extension that is  $\nu$ -equivalent to  $P$  for any environment  $\nu$  and has smaller depth of nesting of *Cube*.

Given any  $P = Cube(Y, ?x, f, E, P')$ , let, for  $Z \subseteq Y$ ,

$$P_Z = Extend(?y_1, all, \dots, Extend(?y_n, all, GroupAgg(Z, ?x, f, E, P')) \dots),$$

where  $\{?y_1, \dots, ?y_n\} = Y \setminus Z$ . Then, the pattern

$$Union(P_{Z_1}, \dots, Union(P_{Z_{m-1}}, P_{Z_m}) \dots),$$

where  $Z_1, \dots, Z_m$  are all the subsets of  $Y$ , is as required.  $\square$

For RDF data conforming to the W3C Data Cube specification, we can define a simplified version of the cube operator where the construction of the pattern argument in the operator is implicit. This simplified operator can thus be seen as a “macro” and its semantics is defined by expanding its definition using the generic cube operator.

**Definition 8.5.** For  $D = \{d_1, \dots, d_k\}$  a set of dimension properties (i.e., IRIs),  $?x$  a variable,  $f$  an aggregate function, and  $m$  a measure property (i.e., an IRI), the pattern  $Cube(D, ?x, f, m)$  is an abbreviation for the pattern

$$Cube(\{?d_1, \dots, ?d_k\}, ?x, f, ?m, \{(?z, m, ?m), (?z, d_1, ?d_1), \dots, (?z, d_k, ?d_k)\}).$$

**Example 8.6.** Consider the pattern  $Cube(\{department, role\}, ?x, Sum, salary)$ , which aggregates *salary* over *department* and *role*. Over graph  $G'_{ex}$ , the pattern evaluates to the multiset consisting of the following mappings, all with multiplicity 1:

$$\begin{aligned} &\{?department \mapsto all, ?role \mapsto all, ?x \mapsto 12000\}, \\ &\{?department \mapsto all, ?role \mapsto Professor, ?x \mapsto 4000\}, \\ &\{?department \mapsto all, ?role \mapsto Assistant, ?x \mapsto 8000\}, \\ &\{?department \mapsto CS, ?role \mapsto all, ?x \mapsto 6500\}, \\ &\{?department \mapsto Physics, ?role \mapsto all, ?x \mapsto 2500\}, \\ &\{?department \mapsto Math, ?role \mapsto all, ?x \mapsto 3000\}, \\ &\{?department \mapsto CS, ?role \mapsto Professor, ?x \mapsto 4000\}, \\ &\{?department \mapsto CS, ?role \mapsto Assistant, ?x \mapsto 2500\}, \\ &\{?department \mapsto Physics, ?role \mapsto Assistant, ?x \mapsto 2500\}, \\ &\{?department \mapsto Math, ?role \mapsto Assistant, ?x \mapsto 3000\}. \end{aligned}$$

### 8.3 Window-based Queries

*Window-based* queries are heavily used in data analysis applications. In the relational case, a window identifies a set of rows “around” each individual row in a relation. Once a window has been identified, we can aggregate over the window for each row and extend the row with the result. There is an important difference between groups and windows: the former partition the rows of a relation and compute a value for each partition, whereas the latter compute a different value for each row according to its associated window. Windows were introduced in the SQL:2003 standard and their semantics and efficient evaluation were studied in [16, 42].

The following definition introduces windows in the context of SPARQL. Intuitively, the expression  $F$  in the definition specifies a window as a multiset of “surrounding” mappings for a specific mapping  $\mu$ ; in turn,  $AggWindow$  is used to compute an aggregate value for each mapping based on its corresponding window, which is then appended using a fresh variable (provided that the computed value is not *error*).

*Definition 8.7.* Given a pattern  $P$ , a renaming  $\theta$  of  $\text{var}(P)$  to fresh variables, an expression  $F$  over  $\text{var}(P) \cup \text{var}(P\theta)$ , a variable  $?x \notin \text{var}(P)$ , an aggregate function  $f$ , and an expression  $E$  over  $\text{var}(P)$ ,  $AggWindow(\theta, F, ?x, f, E, P)$  is also a pattern. For a pattern  $AggWindow(\theta, F, ?x, f, E, P)$ , graph  $G$ , environment  $v$ , and mapping  $\mu$ , let

$$v_\mu = f(\llbracket v \mid \mu' \in [P]_G^v, [F]_G^{\mu \cup \mu'\theta} = \text{true}, v = [E]_G^{\mu'} \rrbracket).$$

Then the *evaluation*  $[AggWindow(\theta, F, ?x, f, E, P)]_G^v$  is

$$\llbracket \mu \mid \mu' \in [P]_G^v, v_{\mu'} \neq \text{error}, \mu = \mu' \cup \{?x \mapsto v_{\mu'}\} \rrbracket \uplus \llbracket \mu \mid \mu \in [P]_G^v, v_\mu = \text{error} \rrbracket.$$

*Example 8.8.* Consider the following pattern, which, for every employee, computes the average of all salaries that are lower than their salary:

$$AggWindow(\theta, ?s > ?s', ?n, \text{Avg}, ?s, (?x, \text{salary}, ?s)),$$

where  $\theta$  is the renaming sending each variable  $?u$  to  $?u'$ . For instance, on the graph  $G_{ex}$ , we obtain the mappings:

$$\begin{aligned} \{?x \mapsto a, ?s \mapsto 4000, ?n \mapsto 3000\}, \\ \{?x \mapsto b, ?s \mapsto 5000, ?n \mapsto 3500\}, \\ \{?x \mapsto c, ?s \mapsto 3000, ?n \mapsto 0\}. \end{aligned}$$

Note that in the last mapping  $?n$  is bound to 0 since there are no salaries lower than 3000 (see Definition 6.5).

Finally, note that the  $AggWindow$  pattern can be equivalently expressed in  $\text{Sparql}_{As}$ :

$$\text{GroupAgg}(\{?x, ?s\}, ?n, \text{Avg}, ?s, \text{Filter}(?s > ?s', \text{Join}((?x, \text{salary}, ?s), (?x', \text{salary}, ?s')))).$$

As suggested by the above example,  $AggWindow$  can also be expressed in our algebra. The difficulty in the proof of this result is to show that the exists operator can be eliminated in the presence of windows.

**PROPOSITION 8.9.** *The extension of  $\text{Sparql}_{PDEAs}^{\exists}$  that allows for  $AggWindow$  has the same expressive power as  $\text{Sparql}_{As}$ .*

**PROOF.** We first prove that the extension of  $\text{Sparql}_{PDEAs}^{\exists}$  with  $AggWindow$  has the same expressive power as the extension of  $\text{Sparql}_{PDEAs}$  with only  $AggWindow$ ; then, we show that any  $P = AggWindow(\theta', F, ?x, f, E, P')$  is  $\emptyset$ -equivalent to a  $\text{Sparql}_{PDEAs}$  pattern, which, together with

Theorems 6.11 and 6.19 as well as Corollaries 6.12 and 6.20, would imply the statement of the theorem.

For the first part, we need to generalise Theorem 4.20 to the extension of  $\text{Sparql}_{PDEAs}$  with  $\text{AggWindow}$ . By Theorem 6.19, it is enough to generalise it to the extension of  $\text{Sparql}_{PDEA}$  with  $\text{AggWindow}$ . Lemmas 4.16 and 4.17 hold for this language, so we only need to generalise Lemma 4.18. By Corollary 6.12, it is already generalised to  $\text{Sparql}_{PDEA}$ , so we just need to extend it to  $\text{AggWindow}$ . Let

$$\text{Sub}_\theta(\text{AggWindow}(\theta', F, ?x, f, E, P'), P_1) = \text{AggWindow}(\theta'', F \wedge \text{eq}(\text{var}(P_1), \theta''), ?x, f, E, \text{Sub}_\theta(P', P_1)),$$

where  $\theta''$  is an extension of  $\theta'$  to  $\text{var}(\text{Sub}_\theta(P', P_1)) \setminus \text{var}(P')$  that sends all these variables to fresh copies. Claim 4.19 then generalises: we just need to consider yet another case for the induction step.

For the second part of the proof, consider a pattern  $P = \text{AggWindow}(\theta', F, ?x, f, E, P')$ . It is  $\emptyset$ -equivalent to the pattern

$$\text{Project}(\text{var}(P') \cup \{?x\}, \text{Filter}(\text{eq}(\text{var}(P'), \theta''), \text{Join}(P', P''\theta''))),$$

with  $\theta''$  another renaming of  $\text{var}(P')$  to fresh variables and

$$P'' = \text{GroupAgg}(\text{var}(P'), ?x, f, E, \text{Filter}(F, \text{Join}(P', P'\theta'))).$$

Note that  $P''$  and  $P$  coincide if interpreted as sets; the additional transformations are applied to  $P''$  to obtain the correct multiplicities.  $\square$

Note that the  $\emptyset$ -equivalent patterns in the second part of the proof of Proposition 8.9 may not be  $\nu$ -equivalent for some environments  $\nu$ , so a straightforward proof such as the one for Proposition 8.4 does not work and the first part is required.

## 9 ADDITIONAL CONSIDERATIONS

We made several simplifying assumptions that made us deviate from the standard. First, we omitted the auxiliary SPARQL 1.0 pattern operator *Diff* and SPARQL 1.1 operator *Minus*. Second, we have considered only expressions already available in SPARQL 1.0, whereas SPARQL 1.1 defines a richer language for expressions. Third, we have assumed that triple patterns contain no blank nodes, which are, however, allowed in the standard. Fourth, we have assumed that the result of a query is a multiset of mappings, where the standard defines it as a list. Fifth, we did not consider counting functions aggregating over a special symbol  $*$  instead of a list of expressions. Finally, we have omitted the non-deterministic aggregate functions *Sample* and *GroupConcat*.

The non-deterministic aggregate functions can be treated similarly to the other aggregate functions. We next discuss how the other five assumptions affect our results.

*Difference and Minus.* The pattern operator *Diff* in the SPARQL 1.0 algebra does not have a counterpart in the SPARQL 1.0 syntax, but plays an auxiliary role in the definition of the semantics of *LeftJoin*. Instead, a slightly different operator *Minus* is introduced in SPARQL 1.1. Next, we give the formal definitions of these two difference operators and present ways to express them via the core Sparql operators.

*Definition 9.1.* If  $E$  is an expression and  $P_1, P_2$  are patterns, then  $\text{Diff}(E, P_1, P_2)$  is a pattern with the following semantics, for a graph  $G$  and environment  $\nu$ :

$$[\text{Diff}(E, P_1, P_2)]_G^\nu = \{\mu \mid \mu \in [P_1]_G^\nu, \forall \mu_2 \in [P_2]_G^\nu. (\mu \approx \mu_2 \text{ or } [E]_G^{\mu \cup \mu_2} \neq \text{true})\}.$$

If  $P_1$  and  $P_2$  are patterns, then  $Minus(P_1, P_2)$  is a pattern with the following semantics, for a graph  $G$  and environment  $\nu$ :

$$[Minus(P_1, P_2)]_G^\nu = \{ \mu \mid \mu \in [P_1]_G^\nu, \forall \mu_2 \in [P_2]_G^\nu, (\mu \approx \mu_2 \text{ or } (\text{dom}(\mu) \cap \text{dom}(\mu_2)) \setminus \text{dom}(\nu) = \emptyset) \}.$$

Note that the semantics of *LeftJoin* can be stated using *Diff* as follows:

$$[LeftJoin(E, P_1, P_2)]_G^\nu = [Filter(E, Join(P_1, P_2))]_G^\nu \cup [Diff(E, P_1, P_2)]_G^\nu.$$

This is the only way in which *Diff* is used in the specification.

The following example illustrates the difference between *Diff*, *Minus*, and operator *SetMinus*, introduced in Section 3.1.

*Example 9.2.* Let  $P_1, P_2$  be patterns and  $G$  a graph such that  $[P_1]_G^0$  and  $[P_2]_G^0$  are the multisets with base sets  $\{\mu_1, \mu'_1\}$  and  $\{\mu_2\}$ , respectively, where

$$\mu_1 = \{?x \mapsto a\}, \quad \mu'_1 = \{?x \mapsto a, ?y \mapsto b\}, \quad \mu_2 = \{?y \mapsto b\},$$

and all mappings have multiplicity 1. Then,  $[Diff(true, P_1, P_2)]_G^0$  is empty,  $[Minus(P_1, P_2)]_G^0$  consists of a single occurrence of  $\mu_1$ , while  $[SetMinus(P_1, P_2)]_G^0 = [P_1]_G^0$ , that is, has both  $\mu_1$  and  $\mu'_1$ .

As we already said, both *Diff* and *Minus* are expressible via the core Sparql operators, that is, these operators do not add any expressive power to the languages considered in this paper. In particular, it is well-known that, for an expression  $E$  and patterns  $P_1, P_2$  in any extension of Sparql,  $Diff(E, P_1, P_2)$  is  $\nu$ -equivalent, for any  $\nu$ , to the pattern

$$Filter(\neg \text{bound}(?x), LeftJoin(E, P_1, Join(P_2, (?x, ?y, ?z)))),$$

where  $?x, ?y$ , and  $?z$  are fresh variables (see, e.g., [9]). In turn,  $Minus(P_1, P_2)$  is  $\emptyset$ -equivalent to  $Diff(E, P_1, P_2\theta)$ , where  $\theta$  is a variable renaming of  $\text{var}(Minus(P_1, P_2))$  and

$$E = \bigwedge_{?x \in \text{dom}(\theta)} (\text{bound}(?x) \wedge \text{bound}(?x\theta) \rightarrow ?x = ?x\theta) \wedge \bigvee_{?x \in \text{dom}(\theta)} (\text{bound}(?x) \wedge \text{bound}(?x\theta)).$$

This equivalence is introduced in [37] for set semantics, but an immediate inspection shows that it holds under bag semantics as well. To deal with *Minus* in non-empty environments, we first have to eliminate exists, as it is done in Theorem 4.20, and then apply the equivalence for the empty environment. To this end, we show the following proposition.

**PROPOSITION 9.3.** *Let Sparql $_\star$  be any language considered in this paper. The extension of Sparql $_\star$  allowing for exists and Minus has the same expressive power as Sparql $_\star$ .*

The proof of this proposition is in the appendix.

*Expressions.* We focus on two constructs due to their potential implications: ternary if, which computes one of two expressions depending on the evaluation of a third one, and coalesce, which allows us to “recover” from errors.

*Definition 9.4.* If  $E_1, E_2$  and  $E_3$  are expressions then  $\text{if}(E_1, E_2, E_3)$  is an expression with the following semantics: for a given  $\mu$  and  $G$ , the *evaluation*  $[\text{if}(E_1, E_2, E_3)]_G^\mu$  is  $[E_2]_G^\mu$  if  $[E_1]_G^\mu = \text{true}$ , it is  $[E_3]_G^\mu$  if  $[E_1]_G^\mu = \text{false}$ , and *error* otherwise. Let also  $\text{var}(\text{if}(E_1, E_2, E_3)) = \text{var}(E_1) \cup \text{var}(E_2) \cup \text{var}(E_3)$ .

For  $\mathbf{E} = \langle E_1, \dots, E_n \rangle$  a non-empty list of expressions,  $\text{coalesce}(\mathbf{E})$  is an expression with the following semantics: for  $\mu$  and  $G$ , the *evaluation*  $[\text{coalesce}(\mathbf{E})]_G^\mu$  is *error* if  $[E_i]_G^\mu = \text{error}$  for each

$1 \leq i \leq n$ , and  $[E_j]_G^\mu$  otherwise, for the smallest  $j$  with  $[E_j]_G^\mu \neq \text{error}$ . Let also  $\text{var}(\text{coalesce}(\mathbf{E})) = \text{var}(E_1) \cup \dots \cup \text{var}(E_n)$ .

*Example 9.5.* The expression  $\text{if}(?h < 12, \text{“am”}, \text{“pm”})$  returns the string “am” if the value of the variable  $?h$  is less than 12, and “pm” otherwise;  $\text{coalesce}(\langle ?x, 0 \rangle)$  returns the value of  $?x$  if  $?x$  is bound and 0 otherwise.

These expressions can be rewritten in terms of Sparql expressions, and hence their introduction is immaterial to our results in this paper. Indeed, for any pattern  $P$  and expression  $E$  with a subexpression  $\text{if}(E_1, E_2, E_3)$ , for expressions  $E_1, E_2$ , and  $E_3$ , all in any extension of Sparql, as well as for any environment  $\nu$ , the pattern  $\text{Filter}(E, P)$  is  $\nu$ -equivalent to

$$\text{Union}(\text{Union}(\text{Filter}(E_1 \wedge E_2', P), \text{Filter}(\neg E_1 \wedge E_3', P)), \\ \text{Filter}(E_1', \text{SetMinus}(P, \text{Filter}(E_1 \vee \neg E_1, P)))),$$

where  $E_i'$ , for  $i = 1, 2, 3$ , is obtained from  $E$  by replacing  $\text{if}(E_1, E_2, E_3)$  by  $E_i$ .

Similarly, for any  $P$  and  $E$  with a subexpression  $\text{coalesce}(\langle E_1, E_2 \rangle)$ , for expressions  $E_1$  and  $E_2$ , all in any extension of Sparql, as well as for any environment  $\nu$ , the pattern  $\text{Filter}(E, P)$  is  $\nu$ -equivalent to

$$\text{Union}(\text{Filter}((E_1 \vee \neg E_1) \wedge E_1', P), \text{Filter}(E_2', \text{SetMinus}(P, \text{Filter}(E_1 \vee \neg E_1, P)))),$$

where  $E_i'$ , for  $i = 1, \dots, 2$ , is obtained from  $E$  by replacing  $\text{coalesce}(\langle E_1, E_2 \rangle)$  by  $E_i$ ;  $\text{coalesce}$  over longer lists is treated analogously (while  $\text{coalesce}$  over singleton lists is vacuous). Similarly,  $\text{if}$  and  $\text{coalesce}$  can be eliminated from other constructs that mention expressions in all the languages considered in this paper, such as *LeftJoin* and *GroupAgg*. Therefore, these expression constructs do not increase the expressive power of any language considered extension of Sparql.

*Blank Nodes.* According to the SPARQL 1.1 standard, triple patterns in BGPs may contain blank nodes in subject and object positions. Roughly speaking, such blank nodes are treated as variables for the purposes of pattern-graph matching; in contrast to variables, however, the scope of blank nodes is confined to the BGP in which they occur. The effect of blank nodes within BGPs can be simulated in our algebra by using *Project*: any BGP  $B$  with variables  $X$  and blank nodes  $N$  is  $\nu$ -equivalent, for any  $\nu$ , to the pattern  $\text{Project}(X, B\xi)$ , where  $\xi$  is a renaming of  $N$  to fresh variables. As shown in Section 3 (Corollary 3.12), projection in patterns does not add expressive power to Sparql, and hence neither does allowing blank nodes in BGPs. Note, however, that in combination with *Distinct* at pattern level, blank nodes do lead to an increase in expressive power since they introduce projection, and hence the proof of Theorem 3.14 can be easily adapted.

*Lists of Solutions.* We have defined the semantics of patterns and queries uniformly in terms of multisets of mappings, which simplifies the algebra by avoiding numerous type conversions. The standard, however, defines query evaluation in terms of lists of mappings. It also allows for solution modifiers, such as *Slice* and *OrderBy*, the semantics of which depends on the order of mappings in the solution sequence of the query.

We next argue that dispensing with lists altogether does not essentially affect any of our results. Given a multiset  $\Omega$  of mappings, let  $\Lambda_\Omega$  be the set of all lists of mappings that coincide with  $\Omega$  when disregarding the order of elements. Also, let  $h$  be the function that translates queries from any language  $\text{Sparql}_\star$  considered in this paper into the normative one by adding the necessary type conversions between multisets and lists. Then for every  $\text{Sparql}_\star$  query  $Q$ , and graph  $G$ , we have  $[Q]_G = \Omega$  if and only if  $[h(Q)]_G^{\text{std}} \in \Lambda_\Omega$ , where  $[\cdot]^{\text{std}}$  is the list evaluation function defined in the SPARQL standard. This correspondence demonstrates an additional benefit of using multisets rather than lists: every query evaluates to a *unique* multiset (except those using the non-deterministic

aggregate Sample), whereas the evaluation to a list of solutions is non-deterministic even for simple queries.

*Counting Aggregates over \**. Apart from aggregates of the form  $Aggregate(E, f, \Gamma)$  with  $E$  a list of expressions, SPARQL 1.1 allows for aggregates  $A = Aggregate(*, f, \Gamma)$ , where  $*$  is a distinguished symbol,  $f$  is Count or CountD, and  $\Gamma = Group(F, P)$  is a group. The evaluation  $[A]_G$  of such  $A$  over a graph  $G$  is the partial function from  $v$ -lists to values defined as follows, for each  $Key$  in the domain of  $[\Gamma]_G$ : if  $f$  is Count, then  $[A]_G(Key)$  is the sum of the cardinalities of all mappings in  $[\Gamma]_G(Key)$ , and if  $f$  is CountD, then it is the number of distinct mappings in  $[\Gamma]_G(Key)$ . These aggregates can be rewritten to the algebra in Section 6.1, since  $Aggregate(*, Count, \Gamma)$  is equivalent to  $Aggregate(\text{bound}(?x) \vee \neg\text{bound}(?x), Count, \Gamma)$ , for  $?x$  an arbitrary variable, while  $Aggregate(*, CountD, Group(F, P))$  is equivalent to  $Aggregate(*, Count, Group(F, Distinct(P)))$ . Note that the expression  $\text{bound}(?x) \vee \neg\text{bound}(?x)$ , which always evaluates to *true*, is needed in the first equivalence to avoid disregarding mappings with unbound  $?x$ . Also,  $?x$  is not required to appear in  $\Gamma$ ; in fact, if  $?x$  is a fresh variable, then the expression can be simplified to  $\neg\text{bound}(?x)$ .

## 10 RELATED WORK

Query nesting, variable assignment and aggregation are key features of the database query language SQL (e.g., see [25] for an excellent textbook covering these topics). In addition to being widely-used in practice, the formal properties of query nesting and aggregation are also well-understood from a formal point of view in the relational setting. These include semantics and expressive power [18, 27, 28, 30, 44–47, 54]; query equivalence [17, 21]; and query rewriting using views [20]. Analytic OLAP queries have also been studied extensively in the relational case (e.g., see [16, 26, 42]).

Out of the above works, the most closely related to ours is arguably the line of research in [27, 28, 46, 47], which studies theoretical properties of bag algebras. In particular, it studies the relative expressive power of bag algebra primitives, the relation to set algebras, with and without aggregation, and establishes a LOGSPACE data complexity bound for query answering. These results are complementary to the ones in our paper, and we expect that many of them can be transferred to the setting of SPARQL.

In contrast to the relational case, query nesting and aggregation are relatively new features of SPARQL that were only standardised in its latest version. To the best of our knowledge, the formal properties of query nesting have only been discussed in [4] and more recently in [5]; a detailed discussion of these works has been provided in previous sections. We are not aware of any papers discussing formal aspects of aggregation in SPARQL 1.1; to the best of our knowledge, the efficient implementation and optimisation of SPARQL 1.1 aggregate queries also constitutes a new and open research area [32].

Analytic and OLAP queries in the context of Semantic Web applications have also been considered only recently. Colazzo et al. [22] proposed a semantics for analytical queries and OLAP operators in the context of RDF; their work, however, does not provide a clear connection with the semantics of aggregation in SPARQL 1.1 and does not study the expressive power of the introduced operators. The RDF Cube specification [23] provides a vocabulary and specific guidelines for representing multidimensional data in RDF; practical aspects and first concrete implementations of the standard have been recently discussed in [33, 35, 55].

This paper extends a conference publication [34] by providing extended discussions and examples as well as the following major technical contributions: (i) complete proofs of all our results; (ii) a principled fix to the normative semantics of exists expressions, and a comprehensive analysis of the impact of this new semantics on the expressive power of query nesting and aggregation

$$\begin{array}{c}
\left\{ \text{Sparql}, \text{Sparql}_P, \text{Sparql}_D, \right\} < \left\{ \text{Sparql}_{PD}, \right\} \\
\left\{ \text{Sparql}^\exists, \text{Sparql}_P^\exists, \text{Sparql}_D^\exists \right\} < \left\{ \text{Sparql}_{PD}^\exists \right\} \\
\wedge \\
\left\{ \text{Sparql}_E, \text{Sparql}_{PE}, \text{Sparql}_{DE}, \right\} < \left\{ \text{Sparql}_{PDE}, \right\} < \left\{ \begin{array}{l} \text{Sparql}_A, \dots, \text{Sparql}_{PDEA}, \\ \text{Sparql}_A^\exists, \dots, \text{Sparql}_{PDEA}^\exists, \\ \text{Sparql}_{As}, \dots, \text{Sparql}_{PDEAs}, \\ \text{Sparql}_{As}^\exists, \dots, \text{Sparql}_{PDEAs}^\exists \end{array} \right\} \\
\left\{ \text{Sparql}_E^\exists, \text{Sparql}_{PE}^\exists, \text{Sparql}_{DE}^\exists \right\} < \left\{ \text{Sparql}_{PDE}^\exists \right\} < \left\{ \text{Sparql}_{PDEAs}^\exists \right\}
\end{array}$$

Fig. 3. Main expressivity results of this paper. Here, Sparql is the core SPARQL 1.0 algebra, while indices abbreviate its extensions as follows:  $P$  stands for subqueries with *Project*,  $D$  for subqueries with *Distinct*,  $\exists$  for subqueries with exists,  $E$  for variable assignment *Extend*,  $A$  for aggregation *AggregateJoin*, and  $As$  for normalised aggregation *GroupAgg*; languages with the same expressive power are grouped together by  $\{\}$ , while “strictly more expressive” relation is denoted by  $<$ .

constructs; (iii) treatment of property paths; and (iv) extended treatment of OLAP queries and multidimensional data representation standards for RDF.

## 11 CONCLUSION AND FUTURE WORK

We have presented an in-depth analysis of the SPARQL 1.1 algebra. Our investigation has shed light on the complex inter-dependencies between the algebraic operators that enable query nesting, variable assignment, and aggregation, which are critical to many emerging applications of semantic technologies. Our main results on relative expressivity of the fragments of SPARQL 1.1 algebra are summarised in Figure 3.

We see many possible avenues for future work. We are planning to study the interaction between aggregation and query nesting operators with other features of SPARQL 1.1 such as entailment regimes and query federation. Furthermore, there have been proposals for an extension of SPARQL with stream reasoning and event processing features [6, 11, 15] as well as with analytical queries [22]; it would be interesting to study the connections between these languages and the SPARQL 1.1 normative algebra.

## ELECTRONIC APPENDIX

The electronic appendix for this article can be accessed in the ACM Digital Library.

## REFERENCES

- [1] Alberto Abelló, Oscar Romero, Torben Bach Pedersen, Rafael Berlanga Llavori, Victoria Nebot, María José Aramburu Cabo, and Alkis Simitsis. 2015. Using Semantic Web Technologies for Exploratory OLAP: A Survey. *IEEE TKDE* 27, 2 (2015), 571–588.
- [2] Shqiponja Ahmetaj, Wolfgang Fischl, Reinhard Pichler, Mantas Simkus, and Sebastian Skritek. 2015. Towards Reconciling SPARQL and Certain Answers. In *WWW 2015*. ACM, 23–33.
- [3] Renzo Angles and Claudio Gutierrez. 2008. The Expressive Power of SPARQL. In *ISWC 2008 (LNCS)*, Vol. 5318. Springer, 114–129.
- [4] Renzo Angles and Claudio Gutierrez. 2011. Subqueries in SPARQL. In *AMW 2011 (CEUR Workshop Proceedings)*, Vol. 749. CEUR-WS.org.
- [5] Renzo Angles and Claudio Gutierrez. 2016. Negation in SPARQL. *CoRR* abs/1603.06053 (2016).
- [6] Darko Anicic, Paul Fodor, Sebastian Rudolph, and Nenad Stojanovic. 2011. EP-SPARQL: A Unified Language for Event Processing and Stream Reasoning. In *WWW 2011*. ACM, 635–644.
- [7] Marcelo Arenas, Sebastián Conca, and Jorge Pérez. 2012. Counting Beyond a Yottabyte, or How SPARQL 1.1 Property Paths will Prevent Adoption of the Standard. In *WWW 2012*. ACM, 629–638.

- [8] Marcelo Arenas, Georg Gottlob, and Andreas Pieris. 2014. Expressive Languages for Querying the Semantic Web. In *PODS 2014*. ACM, 14–26.
- [9] Marcelo Arenas and Jorge Pérez. 2011. Querying Semantic Web Data with SPARQL. In *PODS 2011*. ACM, 305–316.
- [10] Elham Akbari Azirani, François Goasdoué, Ioana Manolescu, and Alexandra Roatis. 2015. Efficient OLAP Operations for RDF Analytics. In *ICDE Workshops 2015*. IEEE Computer Society, 71–76.
- [11] Davide Francesco Barbieri, Daniele Braga, Stefano Ceri, Emanuele Della Valle, and Michael Grossniklaus. 2010. C-SPARQL: a Continuous Query Language for RDF Data Streams. *Int. J. Semantic Comput.* 4, 1 (2010), 3–25.
- [12] Christian Bizer, Tom Heath, and Tim Berners-Lee. 2009. Linked Data - The Story So Far. *Int. J. Semantic Web Inf. Syst.* 5, 3 (2009), 1–22.
- [13] Carlos Buil Aranda, Marcelo Arenas, Óscar Corcho, and Axel Polleres. 2013. Federating Queries in SPARQL 1.1: Syntax, Semantics and Evaluation. *J. Web Sem.* 18, 1 (2013), 1–17.
- [14] Carlos Buil Aranda, Axel Polleres, and Jürgen Umbrich. 2014. Strategies for Executing Federated Queries in SPARQL1.1. In *ISWC 2014, Part II (LNCS)*, Vol. 8797. Springer, 390–405.
- [15] Jean-Paul Calbimonte, Hoyoung Jeung, Óscar Corcho, and Karl Aberer. 2012. Enabling Query Technologies for the Semantic Sensor Web. *Int. J. Semantic Web Inf. Syst.* 8, 1 (2012), 43–63.
- [16] Yu Cao, Chee-Yong Chan, Jie Li, and Kian-Lee Tan. 2012. Optimization of Analytic Window Functions. *PVLDB* 5, 11 (2012), 1244–1255.
- [17] Sara Cohen. 2005. Containment of Aggregate Queries. *SIGMOD Record* 34, 1 (2005), 77–85.
- [18] Sara Cohen. 2006. Equivalence of Queries Combining Set and Bag-Set Semantics. In *PODS 2006*. ACM, 70–79.
- [19] Sara Cohen. 2009. Equivalence of Queries that are Sensitive to Multiplicities. *VLDB J.* 18, 3 (2009), 765–785.
- [20] Sara Cohen, Werner Nutt, and Yehoshua Sagiv. 2006. Rewriting Queries with Arbitrary Aggregation Functions Using Views. *ACM TODS* 31, 2 (2006), 672–715.
- [21] Sara Cohen, Werner Nutt, and Yehoshua Sagiv. 2007. Deciding Equivalences among Conjunctive Aggregate Queries. *J. ACM* 54, 2, Article 5 (2007), 50 pages.
- [22] Dario Colazzo, François Goasdoué, Ioana Manolescu, and Alexandra Roatis. 2014. RDF Analytics: Lenses over Semantic Graphs. In *WWW 2014*. ACM, 467–478.
- [23] Richard Cyganiak and Dave Reynolds. 2014. *The RDF Data Cube Vocabulary*. W3C Recommendation. W3C.
- [24] Lorena Etcheverry and Alejandro A. Vaisman. 2012. Enhancing OLAP Analysis with Web Cubes. In *ESWC 2012 (LNCS)*, Vol. 7295. Springer, 469–483.
- [25] Hector Garcia-Molina, Jeffrey D. Ullman, and Jennifer Widom. 2009. *Database Systems: The Complete Book* (2nd ed.). Pearson Education.
- [26] Jim Gray, Surajit Chaudhuri, Adam Bosworth, Andrew Layman, Don Reichart, Murali Venkatrao, Frank Pellow, and Hamid Pirahesh. 1997. Data Cube: A Relational Aggregation Operator Generalizing Group-by, Cross-Tab, and Sub Totals. *Data Min. Knowl. Discov.* 1, 1 (1997), 29–53.
- [27] Stéphane Grumbach, Leonid Libkin, Tova Milo, and Limsoon Wong. 1996. Query Languages for Bags: Expressive Power and Complexity. *SIGACT News* 27, 2 (1996), 30–44.
- [28] Stéphane Grumbach and Tova Milo. 1996. Towards Tractable Algebras for Bags. *J. Comput. Syst. Sci.* 52, 3 (1996), 570–588.
- [29] Steve Harris and Andy Seaborne. 2013. *SPARQL 1.1 Query Language*. W3C Recommendation. W3C.
- [30] Lauri Hella, Leonid Libkin, Juha Nurmonen, and Limsoon Wong. 2001. Logics with Aggregate Operators. *J. ACM* 48, 4 (2001), 880–907.
- [31] Ivan Herman. 2016. *Errata in SPARQL 1.1*. v1.44. W3C. <http://www.w3.org/2013/sparql-errata>
- [32] Dilshod Ibragimov, Katja Hose, Torben Bach Pedersen, and Esteban Zimányi. 2015. Processing Aggregate Queries in a Federation of SPARQL Endpoints. In *ESWC 2015 (LNCS)*, Vol. 9088. Springer, 269–285.
- [33] Evangelos Kalampokis, Andriy Nikolov, Peter Haase, Richard Cyganiak, Arkadiusz Stasiewicz, Areti Karamanou, Maria Zotou, Dimitris Zeginis, Efthimios Tambouris, and Konstantinos A. Tarabanis. 2014. Exploiting Linked Data Cubes with OpenCube Toolkit. In *ISWC 2014 (Posters & Demonstrations) (CEUR Workshop Proceedings)*, Vol. 1272. CEUR-WS.org, 137–140.
- [34] Mark Kaminski, Egor V. Kostylev, and Bernardo Cuenca Grau. 2016. Semantics and Expressive Power of Subqueries and Aggregates in SPARQL 1.1. In *WWW 2016*. ACM, 227–238.
- [35] Yasar Khan, Muhammad Saleem, Aftab Iqbal, Muntazir Mehdi, Aidan Hogan, Axel-Cyrille Ngonga Ngomo, Stefan Decker, and Ratnesh Sahay. 2014. SAFE: Policy Aware SPARQL Query Federation Over RDF Data Cubes. In *SWAT4LS 2014 (CEUR Workshop Proceedings)*, Vol. 1320. CEUR-WS.org.
- [36] Anthony C. Klug. 1982. Equivalence of Relational Algebra and Relational Calculus Query Languages Having Aggregate Functions. *J. ACM* 29, 3 (1982), 699–717.

- [37] Roman Kontchakov and Egor V. Kostylev. 2016. On Expressibility of Non-Monotone Operators in SPARQL. In *KR 2016*. AAAI Press, 369–379.
- [38] Roman Kontchakov, Martin Rezk, Mariano Rodriguez-Muro, Guohui Xiao, and Michael Zakharyashev. 2014. Answering SPARQL Queries over Databases under OWL 2 QL Entailment Regime. In *ISWC 2014, Part I (LNCS)*, Vol. 8796. Springer, 552–567.
- [39] Egor V. Kostylev and Bernardo Cuenca Grau. 2014. On the Semantics of SPARQL Queries with Optional Matching under Entailment Regimes. In *ISWC 2014, Part II (LNCS)*, Vol. 8797. Springer, 374–389.
- [40] Egor V. Kostylev, Juan L. Reutter, Miguel Angel Romero Orth, and Domagoj Vrgoc. 2015. SPARQL with Property Paths. In *ISWC 2015, Part I (LNCS)*, Vol. 9366. Springer, 3–18.
- [41] Egor V. Kostylev, Juan L. Reutter, and Martín Ugarte. 2015. CONSTRUCT Queries in SPARQL. In *ICDT 2015 (LIPIcs)*, Vol. 31. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, 212–229.
- [42] Viktor Leis, Kan Kundhikanjana, Alfons Kemper, and Thomas Neumann. 2015. Efficient Processing of Window Functions in Analytical SQL Queries. *PVLDB* 8, 10 (2015), 1058–1069.
- [43] Andrés Letelier, Jorge Pérez, Reinhard Pichler, and Sebastian Skritek. 2013. Static Analysis and Optimization of Semantic Web Queries. *ACM TODS* 38, 4, Article 25 (2013), 45 pages.
- [44] Leonid Libkin. 2000. Logics with Counting and Local Properties. *ACM TOCL* 1, 1 (2000), 33–59.
- [45] Leonid Libkin. 2003. Expressive Power of SQL. *Theor. Comput. Sci.* 296, 3 (2003), 379–404.
- [46] Leonid Libkin and Limsoon Wong. 1994. New Techniques for Studying Set Languages, Bag Languages and Aggregate Functions. In *PODS 1994*. ACM, 155–166.
- [47] Leonid Libkin and Limsoon Wong. 1997. Query Languages for Bags and Aggregate Functions. *J. Comput. Syst. Sci.* 55, 2 (1997), 241–272.
- [48] Katja Losemann and Wim Martens. 2012. The complexity of evaluating path expressions in SPARQL. In *PODS 2012*. ACM, 101–112.
- [49] Frank Manola and Eric Miller. 2004. *RDF Primer*. W3C Recommendation. W3C.
- [50] Jorge Pérez, Marcelo Arenas, and Claudio Gutierrez. 2009. Semantics and Complexity of SPARQL. *ACM TODS* 34, 3, Article 16 (2009), 45 pages.
- [51] Axel Polleres and Johannes Peter Wallner. 2013. On the Relation Between SPARQL1.1 and Answer Set Programming. *J. Appl. Non-Class. Log.* 23, 1-2 (2013), 159–212.
- [52] Eric Prud'hommeaux and Andy Seaborne. 2008. *SPARQL Query Language for RDF*. W3C Recommendation. W3C.
- [53] Michael Schmidt, Michael Meier, and Georg Lausen. 2010. Foundations of SPARQL Query Optimization. In *ICDT 2010*. ACM, 4–33.
- [54] Nicole Schweikardt. 2005. Arithmetic, First-Order Logic, and Counting Quantifiers. *ACM TOCL* 6, 3 (2005), 634–671.
- [55] Efthimios Tambouris, Evangelos Kalampokis, and Konstantinos A. Tarabanis. 2015. Processing Linked Open Data Cubes. In *EGOV 2015 (LNCS)*, Vol. 9248. Springer, 130–143.
- [56] Xiaowang Zhang and Jan Van den Bussche. 2014. On the Primitivity of Operators in SPARQL. *Inf. Process. Lett.* 114, 9 (2014), 480–485.
- [57] Xiaowang Zhang and Jan Van den Bussche. 2015. On the Power of SPARQL in Expressing Navigational Queries. *Comput. J.* 58, 11 (2015), 2841–2851.
- [58] Peixiang Zhao, Xiaolei Li, Dong Xin, and Jiawei Han. 2011. Graph Cube: On Warehousing and OLAP Multidimensional Networks. In *SIGMOD 2011*. ACM, 853–864.

Received July 2016; revised February 2017; accepted April 2017

## Online Appendix to: Query Nesting, Assignment, and Aggregation in SPARQL 1.1

Mark Kaminski, Egor V. Kostylev, and Bernardo Cuenca Grau, University of Oxford

---

This appendix contains all the technical proofs omitted in the main part of the paper.

---

### A PROOFS FOR SECTION 3

CLAIM 3.4. *Patterns  $\text{SetMinus}(P_1, P_2)$  and  $P$  are equivalent.*

PROOF. For the forward direction, assume that  $\mu \in \llbracket \text{SetMinus}(P_1, P_2) \rrbracket_G$  for a graph  $G$ . We prove that  $\mu \in \llbracket P \rrbracket_G$  and

$$\text{card}_{\llbracket P \rrbracket_G}(\mu) = \text{card}_{\llbracket \text{SetMinus}(P_1, P_2) \rrbracket_G}(\mu).$$

By the definition of *SetMinus* we have that  $\mu \in \llbracket P_1 \rrbracket_G$  and  $\mu \notin \llbracket P_2 \rrbracket_G$ . This implies  $\llbracket \text{eq}(X, \theta) \rrbracket_G^{\mu \cup \mu'} = \text{false}$  for every  $\mu' \in \llbracket P'_2 \theta \rrbracket_G$ ; but then, the definition of  $P'$  and the semantics of the *LeftJoin* operator ensure that  $\mu \in \llbracket P' \rrbracket_G$ , and, moreover,  $\text{card}_{\llbracket P' \rrbracket_G}(\mu) = \text{card}_{\llbracket P_1 \rrbracket_G}(\mu) = \text{card}_{\llbracket \text{SetMinus}(P_1, P_2) \rrbracket_G}(\mu)$  (note that  $\{?x, ?y, ?z\}$  always matches  $G$ , which is non-empty by definition). Finally, since  $?x$  is fresh, that is, unbound in  $\mu$ , we have that  $\mu \in \llbracket P \rrbracket_G$  and  $\text{card}_{\llbracket P \rrbracket_G}(\mu) = \text{card}_{\llbracket P' \rrbracket_G}(\mu)$ , as required.

For the other direction, assume that  $\mu \in \llbracket P \rrbracket_G$  for a graph  $G$ . We show that  $\mu \in \llbracket \text{SetMinus}(P_1, P_2) \rrbracket_G$  and

$$\text{card}_{\llbracket \text{SetMinus}(P_1, P_2) \rrbracket_G}(\mu) = \text{card}_{\llbracket P \rrbracket_G}(\mu).$$

Since  $\mu \in \llbracket P \rrbracket_G$ , variable  $?x$  is unbound in  $\mu$ . Thus, by the semantics of *LeftJoin* and the fact that  $G$  has at least one triple,  $\mu \in \llbracket P_1 \rrbracket_G$ ,  $\text{card}_{\llbracket P_1 \rrbracket_G}(\mu) = \text{card}_{\llbracket P \rrbracket_G}(\mu)$ , and, for all  $\mu_2 \in \llbracket P'_2 \theta \rrbracket_G$ , either  $\mu \approx \mu_2$ , or  $\llbracket \text{eq}(X, \theta) \rrbracket_G^{\mu \cup \mu_2} \neq \text{true}$ . It remains to show that  $\mu \notin \llbracket P_2 \rrbracket_G$ . For the sake of contradiction, suppose this is not the case. Since  $?x, ?y$ , and  $?z$  are fresh, there are terms  $t_x, t_y$ , and  $t_z$  such that  $\mu \cup \{?x \mapsto t_x, ?y \mapsto t_y, ?z \mapsto t_z\} \in \llbracket P'_2 \rrbracket_G$ , and hence  $\mu \theta \cup \{?x \mapsto t_x, ?y \mapsto t_y, ?z \mapsto t_z\} \in \llbracket P'_2 \theta \rrbracket_G$ . Since  $\mu \sim (\mu \theta \cup \{?x \mapsto t_x, ?y \mapsto t_y, ?z \mapsto t_z\})$  and  $\llbracket \text{eq}(X, \theta) \rrbracket_G^{\mu \cup \mu \theta \cup \{?x \mapsto t_x, ?y \mapsto t_y, ?z \mapsto t_z\}} = \text{true}$ , we obtain a contradiction, which completes the proof of the claim.  $\square$

CLAIM 3.11. *Any Sparql<sub>PD</sub> query  $Q$  that has a subpattern  $\text{Distinct}(P)$  such that  $P$  has in turn a subpattern  $\text{Distinct}(P')$  is equivalent to the query obtained from  $Q$  by replacing  $\text{Distinct}(P')$  with  $P'$ .*

PROOF. First, we can assume that  $Q$  does not have subpatterns of the form  $\text{Distinct}(\text{Distinct}(P))$ , because such subpatterns can clearly be replaced by  $\text{Distinct}(P)$ .

Therefore, it is enough to prove that any pattern

$$\text{Distinct}(\text{Op}(\text{ExpVar}, P_1, \dots, P_i, \dots, P_n))$$

is equivalent to

$$\text{Distinct}(\text{Op}(\text{ExpVar}, P_1, \dots, \text{Distinct}(P_i), \dots, P_n)),$$

for any  $i \in \{1, \dots, n\}$  and any  $\text{Op}(\text{ExpVar}, P_1, \dots, P_i, \dots, P_n)$  among *Join*( $P_1, P_2$ ), *Union*( $P_1, P_2$ ), *Filter*( $E, P_1$ ), *LeftJoin*( $E, P_1, P_2$ ), and *Project*( $X, P_1$ ) such that  $P_1$  and  $P_2$  are patterns in Sparql<sub>PD</sub>. Indeed, by exhaustively applying the equivalence to  $Q$  in the direction introducing extra *Distinct*, we obtain a query  $Q''$  in which each subpattern of  $P$  with some operator different from *Distinct* on top is prefixed by *Distinct*. Moreover,  $Q''$  can be also obtained by exhaustively applying the equivalence to the query  $Q'$  obtained from  $Q$  by replacing  $\text{Distinct}(P')$  with  $P'$ . Hence,  $Q$  is equivalent to  $Q'$  by transitivity of the equivalence relation.

We now prove the equivalence for the case of *LeftJoin*; the remaining cases are similar. We show that pattern  $P^1$  is equivalent to pattern  $P^2$ , for

$$P^1 = \text{Distinct}(\text{LeftJoin}(E, P_1, P_2)) \text{ and}$$

$$P^2 = \text{Distinct}(\text{LeftJoin}(E, \text{Distinct}(P_1), P_2)).$$

By definition, the multisets  $\llbracket \text{LeftJoin}(E, P_1, P_2) \rrbracket_G$  and  $\llbracket \text{LeftJoin}(E, \text{Distinct}(P_1), P_2) \rrbracket_G$  have the same base set for any graph  $G$ . However, this base set is also the base set of  $\llbracket P^1 \rrbracket_G$  and  $\llbracket P^2 \rrbracket_G$ . Therefore,  $\llbracket P^1 \rrbracket_G = \llbracket P^2 \rrbracket_G$ , as required, because these multisets have multiplicity 1 for all mappings by the semantics of *Distinct*.  $\square$

CLAIM 3.15. *For every Sparql pattern  $P$  not mentioning  $d_1, d_2$  and every  $\mu' \in \llbracket P \rrbracket_{G_{1,2}}$ , the mapping  $\mu''$  such that  $\text{dom}(\mu'') = \text{dom}(\mu')$  and, for each  $?u \in \text{dom}(\mu')$ ,*

$$\mu''(?u) = \begin{cases} d_2 & \text{if } \mu'(?u) = d_1, \\ d_1 & \text{if } \mu'(?u) = d_2, \\ \mu'(?u) & \text{otherwise,} \end{cases}$$

*also belongs to  $\llbracket P \rrbracket_{G_{1,2}}$ .*

PROOF. The proof is by induction on the structure of  $P$ .

The base case is when  $P$  is a BGP. Since  $P$  does not mention  $d_1$  or  $d_2$ , for any  $?u$  with  $\mu'(?u) = d_j$ , BGP  $P$  contains a triple pattern of the form  $(s, t, ?u)$  where  $s, t$  are terms such that  $\mu'(s) = c$  and  $\mu'(t) = q$  (assuming that  $\mu'(t') = t'$  for  $t'$  an IRI). Then  $\mu'' \in \llbracket P \rrbracket_{G_{1,2}}$  follows by construction.

In the inductive step,  $P$  must be one of *Join*( $P_1, P_2$ ), *Union*( $P_1, P_2$ ), *Filter*( $E, P_1$ ), and *LeftJoin*( $E, P_1, P_2$ ). As in all four cases the arguments are similar, we consider just the first and the last.

Let  $P = \text{Join}(P_1, P_2)$ . Then, by the definition of *Join*,  $\mu' = \mu'_1 \cup \mu'_2$  such that  $\mu'_i \in \llbracket P_i \rrbracket_{G_{1,2}}$  for both  $i \in \{1, 2\}$ . Also, for any  $?u$  with  $\mu'(?u) = d_j$  there is some  $i$  with  $?u \in \text{dom}(\mu'_i)$ . By the inductive hypothesis, there are  $\mu''_i \in \llbracket P_i \rrbracket_{G_{1,2}}$  for both  $i$  that agree with  $\mu'_i$  on all variables but  $?u$  with  $\mu'_i(?u) = d_j$ , for which  $\mu''_i(?u) = d_{3-j}$ . By construction,  $\mu''_1 \sim \mu''_2$ , and hence  $\mu''_1 \cup \mu''_2 \in \llbracket P \rrbracket_{G_{1,2}}$  as required.

Suppose now  $P = \text{LeftJoin}(E, P_1, P_2)$ . By the definition of *LeftJoin*, there are two cases:

(i)  $\mu' = \mu'_1 \cup \mu'_2$  such that  $\mu'_i \in \llbracket P_i \rrbracket_{G_{1,2}}$ , for  $i \in \{1, 2\}$ , and  $\llbracket E \rrbracket_{G_{1,2}}^{\mu'} = \text{true}$ , or

(ii)  $\mu' \in \llbracket P_1 \rrbracket_{G_{1,2}}$  and for any  $\mu'_2 \in \llbracket P_2 \rrbracket_{G_{1,2}}$  either  $\mu' \sim \mu'_2$  or  $\llbracket E \rrbracket_{G_{1,2}}^{\mu' \cup \mu'_2} \neq \text{true}$ .

Case (i) differs from the one for *Join* only by the presence of  $E$ , which, however, does not make much difference: since  $E$  evaluates to *true* on  $\mu'$  and does not mention any  $d_j$ , it evaluates to *true* on  $\mu''$  as well. In case (ii),  $\mu'' \in \llbracket P_1 \rrbracket_{G_{1,2}}$  by the inductive hypothesis. Therefore, it suffices to show that for every  $\mu''_2 \in \llbracket P_2 \rrbracket_{G_{1,2}}$  either  $\mu'' \sim \mu''_2$  or  $\llbracket E \rrbracket_{G_{1,2}}^{\mu'' \cup \mu''_2} \neq \text{true}$ . Suppose, for contradiction, that this is not the case and there is some  $\mu''_2 \in \llbracket P_2 \rrbracket_{G_{1,2}}$  such that  $\mu'' \sim \mu''_2$  and  $\llbracket E \rrbracket_{G_{1,2}}^{\mu'' \cup \mu''_2} = \text{true}$ . Consider the mapping  $\mu'_2$  that agrees with  $\mu''_2$  on all variables but  $?u$  with  $\mu'_2(?u) = d_j$ , for which  $\mu'_2(?u) = d_{3-j}$ . On the one hand, by the inductive hypothesis,  $\mu'_2 \in \llbracket P_2 \rrbracket_{G_{1,2}}$ . On the other,  $\mu' \sim \mu'_2$  and  $\llbracket E \rrbracket_{G_{1,2}}^{\mu' \cup \mu'_2} = \text{true}$  by construction and the fact that  $E$  does not mention any  $d_j$ , contradicting our assumptions.  $\square$

## B PROOFS FOR SECTION 4

PROPOSITION 4.10. *Let  $P$  be a Sparql $_{PD}^{\exists}$  pattern,  $G$  a graph, and  $\nu$  an environment. Then  $\nu \subseteq \mu$  for every  $\mu \in [P]_G^{\nu}$ .*

PROOF. The proof is by induction on the structure of  $P$ . For the basis of the induction, let  $P$  be a BGP. Then the claim follows from Definition 4.7.

All inductive cases, i.e., the cases when  $P$  is one of *Join*( $P_1, P_2$ ), *Union*( $P_1, P_2$ ), *Filter*( $E, P'$ ), *LeftJoin*( $E, P_1, P_2$ ), *Project*( $X, P'$ ), and *Distinct*( $P'$ ), are straightforward and similar. Let us consider only the most complex case when  $P = \text{Project}(X, P')$ . Let  $\mu$  be a mapping in  $[\text{Project}(X, P')]_G^{\nu}$ . Then, by definition,  $\mu = \mu'|_X \cup \nu|_{\text{dom}(\nu) \setminus X}$  for some  $\mu' \in [P']_G^{\nu|_X}$ . By the inductive hypothesis,  $\nu|_X \subseteq \mu'$  and hence  $\nu|_X \subseteq \mu'|_X$ . Therefore,  $\nu = \nu|_X \cup \nu|_{\text{dom}(\nu) \setminus X} \subseteq \mu'|_X \cup \nu|_{\text{dom}(\nu) \setminus X} = \mu$ .  $\square$

LEMMA 4.12. Let  $G$  be a graph,  $E$  a Sparql $_{PD}^{\exists}$  expression with rectification  $\bar{E}$ , and  $P$  a Sparql $_{PD}^{\exists}$  pattern with rectification  $\bar{P}$ .

1. For every two mappings  $\mu, \nu$  such that the variables introduced by the rectification  $\bar{E}$  of  $E$  are not in  $\text{dom}(\mu) \cup \text{dom}(\nu)$  and  $\llbracket v(\bar{E}) \rrbracket_G^\mu$  is well-defined, if  $\mu \sim \nu$  then  $\llbracket v(\bar{E}) \rrbracket_G^\mu = [E]_G^{\mu \cup \nu}$ .
2. For every two mappings  $\mu, \nu$  such that the variables introduced by the rectification  $\bar{P}$  of  $P$  are not in  $\text{dom}(\nu)$  and  $\llbracket v(\bar{P}) \rrbracket_G$  is well-defined,
  - (a) if  $\mu \in \llbracket v(\bar{P}) \rrbracket_G$  then  $\mu \cup \nu \in [P]_G^\nu$  and  $\text{card}_{\llbracket v(\bar{P}) \rrbracket_G}(\mu) = \text{card}_{[P]_G^\nu}(\mu \cup \nu)$ ;
  - (b) if  $\mu \in [P]_G^\nu$  then  $\mu|_{\text{dom}(\mu) \setminus \text{dom}(\nu)} \in \llbracket v(\bar{P}) \rrbracket_G$  and  $\text{card}_{[P]_G^\nu}(\mu) = \text{card}_{\llbracket v(\bar{P}) \rrbracket_G}(\mu|_{\text{dom}(\mu) \setminus \text{dom}(\nu)})$ .

PROOF. We proceed by simultaneous induction on the structure of  $E$  and  $P$ . We start with  $E$ , that is, with statement 1.

For the base case,  $E$  is either  $?x$ ,  $\ell$ ,  $\text{isIRI}(t)$ ,  $\text{isLiteral}(t)$ ,  $\text{isBlank}(t)$ , or  $\text{bound}(?x)$ . In the last case, the proof is straightforward since, by assumption,  $?x \notin \text{dom}(\nu)$ . The proofs for the other cases are similar, so we consider only the case  $E = ?x$  for a variable  $?x$ . Then, assuming  $\mu \sim \nu$ , we have  $\llbracket v(?x) \rrbracket_G^\mu = \mu(v(?x)) = (\mu \cup \nu)(?x) = \llbracket ?x \rrbracket_G^{\mu \cup \nu}$ .

Consider the inductive step. If  $E$  is  $E_1 \circ E_2$ , for  $\circ \in \{+, -, *, /, \doteq, <, >, \wedge, \vee\}$ , or  $\neg E'$ , then the statement follows directly by the inductive hypothesis. Thus, suppose that  $E = \text{exists}(P)$ . By definition,  $\llbracket v(\text{exists}(\bar{P})) \rrbracket_G^\mu = \llbracket \text{exists}(v(\bar{P})) \rrbracket_G^\mu = \text{true}$  if  $\llbracket \mu(v(\bar{P})) \rrbracket_G = \llbracket (\mu \cup \nu)(\bar{P}) \rrbracket_G \neq \emptyset$  and *false* otherwise. Moreover, since the rectification of  $E$  does not introduce variables in  $\text{dom}(\mu) \cup \text{dom}(\nu) = \text{dom}(\mu \cup \nu)$ , all variables in  $\text{dom}(\mu \cup \nu)$  that occur in  $\bar{P}$  are contained in  $\text{var}(\bar{P})$ . By the inductive hypothesis for statements 2(a) and 2(b),  $\llbracket (\mu \cup \nu)(\bar{P}) \rrbracket_G \neq \emptyset$  if and only if  $[P]_G^{\mu \cup \nu} \neq \emptyset$ , which, in turn, holds if and only if  $[\text{exists}(P)]_G^{\mu \cup \nu} = \text{true}$ . The claim follows.

Let us now prove statement 2. For both 2(a) and 2(b) we concentrate on the first part; the reasoning regarding multiplicities is straightforward.

For the base case, suppose  $P$  is a BGP  $B$ . For statement 2(a), let  $\mu \in \llbracket v(\bar{B}) \rrbracket_G$ . Since  $\bar{B} = B$ ,  $\mu(v(B)) \subseteq G$  and  $\text{dom}(\mu) = \text{var}(B) \setminus \text{dom}(\nu)$  by definition. Consequently,  $\mu \sim \nu$  and  $(\mu \cup \nu)(B) = \mu(v(B))$ . Thus  $(\mu \cup \nu)(B) \subseteq G$  and  $\mu \cup \nu \in [B]_G^\nu$  as required. For statement 2(b), let  $\mu \in [B]_G^\nu$ . Then  $\mu(B) \subseteq G$ ,  $\text{dom}(\mu) = \text{var}(B) \cup \text{dom}(\nu)$ , and  $\mu \sim \nu$ . Hence,  $\text{var}(v(B)) = \text{var}(B) \setminus \text{dom}(\nu) = \text{dom}(\mu) \setminus \text{dom}(\nu)$ . Moreover,  $\mu(v(B)) = \mu(B)$ , so  $\mu(B) \subseteq G$  implies  $\mu|_{\text{dom}(\mu) \setminus \text{dom}(\nu)}(v(B)) \subseteq G$ . Since  $\bar{B} = B$ ,  $\mu|_{\text{dom}(\mu) \setminus \text{dom}(\nu)} \in \llbracket v(\bar{B}) \rrbracket_G$ , as required.

For the inductive step, we consider several cases.

Suppose that  $P = \text{Join}(P_1, P_2)$ . For statement 2(a), consider an arbitrary  $\mu \in \llbracket v(\text{Join}(\bar{P}_1, \bar{P}_2)) \rrbracket_G$ . Then  $\mu = \mu_1 \cup \mu_2$  such that  $\mu_i \in \llbracket v(\bar{P}_i) \rrbracket_G$ . Moreover, the variables introduced by the rectification of each  $P_i$  are not in  $\text{dom}(\nu)$  since this is the case for  $\bar{P}$ . Therefore, by the inductive hypothesis,  $\mu_i \cup \nu \in [P_i]_G^\nu$ . Hence,  $\mu_1 \cup \mu_2 \cup \nu = \mu \cup \nu \in [\text{Join}(P_1, P_2)]_G^\nu$ , as required. The argument for statement 2(b) is analogous.

Suppose  $P = \text{Union}(P_1, P_2)$ . Both statements are shown analogously to case of *Join*.

Suppose  $P = \text{Filter}(E, P')$ . Note that, analogously to the case of *Join*, the variables introduced by the rectifications  $\bar{E}$  and  $\bar{P}'$  are not in  $\text{dom}(\nu)$  since it is the case for  $P$ . For statement 2(a), let  $\mu \in \llbracket v(\text{Filter}(\bar{E}, \bar{P}')) \rrbracket_G$ . By definition,  $\mu \in \llbracket v(\bar{P}') \rrbracket_G$  and  $\llbracket v(\bar{E}) \rrbracket_G^\mu = \text{true}$ . Moreover, we have  $\text{dom}(\mu) \subseteq \text{var}(v(\bar{P}'))$ , and hence, all the variables introduced by the rectification  $\bar{E}$  are not in  $\text{dom}(\mu)$ . By the inductive hypothesis for statement 2(a), we obtain  $\mu \cup \nu \in [P']_G^\nu$ . Thus,  $\mu \sim \nu$ , and, by the inductive hypothesis for statement 1,  $[E]_G^{\mu \cup \nu} = \text{true}$ . Hence  $\mu \cup \nu \in [\text{Filter}(E, P')]_G^\nu$ , as required. For statement 2(b), suppose  $\mu \in [\text{Filter}(E, P')]_G^\nu$ . By definition,  $\mu \in [P']_G^\nu$  and  $[E]_G^\mu = \text{true}$ . Moreover, by rectification, all variables in  $\text{dom}(\mu)$  and  $\bar{E}$  are in  $\text{var}(\bar{E})$ . By the inductive hypothesis for statement 2(b), we then have that  $\mu' = \mu|_{\text{dom}(\mu) \setminus \text{dom}(\nu)} \in \llbracket v(\bar{P}') \rrbracket_G$ . Moreover, by Proposition 4.10,  $\mu = \mu' \cup \nu$ . Hence, by the inductive hypothesis for statement 1,  $\llbracket v(\bar{E}) \rrbracket_G^{\mu'} = \text{true}$ . Therefore,  $\mu' \in \llbracket v(\text{Filter}(\bar{E}, \bar{P}')) \rrbracket_G$ , as required.

Suppose that  $P = \text{LeftJoin}(E, P_1, P_2)$ . Note that, once again, all variables in  $\text{dom}(\nu)$  that occur in  $\bar{E}$ ,  $\bar{P}_1$ , and  $\bar{P}_2$  are contained in  $\text{var}(\bar{E})$ ,  $\text{var}(\bar{P}_1)$ , and  $\text{var}(\bar{P}_2)$ , respectively. For statement 2(a), assume that  $\mu \in \llbracket v(\text{LeftJoin}(\bar{E}, \bar{P}_1, \bar{P}_2)) \rrbracket_G$ . By definition, one of the following holds:

- (i)  $\mu \in \llbracket \text{Filter}(v(\bar{E}), \text{Join}(v(\bar{P}_1), v(\bar{P}_2))) \rrbracket_G$ , or

(ii)  $\mu \in \llbracket v(\bar{P}_1) \rrbracket_G$  and for each  $\mu_2 \in \llbracket v(\bar{P}_2) \rrbracket_G$ , either  $\mu \approx \mu_2$  or  $\llbracket v(\bar{E}) \rrbracket_G^{\mu \cup \mu_2} \neq \text{true}$ .

In case (i), the argument is the same as in the cases for *Join* and *Filter*. In case (ii), by the inductive hypothesis we have  $\mu \cup v \in [P_1]_G^v$ . Suppose, for contradiction, that  $\mu \cup v \notin [P]_G^v$ . By definition, there is some  $\mu'_2 \in [P_2]_G^v$  such that  $\mu \cup v \sim \mu'_2$  and  $[E]_G^{\mu \cup \mu'_2 \cup v} = \text{true}$ . Then,  $\mu \sim \mu_2$  for  $\mu_2 = \mu'_2|_{\text{dom}(\mu'_2) \setminus \text{dom}(v)}$  and, by the inductive hypothesis for statement 2(b),  $\mu_2 \in \llbracket v(\bar{P}_2) \rrbracket_G$ . By rectification, we have that all variables in  $\text{dom}(\mu) \cup \text{dom}(\mu_2)$  and  $\bar{E}$  are contained in  $\text{var}(\bar{E})$ , and hence, by Proposition 4.10 and the inductive hypothesis for statement 1,  $\llbracket v(\bar{E}) \rrbracket_G^{\mu \cup \mu_2} = \text{true}$ , contradicting the requirements in (ii). Thus, our assumption was wrong, and statement 2(a) holds. We proceed to 2(b). Consider  $\mu \in [LeftJoin(E, P_1, P_2)]_G^v$ . One of the following must hold:

(i)  $\mu \in [Filter(E, Join(P_1, P_2))]_G^v$ , or

(ii)  $\mu \in [P_1]_G^v$  and for each  $\mu'_2 \in [P_2]_G^v$  either  $\mu \approx \mu'_2$  or  $[E]_G^{\mu \cup \mu'_2} \neq \text{true}$ .

Case (i) is proved analogously to the case for *Filter* and *Join*. In case (ii), by the inductive hypothesis we have  $\mu_1 \in \llbracket v(\bar{P}_1) \rrbracket_G$  for  $\mu_1 = \mu|_{\text{dom}(\mu) \setminus \text{dom}(v)}$ . Suppose, for contradiction, that  $\mu_1 \notin \llbracket v(\bar{P}) \rrbracket_G$ . Then there is some  $\mu_2 \in \llbracket v(\bar{P}_2) \rrbracket_G$  such that  $\mu_1 \sim \mu_2$  and  $\llbracket v(\bar{E}) \rrbracket_G^{\mu_1 \cup \mu_2} = \text{true}$ . Then, by the inductive hypothesis for statement 2(a), we have  $\mu_2 \cup v \in [P_2]_G^v$ . Moreover,  $\mu \sim \mu_2 \cup v$ , and, by rectification, all variables in  $\text{dom}(\mu) \cup \text{dom}(\mu_2)$  and  $\bar{E}$  are contained in  $\text{var}(\bar{E})$ . Thus, by the inductive hypothesis for statement 1,  $[E]_G^{\mu \cup \mu_2 \cup v} = \text{true}$ , contradicting the requirements in (ii).

Suppose  $P = Project(X, P')$ . Since  $\llbracket v(Project(X, \bar{P}')) \rrbracket_G$  is well-defined, we have that  $\text{dom}(v) \cap X = \emptyset$ . For statement 2(a), let  $\mu \in \llbracket v(Project(X, \bar{P}')) \rrbracket_G$ . Since  $\text{dom}(v) \cap X = \emptyset$ , we have  $v(Project(X, \bar{P}')) = Project(X, v(\bar{P}'))$ , and hence  $\text{dom}(\mu) \subseteq X$ . Hence, there is some  $\mu' \in \llbracket v(\bar{P}') \rrbracket_G$  such that  $\mu = \mu'|_X$ . Furthermore, since  $\text{dom}(v) \cap X = \emptyset$  and, by assumption, all variables in  $\text{dom}(v)$  and  $\bar{P}'$  are contained in  $X$ , we have  $v(\bar{P}') = \bar{P}'$ . Hence, by the inductive hypothesis for statement 2(a),  $\mu' \in [P']_G^0$ . Since  $v|_X = \emptyset$ , we conclude  $\mu' \in [P']_G^{v|_X}$  and  $\mu \cup v = \mu'|_X \cup v = \mu'|_X \cup v|_{\text{dom}(v) \setminus X}$ , and thus  $\mu \cup v \in [Project(X, P')]_G^v$ , as required. For statement 2(b), let  $\mu \in [Project(X, P')]_G^v$ . Then there is some  $\mu' \in [P']_G^{v|_X}$  such that  $\mu = \mu'|_X \cup v|_{\text{dom}(v) \setminus X}$ . Again, since  $\text{dom}(v) \cap X = \emptyset$ , we have  $\mu' \in [P']_G^0$  and  $\mu = \mu'|_X \cup v$ . By the inductive hypothesis, we then have  $\mu' \in \llbracket \bar{P}' \rrbracket_G = \llbracket v(\bar{P}') \rrbracket_G$ . Hence,  $\mu|_{\text{dom}(\mu) \setminus \text{dom}(v)} = \mu'|_X \in \llbracket Project(X, v(\bar{P}')) \rrbracket_G = \llbracket v(Project(X, \bar{P}')) \rrbracket_G$ , as required.

Suppose  $P = Distinct(P')$ . Statements 2(a) and 2(b) are both immediate by the inductive hypothesis since  $\mu \in [Distinct(\bar{P}')]_G^v$  if and only if  $\mu \in [\bar{P}']_G^v$ ,  $\mu \in \llbracket v(Distinct(\bar{P}')) \rrbracket_G$  if and only if  $\mu \in \llbracket v(\bar{P}') \rrbracket_G$ , and  $\text{var}(Distinct(\bar{P}')) = \text{var}(\bar{P}')$  for any  $\mu, v$  and  $G$ .  $\square$

**CLAIM 4.19.** For any Sparql<sub>\*</sub> patterns  $P$  and  $P_1$ , graph  $G$ , and mapping  $\mu_1 \in [P_1]_G^0$ , a mapping  $\mu$  belongs to  $[P]_G^{\mu_1}$  if and only if there is a mapping  $\mu_*$  in  $[Sub_\theta(P, P_1)]_G^0$  for a renaming  $\theta$  of  $\text{var}(P_1)$  such that  $\mu_*|_{\text{var}(P_1) \cup \text{var}(P)} = \mu$  and  $\mu_*|_{\text{var}(P_1)} \theta = \mu_1 \theta$ .

**PROOF.** The proof is by induction on the structure of  $P$ .

For the base case, let  $P$  be a BGP  $B$ .

For the forward direction, consider any mapping  $\mu \in [B]_G^{\mu_1}$ . Let  $\mu_* = \mu \cup \mu_1 \theta$  (the union is well-defined because the domains of  $\mu$  and  $\mu_1 \theta$  are disjoint). On the one hand, by (the alternative) definition,  $\mu = \mu' \cup \mu_1$  for some  $\mu' \in \llbracket B \rrbracket_G$ , so  $\mu_* = \mu' \cup \mu_1 \cup \mu_1 \theta$ . Therefore,  $\mu_* \in [Join(B, Copy_\theta(P_1))]_G^0$ , that is,  $\mu_* \in [Sub_\theta(B, P_1)]_G^0$ . On the other hand, by construction, the restrictions of  $\mu_*$  to  $\text{var}(P_1) \cup \text{var}(B)$  and  $(\text{var}(P_1))\theta$  coincide with  $\mu$  and  $\mu_1 \theta$ , as required.

For the backward direction, consider any mapping  $\mu_*$  in  $[Join(B, Copy_\theta(P_1))]_G^0$  such that  $\mu_*|_{(\text{var}(P_1))\theta} = \mu_1 \theta$ . Let  $\mu = \mu_*|_{\text{var}(P_1) \cup \text{var}(B)}$ . By construction,  $\mu = \mu' \cup \mu_1$  for some  $\mu' \in \llbracket B \rrbracket_G$ , that is,  $\mu \in [B]_G^{\mu_1}$ , as required.

For the inductive step, we have several cases, depending on the top-most operator of  $P$ . Since all of them are similar, we present only the more complex cases, when the operator is *Join*, *Union*, *LeftJoin*, or *Project*.

*Case 1.* Let the claim hold for patterns  $P^1$  and  $P^2$ , and let  $P = Join(P^1, P^2)$ .

For the forward direction, consider any mapping  $\mu$  in  $[Join(P^1, P^2)]_G^{\mu_1}$ , that is,  $\mu = \mu^1 \cup \mu^2$  for some  $\mu^i \in [P^i]_G^{\mu_1}$ ,  $i \in \{1, 2\}$ . By the inductive hypothesis, there are mappings  $\mu_*^i \in [Sub_{\theta^i}(P^i, P_1)]_G^0$  such that

$\mu_*^i|_{\text{var}(P_1) \cup \text{var}(P^i)} = \mu^i$  and  $\mu_*^i|_{(\text{var}(P_1))\theta^i} = \mu_1\theta^i$ . Let  $\mu_* = \mu_*^1 \cup \mu_*^2 \cup \mu_1\theta$ . On the one hand,

$$\mu_* \in [EqJoin_{\theta, \theta^1}^{\text{var}(P_1)}(EqJoin_{\theta^1, \theta^2}^{\text{var}(P_1)}(\text{Sub}_{\theta^1}(P^1, P_1), \text{Sub}_{\theta^2}(P^2, P_1)), P_1\theta)]_G^0,$$

that is,  $\mu_* \in [\text{Sub}_\theta(\text{Join}(P^1, P^2), P_1)]_G^0$  by construction. On the other, again by construction,  $\mu_*|_{\text{var}(P_1) \cup \text{var}(P)} = \mu$  and  $\mu_*|_{(\text{var}(P_1))\theta} = \mu_1\theta$ , as required.

For the backward direction, consider any  $\mu_*$  in  $[\text{Sub}_\theta(\text{Join}(P^1, P^2), P_1)]_G^0$  such that  $\mu_*|_{(\text{var}(P_1))\theta} = \mu_1\theta$ . Let  $\mu'_*$  be the restriction of  $\mu_*$  to all variables except  $(\text{var}(P_1))\theta$ . By the definition of  $\text{Sub}_\theta(\text{Join}(P^1, P^2), P_1)$ ,  $\mu'_* = \mu_*^1 \cup \mu_*^2$ , where each  $\mu_*^i \in [\text{Sub}_{\theta^i}(P^i, P_1)]_G^0$  is such that  $\mu_*^i|_{(\text{var}(P_1))\theta^i} = \mu_1\theta^i$ . Hence, by the inductive hypothesis,  $\mu^i = \mu_*^i|_{\text{var}(P_1) \cup \text{var}(P^i)}$  belongs to  $[P^i]_G^{\mu_1}$  for both  $i$ . Since  $\mu^1$  and  $\mu^2$  are both restrictions of  $\mu_*$ , their union  $\mu = \mu^1 \cup \mu^2$  is well-defined and belongs to  $[\text{Join}(P^1, P^2)]_G^{\mu_1}$ , as required.

*Case 2.* Let the claim hold for patterns  $P^1$  and  $P^2$ , and let  $P = \text{Union}(P^1, P^2)$ .

For the forward direction, consider any mapping  $\mu$  in  $[\text{Union}(P^1, P^2)]_G^{\mu_1}$ . Then  $\mu \in [P^i]_G^{\mu_1}$  for  $i = 1$  or  $i = 2$ . By the inductive hypothesis, there exists  $\mu_* \in [\text{Sub}_\theta(P^i, P_1)]_G^0$  such that  $\mu_*|_{\text{var}(P_1) \cup \text{var}(P^i)} = \mu$  and  $\mu_*|_{(\text{var}(P_1))\theta} = \mu_1\theta$ . Therefore,  $\mu_* \in [\text{Sub}_\theta(\text{Union}(P^1, P^2), P_1)]_G^0$  by definition. Finally, we have  $\mu_*|_{\text{var}(P_1) \cup \text{var}(P)} = \mu_*|_{\text{var}(P_1) \cup \text{var}(P^i)} = \mu$ , as required.

For the backward direction, consider any mapping  $\mu_*$  in  $[\text{Sub}_\theta(\text{Union}(P^1, P^2), P_1)]_G^0$  such that  $\mu_*|_{(\text{var}(P_1))\theta} = \mu_1\theta$ . By definition,  $\mu_* \in [\text{Sub}_\theta(P^i, P_1)]_G^0$  for  $i = 1$  or  $i = 2$ . Hence, by the inductive hypothesis,  $\mu = \mu_*|_{\text{var}(P_1) \cup \text{var}(P^i)}$  belongs to  $[P^i]_G^{\mu_1}$ , and hence  $\mu \in [\text{Union}(P^1, P^2)]_G^{\mu_1}$ . Moreover, we have  $\mu = \mu_*|_{\text{var}(P_1) \cup \text{var}(P)}$ , as required.

*Case 3.* Let the claim hold for patterns  $P^1$  and  $P^2$ , and let  $P = \text{LeftJoin}(E, P^1, P^2)$  for an expression  $E$ .

For the forward direction, consider any mapping  $\mu$  in  $[\text{LeftJoin}(E, P^1, P^2)]_G^{\mu_1}$ , that is, any  $\mu$  such that either

(i)  $\mu = \mu^1 \cup \mu^2$  for some  $\mu^i \in [P^i]_G^{\mu_1}$ ,  $i \in \{1, 2\}$ , and  $[E]_G^\mu = \text{true}$ , or

(ii)  $\mu = \mu^1$  for  $\mu^1 \in [P^1]_G^{\mu_1}$  such that for any  $\mu^2 \in [P^2]_G^{\mu_1}$  either  $\mu^1 \approx \mu^2$  or  $[E]_G^{\mu^1 \cup \mu^2} \neq \text{true}$ .

In case (i), the existence of  $\mu_*$  in  $[\text{Sub}_\theta(P, P_1)]_G^0$  with the required properties can be proved in the same way as for Case 1, when  $P = \text{Join}(P^1, P^2)$ . Therefore, consider case (ii). By the inductive hypothesis, there exists  $\mu_*^1 \in [\text{Sub}_{\theta^1}(P^1, P_1)]_G^0$  such that  $\mu_*^1|_{\text{var}(P_1) \cup \text{var}(P^1)} = \mu^1$  and  $\mu_*^1|_{(\text{var}(P_1))\theta^1} = \mu_1\theta^1$ . Mapping  $\mu_*^1$  is not extendable to any  $\mu_*^2 \in [\text{Sub}_{\theta^2}(P^2, P_1)]_G^0$  under  $E \wedge \text{eq}(Y\theta^1, (\theta^1)^{-1} \circ \theta^2)$ , that is, we have

$$\mu_*^1 \in [EqLeftJoin_{\theta^1, \theta^2}^{\text{var}(P_1)}(E, \text{Sub}_{\theta^1}(P^1, P_1), \text{Sub}_{\theta^2}(P^2, P_1))]_G^0;$$

indeed, otherwise the restriction  $\mu^2 = \mu_*^2|_{\text{var}(P_1) \cup \text{var}(P^2)}$ , which is in  $[P^2]_G^{\mu_1}$  by the inductive hypothesis, would be such that  $[E]_G^{\mu^1 \cup \mu^2} = \text{true}$ , and hence  $\mu$  would not satisfy the conditions of case (ii). Thus,  $\mu_* = \mu^1 \cup \mu_1\theta$  belongs to

$$[EqJoin_{\theta, \theta^1}^{\text{var}(P_1)}(EqLeftJoin_{\theta^1, \theta^2}^{\text{var}(P_1)}(E, \text{Sub}_{\theta^1}(P^1, P_1), \text{Sub}_{\theta^2}(P^2, P_1)), P_1\theta)]_G^0,$$

that is,  $\mu_* \in [\text{Sub}_\theta(\text{LeftJoin}(E, P^1, P^2), P_1)]_G^0$ . Finally, by construction,  $\mu_*|_{\text{var}(P_1) \cup \text{var}(P)} = \mu_*|_{\text{var}(P_1) \cup \text{var}(P^1)} = \mu$ , and  $\mu_*|_{(\text{var}(P_1))\theta} = \mu_1\theta$ , as required.

For the backward direction, consider any  $\mu_*$  in  $[\text{Sub}_\theta(\text{LeftJoin}(E, P^1, P^2), P_1)]_G^0$  such that  $\mu_*|_{(\text{var}(P_1))\theta} = \mu_1\theta$ . Let  $\mu'_*$  be the restriction of  $\mu_*$  to all variables except  $(\text{var}(P_1))\theta$ . By definition of  $\text{Sub}_\theta(\text{LeftJoin}(E, P^1, P^2), P_1)$  there are two cases:

(i)  $\mu'_* = \mu_*^1 \cup \mu_*^2$ , for  $\mu_*^i \in [\text{Sub}_{\theta^i}(P^i, P_1)]_G^0$ , where

$$[E \wedge \text{eq}((\text{var}(P_1))\theta^1, (\theta^1)^{-1} \circ \theta^2)]_G^{\mu_*^1 \cup \mu_*^2} = \text{true}$$

and  $\mu_*^i|_{(\text{var}(P_1))\theta^i} = \mu_1\theta^i$  for both  $i$ , or

(ii)  $\mu'_* = \mu_*^1$  for some  $\mu_*^1 \in [\text{Sub}_{\theta^1}(P^1, P_1)]_G^0$  such that, for each  $\mu_*^2 \in [\text{Sub}_{\theta^2}(P^2, P_1)]_G^0$ , either  $\mu_*^1 \approx \mu_*^2$  or

$$[E \wedge \text{eq}((\text{var}(P_1))\theta^1, (\theta^1)^{-1} \circ \theta^2)]_G^{\mu_*^1 \cup \mu_*^2} \neq \text{true},$$

and, moreover,  $\mu_*^1 \upharpoonright_{(\text{var}(P_1))\theta^1} = \mu_1 \theta^1$ .

Case (i) is similar to Case 1, when  $P = \text{Join}(P^1, P^2)$ . By the inductive hypothesis, each  $\mu^i = \mu_*^i \upharpoonright_{\text{var}(P_1) \cup \text{var}(P^i)}$  belongs to  $[P^i]_G^{\mu_1}$ . Since  $\mu^1$  and  $\mu^2$  are both restrictions of  $\mu$ , their union  $\mu = \mu^1 \cup \mu^2$  is well-defined. Since  $E$  does not mention variables in  $\text{dom}(\mu_*^i) \setminus \text{dom}(\mu^i)$  for any  $i \in \{1, 2\}$ , we have  $\mu \in [\text{LeftJoin}(E, P^1, P^2)]_G^{\mu_1}$ , as required. Consider case (ii). By the inductive hypothesis, the restriction  $\mu^1 = \mu_*^1 \upharpoonright_{\text{var}(P_1) \cup \text{var}(P^1)}$  belongs to  $[P^1]_G^{\mu_1}$ . We need to show that  $\mu^1 \in [\text{LeftJoin}(E, P^1, P^2)]_G^{\mu_1}$ . Suppose, for contradiction, this is not the case. Then there exists  $\mu^2 \in [P^2]_G^{\mu_1}$  such that  $\mu^1 \sim \mu^2$  and  $[E]_G^{\mu^1 \cup \mu^2} = \text{true}$ . Hence, by the inductive hypothesis, there is  $\mu_*^2 \in [\text{Sub}_{\theta^2}(P^2, P_1)]_G^{\theta^2}$  such that  $\mu_*^2 \upharpoonright_{\text{var}(P_1) \cup \text{var}(P^2)} = \mu^2$  and  $\mu_*^2 \upharpoonright_{(\text{var}(P_1))\theta^2} = \mu_1 \theta^2$ . Since  $\mu^1 \sim \mu^2$  and  $\text{dom}(\mu_*^1) \cap \text{dom}(\mu_*^2) = \text{dom}(\mu^1) \cap \text{dom}(\mu^2)$ , we have  $\mu_*^1 \sim \mu_*^2$ . Moreover, since  $\mu_*^1 \cup \mu_*^2$  and  $\mu^1 \cup \mu^2$  coincide on the variables in  $E$ , we have  $[E]_G^{\mu_*^1 \cup \mu_*^2} = \text{true}$ . Finally,  $\mu_*^1 \cup \mu_*^2$  satisfies  $\text{eq}((\text{var}(P_1))\theta^1, (\theta^1)^{-1} \circ \theta^2)$  because, by construction,  $\mu_*^1 \cup \mu_*^2 \upharpoonright_{(\text{var}(P_1))\theta^i} = \mu_1 \theta^i$  for each  $i$ . Hence,

$$[E \wedge \text{eq}((\text{var}(P_1))\theta^1, (\theta^1)^{-1} \circ \theta^2)]_G^{\mu_*^1 \cup \mu_*^2} = \text{true},$$

which contradicts our assumptions about  $\mu_*^i$  in case (ii). Therefore,  $\mu^1 \in [\text{LeftJoin}(E, P^1, P^2)]_G^{\mu_1}$ , as required.

Case 4. Let the claim hold for a pattern  $P'' = P' \theta_{\text{var}(P') \setminus X}$ , for variables  $X$  and pattern  $P'$ , and let  $P = \text{Project}(X, P')$ .

For the forward direction, consider any mapping  $\mu \in [\text{Project}(X, P')]_G^{\mu_1}$ . Then, by the semantics of *Project* and Proposition 4.10, there is some  $\mu' \in [P']_G^{\mu_1 \upharpoonright_X}$  such that  $\mu = \mu' \upharpoonright_X \cup \mu_1$ . Then  $\mu' \theta_{\text{var}(P') \setminus X} \in [P'']_G^{\mu_1 \upharpoonright_X}$ ; furthermore, by construction and Proposition 4.10, the mapping  $\mu'' = \mu' \theta_{\text{var}(P') \setminus X} \cup \mu_1$  is in  $[P'']_G^{\mu_1}$ . Hence, by the inductive hypothesis, there is some  $\mu_*'' \in [\text{Sub}_{\theta}(P'', P_1)]_G^{\theta}$  such that  $\mu_*'' \upharpoonright_{\text{var}(P_1) \cup \text{var}(P'')} = \mu''$  and  $\mu_*'' \upharpoonright_{(\text{var}(P_1))\theta} = \mu_1 \theta$ . Thus, the mapping  $\mu_* = \mu_*'' \upharpoonright_{X \cup \text{var}(P_1) \cup (\text{var}(P_1))\theta}$  is in  $[\text{Sub}_{\theta}(\text{Project}(X, P'), P_1)]_G^{\theta}$ . Furthermore, since  $\text{var}(P) = X$ , we have

$$\mu_* \upharpoonright_{\text{var}(P_1) \cup \text{var}(P)} = \mu_* \upharpoonright_{\text{var}(P_1) \cup X} = \mu_*'' \upharpoonright_{\text{var}(P_1) \cup X} = \mu'' \upharpoonright_{\text{var}(P_1) \cup X} = (\mu' \theta_{\text{var}(P') \setminus X} \cup \mu_1) \upharpoonright_{\text{var}(P_1) \cup X}.$$

Since  $\mu' \theta_{\text{var}(P') \setminus X}$  coincides on  $X$  with  $\mu'$  and is undefined on all other variables in  $\text{var}(P_1)$ , we obtain that  $\mu_* \upharpoonright_{\text{var}(P_1) \cup \text{var}(P)} = (\mu' \upharpoonright_X \cup \mu_1) \upharpoonright_{\text{var}(P_1) \cup X} = \mu' \upharpoonright_X \cup \mu_1 = \mu$ . Finally, we have  $\mu_* \upharpoonright_{(\text{var}(P_1))\theta} = \mu_*'' \upharpoonright_{(\text{var}(P_1))\theta} = \mu_1 \theta$ .

For the backward direction, consider any mapping  $\mu_* \in [\text{Sub}_{\theta}(\text{Project}(X, P'), P_1)]_G^{\theta}$  such that  $\mu_* \upharpoonright_{(\text{var}(P_1))\theta} = \mu_1 \theta$ . Then  $\mu_* = \mu_*'' \upharpoonright_{X \cup \text{var}(P_1) \cup (\text{var}(P_1))\theta}$  for some mapping  $\mu_*'' \in [\text{Sub}_{\theta}(P'', P_1)]_G^{\theta}$ . Since, by construction,  $\mu_*'' \upharpoonright_{(\text{var}(P_1))\theta} = \mu_1 \theta$ , by the inductive hypothesis, the mapping  $\mu'' = \mu_*'' \upharpoonright_{\text{var}(P_1) \cup \text{var}(P'')}$  belongs to  $[P'']_G^{\mu_1}$ . Since all variables shared by  $P''$  and  $P_1$  are contained in  $X$  and  $\text{dom}(\mu_1) \subseteq \text{var}(P_1)$ , there is some  $\mu' \in [P']_G^{\mu_1 \upharpoonright_X}$  such that  $\mu' \upharpoonright_X = \mu'' \upharpoonright_X$ . Let  $\mu = \mu' \upharpoonright_X \cup \mu_1 \upharpoonright_{\text{dom}(\mu_1) \setminus X}$ . By definition,  $\mu \in [\text{Project}(X, P')]_G^{\mu_1}$ . Moreover, by definition and Proposition 4.10,

$$\mu_* \upharpoonright_{\text{var}(P_1) \cup X} = \mu_*'' \upharpoonright_{\text{var}(P_1) \cup X} = \mu'' \upharpoonright_{\text{var}(P_1) \cup X} = \mu' \upharpoonright_X \cup \mu_1 = \mu,$$

as required. This concludes the proof of the original version of the claim.

Next, we extend the proof of the claim to languages with *Extend*, which we need for Theorem 5.6. To this end, we need to consider one more case for the inductive step.

Case 5. Let the claim hold for a pattern  $P'$ , and let  $P = \text{Extend}(?x, E, P')$  for an expression  $E$ .

For the forward direction, consider any mapping  $\mu$  in  $[\text{Extend}(?x, E, P')]_G^{\mu_1}$ , that is, a mapping such that either

- (i)  $\mu = \mu' \cup \{?x \mapsto [E]_G^{\mu'}\}$  for some  $\mu' \in [P']_G^{\mu_1}$  with  $[E]_G^{\mu'} \neq \text{error}$ , or
- (ii)  $\mu \in [P']_G^{\mu_1}$  such that  $[E]_G^{\mu} = \text{error}$ .

Consider case (i). By the inductive hypothesis, there exists a mapping  $\mu_*'$  in  $[\text{Sub}_{\theta}(P', P_1)]_G^{\theta}$  such that  $\mu_*' \upharpoonright_{\text{var}(P_1) \cup \text{var}(P')} = \mu'$  and  $\mu_*' \upharpoonright_{(\text{var}(P_1))\theta} = \mu_1 \theta$ . The mapping  $\mu_* = \mu_*' \cup \{?x \mapsto [E]_G^{\mu_*'}\}$  satisfies all the required conditions. Case (ii) is similar.

For the backward direction, let  $\mu_*$  be a mapping in  $[\text{Sub}_\theta(\text{Extend}(\?x, E, P'), P_1)]_G^0$  such that  $\mu_* \upharpoonright_{\text{var}(P_1)} \theta = \mu_1 \theta$ . By definition, either

- (i)  $\mu_* = \mu'_* \cup \{\?x \mapsto [E]_G^{\mu'_*}\}$  for some  $\mu'_* \in [\text{Sub}_\theta(P', P_1)]_G^0$  with  $[E]_G^{\mu'_*} \neq \text{error}$ , or
- (ii)  $\mu_* \in [\text{Sub}_\theta(P', P_1)]_G^0$  such that  $[E]_G^{\mu_*} = \text{error}$ .

Consider case (i). By the inductive hypothesis, the mapping  $\mu' = \mu'_* \upharpoonright_{\text{var}(P_1) \cup \text{var}(P')}$  belongs to  $[P']_G^{\mu_1}$ . Then, the mapping  $\mu = \mu' \cup \{\?x \mapsto [E]_G^{\mu'}\}$  satisfies all the required conditions. Case (ii) is similar.  $\square$

### C PROOFS FOR SECTIONS 5 TO 9

**PROPOSITION 5.5.** *For any  $\text{Sparql}_{PDE}^{\exists}$  pattern  $P$  and graph  $G$  such that  $\llbracket P \rrbracket_G$  is well-defined it holds that  $\llbracket P \rrbracket_G = [P]_G^0$ .*

**PROOF.** The proof is based on the fact that Lemma 4.12 is still true for  $\text{Sparql}_{PDE}^{\exists}$ . To extend the lemma, we consider another case in the inductive step for statement 2.

Suppose  $P = \text{Extend}(\?x, E, P')$ . Since  $\llbracket v(\text{Extend}(\?x, \bar{E}, \bar{P}')) \rrbracket_G$  is well-defined, we have that  $\?x \notin \text{dom}(v)$ . For statement 2(a), let  $\mu \in \llbracket v(\text{Extend}(\?x, \bar{E}, \bar{P}')) \rrbracket_G$ . Since  $\?x \notin \text{dom}(v)$ , we have  $v(\text{Extend}(\?x, \bar{E}, \bar{P}')) = \text{Extend}(\?x, v(\bar{E}), v(\bar{P}'))$ . Therefore, by the normative semantics of *Extend*, one of the following holds:

- (i)  $\mu = \mu' \cup \{\?x \mapsto \llbracket v(\bar{E}) \rrbracket_G^{\mu'}\}$  for  $\mu' \in \llbracket v(\bar{P}') \rrbracket_G$  and  $\llbracket v(\bar{E}) \rrbracket_G^{\mu'} \neq \text{error}$ , or
- (ii)  $\mu \in \llbracket v(\bar{P}') \rrbracket_G$  and  $\llbracket v(\bar{E}) \rrbracket_G^{\mu'} = \text{error}$ .

The arguments for these cases are similar, therefore we consider only case (i). The variables introduced by the rectification of  $P$  are not in  $\text{dom}(v)$ , so this is true for  $P'$  and  $E$  as well. Hence, by the inductive hypothesis for statement 2(a),  $\mu' \cup v \in [P']_G^v$ . In particular,  $\mu' \sim v$ . Since  $\text{dom}(\mu') \subseteq \text{var}(v(\bar{P}')) \subseteq \text{var}(\bar{P}')$ , the variables introduced by the rectification of  $E$  are not in  $\text{dom}(\mu')$ . By the inductive hypothesis for statement 1, it then holds that  $\llbracket v(\bar{E}) \rrbracket_G^{\mu'} = [E]_G^{\mu' \cup v}$ . Since  $\?x \notin \text{dom}(v)$  and  $\mu' \sim \?x \mapsto \llbracket v(\bar{E}) \rrbracket_G^{\mu'}$ , we have that  $\mu' \cup v \sim \{\?x \mapsto [E]_G^{\mu' \cup v}\}$ . Therefore, by definition,  $\mu \cup v \in [\text{Extend}(\?x, E, P')]_G^v$ , as required.

For statement 2(b), let  $\mu \in [\text{Extend}(\?x, E, P')]_G^v$ . Again, we have two cases:

- (i)  $\mu = \mu' \cup \{\?x \mapsto [E]_G^{\mu'}\}$  for  $\mu' \in [P']_G^v$  and  $[E]_G^{\mu'} \neq \text{error}$ , or
- (ii)  $\mu \in [P']_G^v$  and  $[E]_G^{\mu'} = \text{error}$ .

We consider only case (i), because case (ii) is similar. As noted before, the variables introduced by the rectification of  $P'$  are not in  $\text{dom}(v)$ , so, by the inductive hypothesis for statement 2(b), it holds that  $\mu' \upharpoonright_{\text{dom}(\mu') \setminus \text{dom}(v)} \in \llbracket v(\bar{P}') \rrbracket_G$ . Moreover, since  $\text{dom}(\mu') \subseteq \text{var}(P')$ , the variables introduced by the rectification of  $P'$  are not in  $\text{dom}(\mu')$ . Therefore, by the inductive hypothesis for statement 1 and Proposition 4.10, it holds that  $\llbracket v(\bar{E}) \rrbracket_G^{\mu' \upharpoonright_{\text{dom}(\mu') \setminus \text{dom}(v)}} = [E]_G^{\mu'}$ . Since  $\?x \notin \text{dom}(v)$ , by definition of the normative semantics of *Extend* we have that  $\mu \upharpoonright_{\text{dom}(\mu) \setminus \text{dom}(v)} \in \llbracket v(\text{Extend}(\?x, \bar{E}, \bar{P}')) \rrbracket_G$ , as required.  $\square$

**PROPOSITION 6.9.** *For any  $\text{Sparql}_{PDEA}^{\exists}$  pattern  $P$  and graph  $G$  such that  $\llbracket P \rrbracket_G$  is well-defined it holds that  $\llbracket P \rrbracket_G = [P]_G^0$ .*

**PROOF.** Similarly to the case of *Extend*, the proof is based on the fact that Lemma 4.12 holds for  $\text{Sparql}_{PDEA}^{\exists}$ . To extend the lemma, we consider another case in the inductive step for statement 2.

Suppose  $P = \text{AggregateJoin}_x(\mathbf{A})$  with  $\mathbf{x} = \langle \?x_1, \dots, \?x_n \rangle$  a non-empty list of variables and  $\mathbf{A} = \langle A_1, \dots, A_n \rangle$  a list of aggregates of the same size. Since  $\llbracket v(\bar{P}) \rrbracket_G$  is well-defined, we have  $\{\?x_1, \dots, \?x_n\} \cap \text{dom}(v) = \emptyset$ . Therefore, by the definition of rectification,  $v(\bar{P}) = \bar{P}$ . The claim is then straightforward. Indeed, all top-level subpatterns and subexpressions in  $\mathbf{A}$  are evaluated in the empty environment, meaning that  $v$  is irrelevant for them. Therefore, the normative and environment semantics of these subpatterns and subexpressions coincide by the inductive hypothesis. The claim follows, because the manipulation of groups and aggregates is the same in both cases, and the final step of evaluating *AggregateJoin* differs precisely as required.  $\square$

**PROPOSITION 7.8.** *For any extension  $\text{Sparql}_\star$  of  $\text{Sparql}_p$ , language  $\text{Sparql}_{\star spp}$  has the same expressive power as  $\text{Sparql}_{\star pp}$ .*

PROOF. Recall that we may consider only queries whose BGPs are either empty or consist of single triple patterns.

On the one hand, by definition, triple patterns  $Path(t_1, PP, t_2)$  with  $PP$  of the form  $Link(u, Inv(PP'))$ ,  $Seq(PP_1, PP_2)$ , or  $Alt(PP_1, PP_2)$  are  $\nu$ -equivalent, for any  $\nu$ , to patterns whose triple patterns involve only subexpressions of  $PP$  that have  $NPS$ ,  $ZeroOrOnePath$ ,  $ZeroOrMorePath$ , or  $OneOrMorePath$  as their top-most operator. On the other hand, again by definition, for any graph  $G$ , environment  $\nu$ , and path expression  $PP$  with any of these four operators on top, we have

$$[Path(t_1, PP, t_2)]_G^\nu = [sPath(t_1, sPP, t_2)]_G^\nu,$$

where  $sPP$  is obtained from  $PP$  by replacing each path expression construct by its set path counterpart. Therefore, given a  $Sparql_\star$  query, we can obtain an equivalent  $Sparql_\star$  query by first iteratively decomposing all triple patterns with  $Link$ ,  $Inv$ ,  $Seq$  or  $Alt$  as the top-most operator to patterns built of triple patterns over  $NPS$ ,  $ZeroOrOnePath$ ,  $ZeroOrMorePath$ , and  $OneOrMorePath$ , and then replacing all path constructs by their set counterparts. The translation in the other direction is the same, except that the first step is not needed.  $\square$

PROPOSITION 9.3. *Let  $Sparql_\star$  be any language considered in this paper. The extension of  $Sparql_\star$  allowing for exists and Minus has the same expressive power as  $Sparql_\star$ .*

PROOF. Let us first consider the case when  $Sparql_\star$  allows for *Extend*. Then, for every environment  $\nu$ ,  $Minus(P_1, P_2)$  is  $\nu$ -equivalent to  $Diff(E, P_1, Join(P'_2, P'_3))$ , where  $P'_2$ ,  $P'_3$ , and  $E$  are defined as follows, for  $\theta_1$  and  $\theta_2$  renamings of  $\text{var}(P_1) \cup \text{var}(P_2) = \{?x_1, \dots, ?x_n\}$  to disjoint sets of variables:

$$\begin{aligned} P'_2 &= \text{Extend}(?x_1\theta_1, ?x_1, \dots, \text{Extend}(?x_n\theta_1, ?x_n, P_2) \dots), \\ P'_3 &= \text{Extend}(?x_1\theta_2, ?x_1, \dots, \text{Extend}(?x_n\theta_2, ?x_n, \emptyset) \dots), \\ E &= \bigvee_{?x \in \text{dom}(\theta_1)} (\text{bound}(?x) \wedge \text{bound}(?x\theta_1) \wedge \neg \text{bound}(?x\theta_2)). \end{aligned}$$

Let now  $Sparql_\star$  allow for *Aggregate*. Then we can rewrite *Minus* using *Extend* as in the previous case, and then apply Theorem 6.11 to rewrite *Extend* by means of aggregation.

Finally, consider the case when  $Sparql_\star$  does not allow for value invention. It is enough to generalise Theorem 4.20 to the extension of  $Sparql_\star$  with *Minus*. Lemmas 4.16 and 4.17 hold for the language, so we only need to generalise Lemma 4.18. We will use the following abbreviations. First, let

$$\begin{aligned} All(?y) &= Union((?y, ?y^1, ?y^2), Union((?y^3, ?y, ?y^4), Union((?y^5, ?y^6, ?y), \{\}))), \\ All(\{y_1, \dots, y_n\}) &= Join(All(?y_1), \dots, Join(All(?y_{n-1}), All(?y_n)) \dots), \end{aligned}$$

where  $?y^i$  are fresh variables, different for each  $?y$ . The evaluation of  $All(\{y_1, \dots, y_n\})$  (in the empty environment) gives all possible mappings over variables  $y_1, \dots, y_n$  bound to all possible terms in the input graph (these mappings also bind some of the fresh variables  $?y^i$ , but these bindings are harmless for the rest of the construction). Note that it is crucial here that  $Sparql_\star$  does not allow for value invention, so  $All(\{y_1, \dots, y_n\})$  indeed evaluates to all mappings that are possible in the evaluation of any pattern. Second, for a pattern  $P$  and two renamings  $\theta$  and  $\theta'$  such that  $\text{dom}(\theta) \subseteq \text{dom}(\theta')$ , let

$$TrunkCopy_{\theta, \theta'}(P) = Filter(E_{\theta, \theta'}, Join(P, All(Y))),$$

where  $Y$  is the range of  $\theta'$  and

$$\begin{aligned} E_{\theta, \theta'} &= \bigwedge_{?x \in \text{dom}(\theta)} \left( (\text{bound}(?x\theta) \rightarrow \neg \text{bound}(?x\theta')) \wedge \right. \\ &\quad \left. (\neg \text{bound}(?x\theta) \rightarrow ((\text{bound}(?x) \leftrightarrow \text{bound}(?x\theta')) \wedge (\text{bound}(?x) \rightarrow ?x = ?x\theta'))) \right) \wedge \\ &\quad \bigwedge_{?x \in \text{dom}(\theta') \setminus \text{dom}(\theta)} \left( (\text{bound}(?x) \leftrightarrow \text{bound}(?x\theta')) \wedge (\text{bound}(?x) \rightarrow ?x = ?x\theta') \right). \end{aligned}$$

Intuitively,  $TrunkCopy_{\theta, \theta'}(P)$  extends all the mappings in the evaluation of  $P$  to the copies  $?x\theta'$  of all variables  $?x$  in the domain of  $\theta'$  such that  $?x\theta$  is not bound (or  $?x$  is not in the domain of  $\theta$ ).

Having these abbreviations at hand, let, for a renaming  $\theta_2$  of  $\text{var}(P_1)$ , a renaming  $\theta'$  of  $\text{var}(P^1) \cup \text{var}(P_1)$ , and a renaming  $\theta'_2$  of  $\text{var}(P^2) \cup \text{var}(P_1)$ , all to fresh variables,

$$\text{Sub}_\theta(\text{Minus}(P^1, P^2), P_1) =$$

$$\text{Diff}(E_{\theta, \theta', \theta_2, \theta'_2}, \text{TrunkCopy}_{\theta, \theta'}(\text{Sub}_\theta(P^1, P_1)), \text{TrunkCopy}_{\theta_2, \theta'_2}(\text{Sub}_{\theta_2}(P^1, P_1))),$$

for

$$E_{\theta, \theta', \theta_2, \theta'_2} = \text{eq}((\text{var}(P_1))\theta, \theta^{-1} \circ \theta_2) \wedge$$

$$\bigwedge_{?x \in \text{dom}(\theta') \cup \text{dom}(\theta'_2)} (\text{bound}(?x\theta') \wedge \text{bound}(?x\theta'_2) \rightarrow ?x\theta' = ?x\theta'_2) \wedge$$

$$\bigvee_{?x \in \text{dom}(\theta') \cup \text{dom}(\theta'_2)} (\text{bound}(?x\theta') \wedge \text{bound}(?x\theta'_2)).$$

Intuitively, this translation is similar to the translation of *Minus* in the empty environment: the difference is that the possible extension is checked not for mappings in the evaluations of  $P^1$  and  $P^2$ , but for their copies with all the bound variables in the environment truncated off; of course, the copies of the environment, attached to the mappings by renamings  $\theta$  and  $\theta_2$ , should coincide as well.

Claim 4.19 then generalises: we just need to consider yet another case for the inductive step.  $\square$