



Optimizing SPARQL Queries with SHACL

Ratan Bahadur Thapa^(✉) and Martin Giese^{ID}

Department of Informatics, University of Oslo, Oslo, Norway
{ratanbt,martingi}@ifi.uio.no

Abstract. We propose a set of optimizations that can be applied to a given SPARQL query, and that guarantee that the optimized query has the same answers under bag semantics as the original query, provided that the queried RDF graph validates certain SHACL constraints. Our optimizations exploit the relationship between graph patterns in the SPARQL queries and the SHACL constraints that describe those patterns in the RDF graph. We prove the correctness of these optimizations and show how they can be propagated to larger queries while preserving answers. Further, we prove the confluence of rewritings that employ these optimizations, guaranteeing convergence to the same optimized query regardless of the rewriting order.

1 Introduction

The Resource Description Framework (RDF) [37] has seen increasing adoption in recent years, prompting the W3C standardization of the SPARQL [24,45] query language for RDF. Since the W3C standardization in 2008, SPARQL has been recognised as a key technology for the Semantic Web and the current version SPARQL 1.1 [24] is well-adopted in both academia and industries, e.g., query engines for AllegroGraph [22], Apache Jena [16], Sesame [11] and OpenLink Virtuoso [21].

The RDF model, as an abstract knowledge representation, doesn't explicitly distinguish between data and metadata, such as schema and constraints. To address this, the W3C has recommended the SHACL [32] constraint language for RDF. SHACL validation relies on “shapes,” which define a set of constraints and indicate which nodes in an RDF graph should be validated against these constraints. Since its recommendation in 2017, SHACL has been widely used to verify compliance of RDF datasets w.r.t. certain policies, such as GDPR requirements [39], and to generate guarantees for data transformation [53,54], facilitating data repairs [4] and others [40]. Additionally, a set of SHACL constraints can function as a “schema” for RDF datasets that satisfy these constraints, enhancing the understandability and usability of the data represented in RDF.

However, the potential of utilizing SHACL to optimize SPARQL queries, similar to how relational constraints optimize queries in relational databases [3], remains largely unexplored. While hand-written SPARQL queries are often tailored to the structure of the RDF graph being queried, queries generated by user interfaces [31,38] or query translators [28,57] and rewriting engines [26,36,52] may contain redundant parts that are not necessary for a specific graph. SHACL can capture information about the structure of an RDF graph, similar to relational database constraints. As a result, triple stores

that have built-in support for SHACL can then use this information to enforce SHACL constraints as well as incorporate them into query processing to optimize queries. To address this issue, we propose a set of optimizations that maintain the correctness of query results and can be applied to a SPARQL query prior to processing, provided that the RDF graph conforms to a specified set of SHACL shapes.

SPARQL queries consist of three main components: (a) pattern matching, which includes optional parts, unions, joins, nesting, and filtering values to find possible matches, (b) solution modifiers, which allow for modifying matched values using operators such as projection and distinct, and (c) output mode, which can be a boolean answer or selections of values from variables that match the patterns, constructing new values, or describing existing data through constraints. To optimize SPARQL queries, we focus on those containing semantically redundant graph patterns that can be minimized or restructured in the query definition based on SHACL descriptions of these patterns in RDF.

The optional operator is a crucial feature in SPARQL queries, commonly used to handle missing or unavailable information [6]. By using optional, query answers can include available information instead of failing to give an answer when parts of the pattern do not match. This feature is particularly important in semantic web applications, where partial knowledge about data is often assumed. However, optional is also a source of complexity in query answering, as it is PSPACE-hard for optional alone [43, 50]. As a result, it has been actively studied [29, 34] and optimized [15, 56] in various query circumstances. In this work, we focus on the Left-Join intuition applied by Xiao et al. [56] for the treatment of optional. However, instead of relying on SQL not null constraint information in their SPARQL-to-SQL translation, we use SHACL constraints to reduce optional patterns to joins. We illustrate this approach using an example below.

Example 1. Consider an RDF graph on the left that validates a SHACL shape s on the right, written in Turtle syntax:

```
:Ida a :Student;
  :hasID "001"^^xsd:int;
  :hasAddress "Oslo".
:Ingrid a :Student;
  :hasID "002"^^xsd:int;
  :hasAddress "Bergen". :StudentNode a sh:NodeShape;
sh:targetClass :Student;
sh:property [ sh:path :hasAddress;
sh:nodeKind sh:Literal;
sh:maxCount 1; sh:minCount 1 ; ];
sh:property [ sh:path :hasID;
dash:uniqueValueForClass
:Student ; ].
```

The node shape, ‘:StudentNode’ declares (a) a class-based target, indicating that all the members of `Student` class are target nodes, (b) cardinality path constraint requiring all students to have exactly one address, and (c) `dash:uniqueValueForClass`¹ path constraint stating that no two students can share the same id. Instances of `Student` validate against the shape if they satisfy both stated constraints.

Consider an information need to retrieve the ids of students including their addresses if available, which could be expressed as the following SPARQL query Q .

```
SELECT ?y ?z WHERE {?x a Student . ?x :hasID ?y . Optional {?x :hasAddress ?z} } .
```

¹ <https://datashapes.org/constraints.html>.

*Q retrieves **ids** of students including their **addresses** if possible from the graph, i.e., in the absence of address, **ids** are still retrieved with variable $?z$ left undefined in the answer. Observe a scenario, where Q is to be evaluated over a graph that satisfies the shape s. Since the shape s guarantees that for each student $?x$, irrespective of their **ids** $?y$, variable $?z$ will always bind to an address (i.e., will never be left unbound), the query Q can be rewritten as follows without the **Optional**,*

```
SELECT ?y ?z WHERE {?x a Student . ?x :hasID ?y . ?x :hasAddress ?z .}
```

*Similarly, consider finding the **ids** of those students who have addresses, such as,*

```
SELECT ?y WHERE {?x a Student . ?x :hasID ?y . ?x :hasAddress ?z .}
```

Since shape s guarantees an address $?z$ for every student $?x$, the original query can be rewritten as follows without changing the solutions,

```
SELECT ?y WHERE {?x a Student . ?x :hasID ?y .}
```

Consider the query $\text{SELECT DISTINCT } ?y \text{ WHERE } \{?x \text{ a Student . } ?x :hasID ?y\}.$ Then, the distinct operator can be removed from the query definition since the shape s guarantees a unique id $?y$ for every student $?x$.

Our approach to optimizing joins in SPARQL is similar to that of query containment concerning inclusion dependencies [27] in relational databases. In particular, when dealing with two query patterns \mathcal{P}_1 and \mathcal{P}_2 under set (resp., bag) semantics, we can simplify their join by reducing it to \mathcal{P}_2 if the answers to \mathcal{P}_1 are included in the answers to \mathcal{P}_2 for any graph. To achieve this, we draw on intuitions from previous literature [9, 13, 33, 49, 51] on join optimization. Specifically, we focus on the task of pruning join patterns that do not affect the output of a query, resulting in reduced query execution costs. Finally, eliminating distinct patterns can also improve query performance significantly as it can be a bottleneck in SPARQL processing, requiring substantial resources to eliminate duplicate rows, particularly with large datasets [7]. We remove distinct if the original query's result is guaranteed to be free of duplicate rows by SHACL. For instance, when joining two triple patterns with a unique one-to-one (resp., one to at most one) relation, applying distinct over the join may not have any effect. An example of this optimization is illustrated in Example 1.

The rest of the paper is organized as follows: In Sect. 2, we review fundamental concepts of SPARQL algebra and SHACL constraints. Query optimizations are presented in Sect. 3. Section 4 contains a discussion of our results, and Sect. 5 concludes the paper. Due to space constraints, the paper includes only partial proofs of the results. For details proofs, we refer the reader to [55].

2 Preliminaries

We next recapitulate the SPARQL algebra and SHACL constraints that we deal with in this paper.

RDF Graph. Assume that \mathbf{I} , \mathbf{B} and \mathbf{L} are countably infinite disjoint sets of *Internationalized Resource Identifiers* (IRIs), *Blank nodes* and *Literals*, respectively. The set of RDF terms \mathbf{T} is $\mathbf{I} \cup \mathbf{B} \cup \mathbf{L}$. An *RDF triple* is an element $\langle s, p, o \rangle$ of $(\mathbf{I} \cup \mathbf{B}) \times \mathbf{I} \times \mathbf{T}$, where s is called the subject, p the predicate and o the object. An RDF graph $G \subseteq (\mathbf{I} \cup \mathbf{B}) \times \mathbf{I} \times \mathbf{T}$ is a finite set of RDF triples. We assume that all RDF graphs are non-empty, which we do not see as a limitation in practice since querying empty graphs has no use in most applications. For simplicity, we represent RDF triples as binary relations in first-order logic, except the triple of the form $\langle s, \text{rdf:type}, C \rangle$. We write $P(s, o)$ (resp., $C(s)$) for the RDF triple $\langle s, P, o \rangle$ (resp., $\langle s, \text{rdf:type}, C \rangle$) and $P^-(s, o)$ for the inverse of triple $P(o, s)$.

Definition 1. The set N_G of nodes of an RDF graph G is the union of sets of subjects and objects of triples of the form $\langle s, P, o \rangle$ and subjects of triples of the form $\langle s, \text{rdf:type}, C \rangle$ in the graph, i.e., $\{s, o \mid P(s, o) \in G \text{ or } C(s) \in G\}$.

SPARQL Algebra. We formally adopt the bag-based (rather than the set-based [43]) semantics for SPARQL as in the W3C specification [24, 45] and studied in the literature [5, 30].

Assume a countably infinite set \mathbf{V} of variables disjoint from \mathbf{T} . A triple pattern is defined as a triple in $(\mathbf{I} \cup \mathbf{L} \cup \mathbf{V}) \times (\mathbf{I} \cup \mathbf{V}) \times (\mathbf{I} \cup \mathbf{L} \cup \mathbf{V})$. We use the same notations $P(x, y)$, $P^-(x, y)$, $C(x)$ for triple patterns as we do for RDF triples. A *basic graph pattern* (*BGP*) is a finite set of triple patterns. For simplicity, we only consider project and distinct solution modifiers and treat SPARQL query as a graph pattern \mathcal{P} , defined by the grammar

$$\begin{aligned} \mathcal{P} ::= & \mathcal{B} \mid \text{FILTER}_{\mathcal{F}}(\mathcal{P}) \mid \text{UNION}(\mathcal{P}_1, \mathcal{P}_2) \mid \text{JOIN}(\mathcal{P}_1, \mathcal{P}_2) \mid \text{MINUS}(\mathcal{P}_1, \mathcal{P}_2) \mid \\ & \text{DIFF}_{\mathcal{F}}(\mathcal{P}_1, \mathcal{P}_2) \mid \text{OPT}_{\mathcal{F}}(\mathcal{P}_1, \mathcal{P}_2) \mid \text{PROJ}_L(\mathcal{P}) \mid \text{DIST}(\mathcal{P}), \end{aligned}$$

where \mathcal{B} is a BGP, $L \subseteq \mathbf{V}$ and \mathcal{F} , called *filter*, is a formula constructed using the logical connectives \wedge and \neg from atoms of the form $\text{bound}(v)$, $(v = c)$, $(v = v')$ for $v, v' \in \mathbf{V}$ and $c \in \mathbf{T}$. Let $\text{var}(\mathcal{P})$ (resp., $\text{var}(\mathcal{F})$) be the set of variables occurring in a graph pattern \mathcal{P} (resp. filter \mathcal{F}). In the rest of the article, we assume that all filter constraints are safe, i.e., $\text{var}(\mathcal{F}) \subseteq \text{var}(\mathcal{P})$ for every pattern $\text{FILTER}_{\mathcal{F}}(\mathcal{P})$.

Given a SPARQL query Q and graph pattern \mathcal{P} , we write $\mathcal{P} \trianglelefteq Q$ if either \mathcal{P} is Q or \mathcal{P} appears in Q . Similarly, we denote by $Q^{\mathcal{P} \rightarrow \mathcal{P}'}$ the rewriting of $\mathcal{P} \trianglelefteq Q$ to graph pattern \mathcal{P}' in Q , $Q^{X \mapsto Y}$ the renaming of $X \subseteq \text{var}(Q)$ to a set Y of variables in Q (i.e., renaming $\rho : X \rightarrow Y$ s.t. $\forall y \in Y, \exists x \in X \text{ and } \rho(x) = y$), $\text{var}(Q \setminus \mathcal{P})$ the set of variables that occur in Q excluding the pattern $\mathcal{P} \trianglelefteq Q$ and $\text{var}(Q \cup \mathcal{P}) := \text{var}(Q) \cup \text{var}(\mathcal{P})$. For readability, we also write nested Join expressions as *concatenation* whenever required, e.g., the graph pattern $\mathcal{P} = \text{JOIN}(\dots \text{JOIN}(\mathcal{P}_1, \mathcal{P}_2), \dots \mathcal{P}_n)$ may be written as $\mathcal{P}_1 \mathcal{P}_2 \dots \mathcal{P}_n$, and the \mathcal{P}_i are called inner graph patterns of \mathcal{P} .

The semantics of graph patterns is defined in terms of (solution) mappings, partial functions, $\mu : \mathbf{V} \rightarrow \mathbf{T}$ with (possibly empty) domain dom_{μ} . Let μ_L (resp., $\mu_{\bar{L}}$) be the restriction of mapping μ to $L \subseteq \mathbf{V}$ (resp. $\mathbf{V} \setminus L$). Two mappings μ_1 and μ_2 are called compatible, denoted by $\mu_1 \sim \mu_2$, if $\mu_1(v) = \mu_2(v)$, for all $v \in \text{dom}_{\mu_1} \cap \text{dom}_{\mu_2}$, in which case $\mu_1 \oplus \mu_2$ denotes a solution mapping with domain $\text{dom}_{\mu_1} \cup \text{dom}_{\mu_2}$ s.t. $\mu_1 \oplus \mu_2 : v \rightarrow \mu_1(v)$ for $v \in \text{dom}_{\mu_1}$, and $\mu_1 \oplus \mu_2 : v \rightarrow \mu_2(v)$ for $v \in \text{dom}_{\mu_2}$. The

truth-value $\mathcal{F}^\mu \in \{\top, \perp, \varepsilon\}$ (i.e., ε stands for “error”) of a filter \mathcal{F} under a mapping μ is defined inductively,

$$\begin{aligned}\square (\text{bound}(v))^\mu &= \begin{cases} \top, & \text{if } v \in \text{dom}_\mu, \\ \perp, & \text{otherwise.} \end{cases} \\ \square (v = c)^\mu &= \begin{cases} \top, & \text{if } \mu(v) = c, \\ \perp, & \text{if } \mu(v) \neq c, \\ \varepsilon \text{ (error),} & \text{if } v \notin \text{dom}_\mu, \end{cases} \quad \text{and } (v = v')^\mu = \begin{cases} \top, & \text{if } \mu(v) = \mu(v'), \\ \perp, & \text{if } \mu(v) \neq \mu(v'), \\ \varepsilon, & \text{if } \{v, v'\} \notin \text{dom}_\mu. \end{cases} \\ \square (\neg \mathcal{F})^\mu &= \begin{cases} \top, & \text{if } \mathcal{F}^\mu = \perp, \\ \perp, & \text{if } \mathcal{F}^\mu = \top, \\ \varepsilon, & \text{if } \mathcal{F}^\mu = \varepsilon, \end{cases} \quad \text{and } (\mathcal{F}_1 \wedge \mathcal{F}_2)^\mu = \begin{cases} \top, & \text{if } \mathcal{F}_1^\mu = \mathcal{F}_2^\mu = \top, \\ \perp, & \text{if } \mathcal{F}_1^\mu = \perp \text{ or } \mathcal{F}_2^\mu = \perp, \\ \varepsilon, & \text{if } \mathcal{F}_1^\mu = \varepsilon \text{ or } \mathcal{F}_2^\mu = \varepsilon. \end{cases}\end{aligned}$$

The evaluation of a SPARQL graph pattern \mathcal{P} over an RDF graph G , denoted by \mathcal{P}^G , returns a multiset (i.e., bag) of mappings. Let $|\mu, \mathcal{P}^G|$ be the multiplicity of a mapping μ in the multiset \mathcal{P}^G . In this sense, we also write $\mu \in \mathcal{P}^G$ when $|\mu, \mathcal{P}^G| > 0$ and $\mu \notin \mathcal{P}^G$ when $|\mu, \mathcal{P}^G| = 0$. Then, the $|\mu, \mathcal{P}^G|$ is defined recursively as follows:

1. Let \mathcal{P} be a BGP. Then, $|\mu, \mathcal{P}^G| = 1$ for every μ if $\text{dom}_\mu = \text{var}(\mathcal{P})$ and $\mu(\mathcal{P}) \subseteq G$, and 0 otherwise;
2. $|\mu, \text{PROJ}_L(\mathcal{P})^G| = \sum_{\mu'=\mu'_L} |\mu', \mathcal{P}^G|$;
3. $|\mu, \text{DIST}(\mathcal{P})^G| = \begin{cases} 1, & \text{if } |\mu, \mathcal{P}^G| > 0, \\ 0, & \text{otherwise;} \end{cases}$;
4. $|\mu, \text{FILTER}_{\mathcal{F}}(\mathcal{P})^G| = \begin{cases} |\mu, \mathcal{P}^G|, & \text{if } \mathcal{F}^\mu = \top, \\ 0, & \text{otherwise;} \end{cases}$;
5. $|\mu, \text{UNION}(\mathcal{P}_1, \mathcal{P}_2)^G| = |\mu, \mathcal{P}_1^G| + |\mu, \mathcal{P}_2^G|$;
6. $|\mu_1 \oplus \mu_2, \text{JOIN}(\mathcal{P}_1, \mathcal{P}_2)^G| = \sum_{\mu_1 \in \mathcal{P}_1^G, \mu_2 \in \mathcal{P}_2^G \text{ with } \mu_1 \sim \mu_2} |\mu_1, \mathcal{P}_1^G| \times |\mu_2, \mathcal{P}_2^G|$;
7. $|\mu, \text{MINUS}(\mathcal{P}_1, \mathcal{P}_2)^G| = \begin{cases} |\mu, \mathcal{P}_1^G|, & \text{if } \forall \mu_2 \in \mathcal{P}_2^G. (\mu \not\sim \mu_2 \text{ or } \text{dom}_\mu \cap \text{dom}_{\mu_2} = \emptyset), \\ 0, & \text{otherwise;} \end{cases}$;
8. $|\mu, \text{DIFF}_{\mathcal{F}}(\mathcal{P}_1, \mathcal{P}_2)^G| = \begin{cases} |\mu, \mathcal{P}_1^G|, & \text{if } \forall \mu_2 \in \mathcal{P}_2^G. (\mu \not\sim \mu_2 \text{ or } \mathcal{F}^{\mu \oplus \mu_2} \neq \top), \\ 0, & \text{otherwise;} \end{cases}$;
9. $|\mu, \text{OPT}_{\mathcal{F}}(\mathcal{P}_1, \mathcal{P}_2)^G| = |\mu, \text{FILTER}_{\mathcal{F}}(\text{JOIN}(\mathcal{P}_1, \mathcal{P}_2))^G| + |\mu, \text{DIFF}_{\mathcal{F}}(\mathcal{P}_1, \mathcal{P}_2)^G|$.

The support of a multiset M is the underlying set $\text{sup}(M) = \{\mu \mid |\mu, M| > 0\}$. The domain of a multiset M of mappings is defined as $\text{dom}_M = \bigcup_{\mu \in M} \text{dom}_\mu$. We denote by $M|_{\bar{X}}$ the multiset of mappings $\mu \in M$ restricted to $V \setminus X$, i.e., $|\mu, M|_{\bar{X}} = \sum_{\mu'=\mu'_{\bar{X}}} |\mu', M|$.

Finally, `Proj` and `Dist` operators capture the query nesting functionality, called *sub-queries* [30, Sect. 3], in SPARQL, which can be arbitrarily deep and may lead to a complex layered structure. We thus adopt the *simple normal form* (SNF) of SPARQL queries from [30, Defn. 3.8] that is suitable for optimizations presented in Sect. 3.

Definition 2. A SPARQL query is in SNF if it has the form `DIST(ProjX(P))` with subquery-free pattern \mathcal{P} or the form `ProjX(P)`, where all subquery patterns in \mathcal{P} are of the form `DIST(ProjX'(P'))` with \mathcal{P}' subquery-free.

Any SPARQL query Q can be brought into SNF by following two steps of normalisation as illustrated below, where $n \geq i \geq 1$:

1. If Q is a query of the form $\text{DIST}(\text{PROJ}_X(\mathcal{P}))$ s.t. \mathcal{P} is subquery $\text{DIST}(\mathcal{P}')$, then Q can be simplified to $\text{DIST}(\text{PROJ}_X(\mathcal{P}'))$.
2. If Q is a query of the form $\text{PROJ}_X(\mathcal{P})$ s.t. $\mathcal{P} = \{\mathcal{P}_1 \dots \text{PROJ}_Y(\mathcal{P}_i) \dots \mathcal{P}_n\}$, then Q can be simplified to $\text{PROJ}_X(\mathcal{P}')$, where $\mathcal{P}' = \{\mathcal{P}_1 \dots \mathcal{P}_i^{W \mapsto W'} \dots \mathcal{P}_n\}$ s.t. $W = (\text{var}(\mathcal{P}_i) \setminus Y) \cap \text{var}(Q \setminus \mathcal{P}_i)$ and $w' \notin \text{var}(Q)$ for all $w' \in W'$.

Step 1 removes redundant `Dist` from inner graph pattern \mathcal{P}' since the outermost `Dist` already guarantees a duplicate-free result. To avoid conflicts with variables outside of \mathcal{P}_i in Q , step 2 renames variables in inner graph pattern \mathcal{P}_i before removing the inner `Proj`.

Henceforth, we assume that all SPARQL queries are in their *SNF*.

Example 2. Consider a SPARQL query that retrieves the name of employees and their office addresses, such as,

$\text{PROJ}_{yz}(\text{JOIN}(\text{hasName}(x, y), \text{PROJ}_{xz}(\text{JOIN}(\text{hasOffice}(x, y), \text{hasAddress}(y, z))))),$

which can be brought into SNF by renaming the variable y of subquery to n , thus removing clashes with any variables outside of the subquery, as follows,

$\text{PROJ}_{yz}(\text{hasName}(x, y) \text{ hasOffice}(x, n) \text{ hasAddress}(n, z)).$

SHACL Constraint. We adopt the abstract syntax for the *core constraints* of SHACL [32] introduced by Corman et al. [19]. In addition, we incorporate a non-standard SPARQL-based constraint component called `dash:uniqueValuesForClass`. Let S , C and P be countably infinite and mutually disjoint sets of SHACL *shape*, *class* and *property* names, respectively. Each SHACL constraint is a set of conditions, usually referred to as *shape*, defined as a triple $\langle s, \tau_s, \phi_s \rangle$ consisting of:

- *name* s ,
- *target definition* τ_s is a SPARQL query with a single output variable, which retrieves the target entities of s from the RDF graph G that is being validated. There are two types of τ_s that we consider:
 - The first type corresponds to the ‘sh:targetClass’ of the SHACL specification. In this case, τ_s is a SPARQL query of the form:

“SELECT ?x WHERE ?x rdf:type/subClassOf* C”.

- The second type corresponds to the ‘sh:targetSubjectOf’ or ‘sh:targetObjectOf’ of the SHACL specification, as shown in Example 3. In this case, τ_s is a SPARQL query of the form:

“SELECT ?x WHERE ?x P ?y” or “SELECT ?x WHERE ?y P ?x”.

In abstract syntax, we express τ_s using the following grammar:

$$\tau_s ::= C \mid \exists P \mid \exists P^-$$

where τ_s being C , $\exists P$ and $\exists P^-$ represent ‘sh:targetClass C’, ‘sh:targetSubjectOf P’, and ‘sh:targetObjectOf P’ in the SHACL specification, respectively. Given a target expression τ_s , we say that $n \in \mathcal{N}_G$ is a target of the shape s , written $G \models \tau_s(n)$, iff $\exists \mu \in \tau_s^G$ s.t. $n \in \mu(var(\tau_s))$.

- *constraint definition* ϕ_s , which can be represented as a boolean SPARQL query as described in [10, 17, 18], indicating whether the target node under validation violates the graph pattern defined by the ϕ_s or not. The constraint ϕ_s is an expression defined according to the following grammar:

$$\begin{aligned}\phi_s &::= \geq_n \alpha. \beta \mid \leq_n \alpha. \beta \mid \triangleright_{\tau_s} \alpha \mid \alpha_1 = \alpha_2 \mid \phi_s \wedge \phi_s \\ \beta &::= \top \mid C \mid s' \mid \neg \beta\end{aligned}$$

where \top stands for the *Boolean* true value, $\alpha, \alpha_1, \alpha_2 \in (\mathbf{P} \cup \{P^- \mid P \in \mathbf{P}\})$, $C \in \mathbf{C}$, $s' \in \mathbf{S}$, \neg stands for *negation*, $n \in \mathbb{N}$, $(\geq_n \alpha. \beta)$ requires that there must be at least n α -successors verifying β , $(\triangleright_{\tau_s} \alpha)$ that the value node of α -successor must be unique among target nodes defined by the τ_s , e.g., `dash:uniqueValueForClass` constraint for the class-based target, and $(\alpha_1 = \alpha_2)$ means equality of the sets of nodes reachable via the respective properties α_1 and α_2 . As syntactic sugar, we write $\phi_1 \vee \phi_2$ for $\neg(\neg \phi_1 \wedge \neg \phi_2)$ and $=_n \alpha. \phi$ for $(\geq_n \alpha. \phi) \wedge (\leq_n \alpha. \phi)$. For a constraint expression ϕ_s and a node $n \in \mathcal{N}_G$, we say that n validates against ϕ_s , written $G \models \phi_s(n)$, iff n does not violate the graph pattern, aka constraints, defined by ϕ_s . For any constraints ϕ and ϕ' , we say that ϕ implies ϕ' , rewritten $\phi \longrightarrow \phi'$, iff for all graphs G and nodes $n \in \mathcal{N}_G$, if $G \models \phi(n)$ then $G \models \phi'(n)$.

A SHACL document is a set of SHACL shapes. An RDF graph G validates against a shape $\langle s, \tau_s, \phi_s \rangle$, if $G \models \phi_s(n)$ for all $n \in \mathcal{N}_G$ with $G \models \tau_s(n)$. An RDF graph G validates against a SHACL document \mathcal{S} , written $G \models \mathcal{S}$, iff G validates against all shapes in \mathcal{S} .

Example 3. Assume an RDF graph on the left that validates a shape with a ‘property-based target,’ i.e., sh:targetSubjectOf :hasID, on the right written in Turtle syntax:

```
:Ida a :Student;
  :hasID "001"^^xsd:int;
  :hasAddress "Oslo".
:Nora a :Student;
  :hasAddress "Oslo".
:Ingrid a :Student;
  :hasID "002"^^xsd:int;
  :hasAddress "Bergen". :TargetSubjectShape a sh:NodeShape;
  sh:targetSubjectOf :hasID;
  sh:property [ sh:path :hasAddress;
    sh:nodeKind sh:Literal;
    sh:minCount 1; ];
  sh:property [ sh:path :hasID;
    dash:uniqueValueForClass
    :Student ; ].
```

Observe that ‘Nora’ is not subject of a ‘hasID’ triple, and therefore, is not a target node of the shape ‘TargetSubjectShape’. Since the target nodes, i.e., ‘Ida’ and ‘Ingrid’, of the shape have an address and unique ids, they validate against the shape.

3 Optimizations

The goal is to optimize a SPARQL query with respect to a SHACL document, by identifying smaller or more efficient queries that have the same answers for all RDF graphs that satisfy the SHACL document.

Next we will present SPARQL query optimizations that resemble SQL query rewriting using relational constraints. They only differ in how SPARQL graph patterns interact with SHACL constraints describing these graph patterns in the RDF graph. For instance, when property path cardinality constraints imply that the first argument of an Opt-pattern will always find at least one solution on the second, we can substitute Opt with a more straightforward join operator, as shown in Lemmas 1.1, 3.1 and 6.1. In the case of a distinct projection, with an ‘at least one’ property path cardinality constraint, we can infer that every solution mapping from the first argument of the join expression is guaranteed to find at least one solution on the second argument. This guarantee on the graph patterns allows us to reduce the entire joint expression to the first argument, as shown in Lemmas 1.2 and 3.2. In the case of an ‘exactly one’ property path cardinality constraint, the described optimization of join patterns occurs even in the absence of distinct projection, as shown in Lemmas 2 and 4. If all possible mappings for the first argument of the optional pattern can find a compatible solution with the second argument, and if the variables used in filter are limited to the first argument, then it is possible to replace the entire optional pattern with the filter expression applied to the first argument, as shown in Lemma 5.2, Corollaries 2.2, 4.2 and 5.2.

Let \mathcal{U} and \mathcal{V} be two graph patterns and \mathcal{S} a SHACL document.

Definition 3. $\mathcal{U} \equiv_{\mathcal{S}} \mathcal{V}$ iff $\mathcal{U}^G = \mathcal{V}^G$ for all graphs G with $G \models \mathcal{S}$.

Definition 4. $\mathcal{U} \equiv_{\mathcal{S},y} \mathcal{V}$ iff, for all graphs G with $G \models \mathcal{S}$,

1. $\mathcal{U}^G|_{\bar{y}} = \mathcal{V}^G|_{\bar{y}}$ and
2. $(\mu|_{\bar{y}} = \mu'|_{\bar{y}}) \longrightarrow (\mu = \mu')$ for all $\mu, \mu' \in \mathcal{U}^G$ and $\mu, \mu' \in \mathcal{V}^G$.

Definition 5. $\mathcal{U} \cong_{\mathcal{S},y} \mathcal{V}$ iff $\text{sup}(\mathcal{U}^G|_{\bar{y}}) = \text{sup}(\mathcal{V}^G|_{\bar{y}})$ for all graphs G with $G \models \mathcal{S}$.

In Theorem 1, we will prove that an $\equiv_{\mathcal{S}}$ equivalence is applicable to all queries, while $\equiv_{\mathcal{S},y}$ is restricted to the placement of variable y within the query, and $\cong_{\mathcal{S},y}$ is exclusive to the distinct graph patterns with restrictions on the placement of variable y . In the subsequent Lemmas 1 to 7, we will establish these equivalences under various constraints, and illustrate their application along with counterexamples when restrictions on the placement of variables are violated in the query.

Let \mathcal{B} and \mathcal{B}' be two BGPs and \mathcal{P} a graph pattern. Then, we write $\mathcal{B} \blacktriangleleft \mathcal{P}$ (resp., $\mathcal{B}, \mathcal{B}' \blacktriangleleft \mathcal{P}$) if $\mathcal{P} = \mathcal{B}$ or $\mathcal{P} = \mathcal{B}\mathcal{P}_1 \dots \mathcal{P}_n$ (resp., $\mathcal{P} = \mathcal{B}\mathcal{B}'$ or $\mathcal{P} = \mathcal{B}\mathcal{B}'\mathcal{P}_1 \dots \mathcal{P}_n$), where \mathcal{P}_i s.t. $n \geq i \geq 1$ are inner graph patterns. Let $C, C' \in \mathbf{C}$ and $P, R \in (\gamma \cup \{\gamma^- \mid \gamma \in \mathbf{P}\})$. Then, in the lemmas that follow, given an arbitrary shape $\langle s, \tau_s, \phi_s \rangle$ with target τ_s definition, let T be a triple pattern s.t.,

$$T = \begin{cases} C(x), & \text{if } \tau_s = C, \\ R(x, z), & \text{if } \tau_s = \exists R, \\ R^-(x, z), & \text{if } \tau_s = \exists R^-. \end{cases}$$

Lemma 1. Let $\langle s, \tau_s, \phi_s \rangle \in \mathcal{S}$ with $(\geq_n P.\top) \in \phi_s$ s.t. $n \geq 1$, and \mathcal{P} a graph pattern s.t. $T \blacktriangleleft \mathcal{P}$. If $y \notin \text{var}(\mathcal{P})$, then

1. $\text{OPT}_{\mathcal{F}}(\mathcal{P}, P(x, y)) \equiv_{\mathcal{S}} \text{FILTER}_{\mathcal{F}}(\text{JOIN}(\mathcal{P}, P(x, y)))$
2. $\text{JOIN}(\mathcal{P}, P(x, y)) \cong_{\mathcal{S}, y} \mathcal{P}$

Proof. Assuming $T \blacktriangleleft \mathcal{P}$ w.r.t. the target τ_s of the shape $\langle s, \tau_s, \phi_s \rangle \in \mathcal{S}$, we know that x is the only variable shared between \mathcal{P} and $P(x, y)$, and y is not in $\text{var}(\mathcal{P})$. Let G be a graph s.t. $G \models \mathcal{S}$.

Clause 1: Using the definition of OPT , we can split the solutions of $\text{OPT}_{\mathcal{F}}(\mathcal{P}, P(x, y))$ into two parts: finding solutions that satisfy both \mathcal{P} and $P(x, y)$, and finding solutions that satisfy only \mathcal{P} . For any mapping in the first part \mathcal{P}^G , we can find at least n mappings in the second part $P(x, y)^G$, where $n \geq 1$ is guaranteed by the constraint $(\geq_n P.\top) \in \phi_s$. Therefore, any solution that satisfies \mathcal{P} and $P(x, y)$ will also be part of the solutions that satisfy only \mathcal{P} w.r.t. G and vice-versa, which means that $|\mu, \text{OPT}_{\mathcal{F}}(\mathcal{P}, P(x, y))^G| = |\mu, \text{FILTER}_{\mathcal{F}}(\text{JOIN}(\mathcal{P}, P(x, y)))^G|$. Thus, $\text{OPT}_{\mathcal{F}}(\mathcal{P}, P(x, y)) \equiv_{\mathcal{S}} \text{FILTER}_{\mathcal{F}}(\text{JOIN}(\mathcal{P}, P(x, y)))$ by Definition 3 if $y \notin \text{var}(\mathcal{P})$.

Clause 2: Starting from the left-hand side $\text{JOIN}(\mathcal{P}, P(x, y))$, we have $y \notin \text{var}(\mathcal{P})$. Using the definition of JOIN , the problem is to find all possible combinations of solutions that satisfy \mathcal{P} and $P(x, y)$. Since constraint $(\geq_n P.\top) \in \phi_s$ together with conditions $T \blacktriangleleft \mathcal{P}$ and $y \notin \text{var}(\mathcal{P})$ guarantees that there exist least n mappings $\mu_2 \in P(x, y)^G$ for $\forall \mu_1 \in \mathcal{P}^G$ s.t. $\mu_1 \sim \mu_2$, we can deduce that any solution that satisfies \mathcal{P} will also satisfy $\text{JOIN}(\mathcal{P}, P(x, y))$ w.r.t. G . Thus, we can claim that $\text{sup}(\text{JOIN}(\mathcal{P}, P(x, y))^G|_{\bar{y}})$ is equal to $\text{sup}(\mathcal{P}^G|_{\bar{y}})$, i.e., $\text{JOIN}(\mathcal{P}, P(x, y)) \cong_{\mathcal{S}, y} \mathcal{P}$ by Definition 4.

Corollary 1 follows from Lemma 1.

Corollary 1. Let $\langle s, \tau_s, \phi_s \rangle \in \mathcal{S}$ with $(\geq_n P.\top) \in \phi_s$ s.t. $n \geq 1$, and \mathcal{P} a graph pattern s.t. $T \blacktriangleleft \mathcal{P}$. If $y \notin \text{var}(\mathcal{P} \cup \mathcal{F})$, then

1. $\text{FILTER}_{\mathcal{F}}(\text{JOIN}(\mathcal{P}, P(x, y))) \cong_{\mathcal{S}, y} \text{FILTER}_{\mathcal{F}}(\mathcal{P})$
2. $\text{OPT}_{\mathcal{F}}(\mathcal{P}, P(x, y)) \cong_{\mathcal{S}, y} \text{FILTER}_{\mathcal{F}}(\mathcal{P})$

Example 4. Let G be a graph such that $G \models \langle \text{Student}, \tau_{\text{Student}}, \phi_{\text{Student}} \rangle$ with $(\geq_n : \text{hasAddress}.\top) \in \phi_{\text{Student}}$ s.t. $n \geq 1$. Then, let Q be a query that asks for all students who do not have a postal address, where \top denotes the tautological filter (true),

$$\text{PROJ}_x(\text{FILTER}_{\neg \text{bound}(y)}(\text{OPT}_{\top}(\text{Student}(x), \text{hasAddress}(x, y))))$$

over G . As per the constraint ϕ_{Student} , all students in G have a postal address. Hence, the optimal solution would be to replace the entire query with an empty query. However, in our setting, the OPT -pattern of Q can be reduced to an equivalent JOIN -pattern by following the clause 1 of Lemma 1 since y does not appear in the first part of the OPT -pattern,

$$\text{PROJ}_x(\text{FILTER}_{\neg \text{bound}(y)}(\text{JOIN}(\text{Student}(x), \text{hasAddress}(x, y)))).$$

In a scenario where y also appears in the first part of the OPT -pattern, such as

$$\text{PROJ}_x(\text{FILTER}_{\neg \text{bound}(y)}(\text{OPT}_{\top}(\text{Student}(x), \text{enrolledIn}(x, y), \text{hasAddress}(x, y)))).$$

A counter-example $G = \{\text{Student}(a), \text{enrolledIn}(a, b), \text{hasAddress}(a, d)\}$ s.t. $G \models \langle \text{Student}, \tau_{\text{Student}}, \phi_{\text{Student}} \rangle$ can be constructed, where a solution can be found that yields “ $\text{Student}(a)$ ” for the Opt pattern but none when Opt is reduced to Join-pattern.

Similarly, assume a query that asks for students who have postal addresses, such as

$$\text{DIST}(\text{PROJ}_{xy}(\text{JOIN}(\text{Student}(x) \text{ hasName}(x, y), \text{hasAddress}(x, z))))$$

over the graph G . Then, clause 2 of Lemma 1 allows us to safely remove the pattern $\text{hasAddress}(x, z)$ from the ‘Dist’ query, given that variable z only appears in this pattern and all students in the graph G satisfying $\langle \text{Student}, \tau_{\text{Student}}, \phi_{\text{Student}} \rangle$ have at least n postal addresses., i.e.,

$$\text{DIST}(\text{PROJ}_{xy}(\text{Student}(x) \text{ hasName}(x, y))).$$

However, a counter-example can be found for the $\equiv_{S,y}$ equivalence when variable z appears outside of the pattern $\text{hasAddress}(x, z)$ in the Dist query above, such as

$$\text{DIST}(\text{PROJ}_{xy}(\text{FILTER}_{\neg \text{bound}(z)}(\text{JOIN}(\text{Student}(x) \text{ hasName}(x, y), \text{hasAddress}(x, z))))).$$

Finally, assume a query scenario over the G as per required by the ‘if-condition’ in Corollary 1,

$$\text{DIST}(\text{PROJ}_{xz}(\text{OPT}_{z=A'}(\text{Student}(x) \text{ hasGrade}(x, z), \text{hasAddress}(x, y)))).$$

Then, the equivalent “ $\text{DIST}(\text{PROJ}_{xz}(\text{FILTER}_{z=A'}(\text{Student}(x) \text{ hasGrade}(x, z))))$ ” query can be deduced by using Corollary 1.2, which reduces the entire ‘Opt’-pattern to the first part of the pattern with the filter applied (i.e., $\equiv_{S,y}$ equivalence).

Lemma 2 follows from the same reasoning as Lemma 1.

Lemma 2. Let $\langle s, \tau_s, \phi_s \rangle \in \mathcal{S}$ with $(=_1 P.T) \in \phi_s$ and \mathcal{P} a graph pattern s.t. $T \blacktriangleleft \mathcal{P}$. If $y \notin \text{var}(\mathcal{P})$, then $\text{JOIN}(\mathcal{P}, P(x, y)) \equiv_{S,y} \mathcal{P}$.

Corollary 2 follows from Lemma 1 and Lemma 2.

Corollary 2. Let $\langle s, \tau_s, \phi_s \rangle \in \mathcal{S}$ with $(=_1 P.T) \in \phi_s$ s.t. $n \geq 1$, and \mathcal{P} a graph pattern s.t. $T \blacktriangleleft \mathcal{P}$. If $y \notin \text{var}(\mathcal{P} \cup \mathcal{F})$, then

1. $\text{FILTER}_{\mathcal{F}}(\text{JOIN}(\mathcal{P}, P(x, y))) \equiv_{S,y} \text{FILTER}_{\mathcal{F}}(\mathcal{P})$
2. $\text{OPT}_{\mathcal{F}}(\mathcal{P}, P(x, y)) \equiv_{S,y} \text{FILTER}_{\mathcal{F}}(\mathcal{P})$

Example 5. consider a SPARQL query that asks for the names of those who have postal addresses, such as

$$\text{PROJ}_y(\text{JOIN}(\text{hasName}(x, y), \text{hasAddress}(x, z))),$$

over a graph G s.t. $G \models \langle \exists \text{hasName}, \tau_{\exists \text{hasName}}, \phi_{\exists \text{hasName}} \rangle$ and $(=_1 :hasAddress. \top) \in \phi_{\exists \text{hasName}}$. Since the constraint $\phi_{\exists \text{hasName}}$ ensures that individuals with names have exactly one postal address, we can deduce the following equivalent query based on Lemma 2:

$$\text{PROJ}_y(\text{hasName}(x, y)).$$

In Lemmas 3 to 5,

- let β be either C' or s' s.t. $s' \in \mathcal{S}$, and
- let \mathcal{P}' be a graph pattern s.t. \mathcal{P}' is either $P(x, y)$ or $P(x, y)C'(y)$ if $\beta = s'$, and $P(x, y)C'(y)$ if $\beta = C'$.

Lemma 3. Let $\langle s, \tau_s, \phi_s \rangle \in S$ with $(\geq_n P. \beta) \in \phi_s$ s.t. $n \geq 1$, and \mathcal{P} a graph pattern s.t. $T \blacktriangleleft \mathcal{P}$. If $y \notin var(\mathcal{P})$, then

1. $OPT_{\mathcal{F}}(\mathcal{P}, \mathcal{P}') \equiv_S FILTER_{\mathcal{F}}(JOIN(\mathcal{P}, \mathcal{P}'))$
2. $JOIN(\mathcal{P}, \mathcal{P}') \cong_{S,y} \mathcal{P}$

In contrast to Lemma 1, Lemma 3 also handles type constraints on the value node of a property predicate P . Corollary 3 follows from Lemma 3.

Corollary 3. Let $\langle s, \tau_s, \phi_s \rangle \in S$ with $(\geq_n P. \beta) \in \phi_s$ s.t. $n \geq 1$, and \mathcal{P} a graph pattern s.t. $T \blacktriangleleft \mathcal{P}$. If $y \notin var(\mathcal{P} \cup \mathcal{F})$, then

1. $FILTER_{\mathcal{F}}(JOIN(\mathcal{P}, \mathcal{P}')) \cong_{S,y} FILTER_{\mathcal{F}}(\mathcal{P})$
2. $OPT_{\mathcal{F}}(\mathcal{P}, \mathcal{P}') \cong_{S,y} FILTER_{\mathcal{F}}(\mathcal{P})$

Example 6. Suppose we have a SPARQL query that asks the names of individuals who have valid addresses, such as

$$DIST(PROJ_y(JOIN(hasName(x, y), hasAddress(x, z) Address(z)))) ,$$

over a graph G s.t. $G \models \langle \exists hasName, \tau_{hasName}, \phi_{hasName} \rangle$ and $(\geq_1 :hasAddress. Address) \in \phi_{people}$. Then, we can use Lemma 3, specifically clause 2, to simplify the query to:

$$DIST(PROJ_y(hasName(x, y))) .$$

Lemma 4. Let $\langle s, \tau_s, \phi_s \rangle \in S$ with $(=_1 P. \beta) \in \phi_s$, and \mathcal{P} a graph pattern s.t. $T \blacktriangleleft \mathcal{P}$. If $y \notin var(\mathcal{P})$, then $JOIN(\mathcal{P}, \mathcal{P}') \cong_{S,y} \mathcal{P}$.

Lemma 4 follows from the same reasoning as Lemma 3, and Corollary 4 follows from Lemma 3 and Lemma 4.

Corollary 4. Let $\langle s, \tau_s, \phi_s \rangle \in S$ with $(=_1 P. \beta) \in \phi_s$ s.t. $n \geq 1$, and \mathcal{P} a graph pattern s.t. $T \blacktriangleleft \mathcal{P}$. If $y \notin var(\mathcal{P} \cup \mathcal{F})$, then

1. $FILTER_{\mathcal{F}}(JOIN(\mathcal{P}, \mathcal{P}')) \cong_{S,y} FILTER_{\mathcal{F}}(\mathcal{P})$
2. $OPT_{\mathcal{F}}(\mathcal{P}, \mathcal{P}') \cong_{S,y} FILTER_{\mathcal{F}}(\mathcal{P})$

Lemma 5. Let $\langle s, \tau_s, \phi_s \rangle \in S$ with $(\geq_n P. \beta) \in \phi_s$ s.t. $n \geq 1$ or $(\leq_0 P. \neg\beta) \in \phi_s$, and \mathcal{P} a graph pattern s.t. $T, P(x, y) \blacktriangleleft \mathcal{P}$. Then,

1. $JOIN(\mathcal{P}, C'(y)) \equiv_S \mathcal{P}$
2. $OPT_{\mathcal{F}}(\mathcal{P}, C'(y)) \equiv_S FILTER_{\mathcal{F}}(\mathcal{P})$

Proof. Assuming $T, P(x, y) \blacktriangleleft \mathcal{P}$ w.r.t. the target τ_s of the shape $\langle s, \tau_s, \phi_s \rangle \in \mathcal{S}$, we have $\text{var}(\mathcal{P}) \cap \text{var}(C'(y)) = \{y\}$. Let G be a graph s.t. $G \models \mathcal{S}$. Then,

Clause 1: Starting from the left-hand side pattern $\text{JOIN}(\mathcal{P}, C'(y))$, the constraint $(\geq_n P.\beta) \in \phi_s$ (resp., $(\leq_0 P.\neg\beta) \in \phi_s$) s.t. $n \geq 1$ together with the condition $P(x, y) \blacktriangleleft \mathcal{P}$ guarantees that there exists exactly one mapping $\mu_2 \in C'(y)^G$ for $\forall \mu_1 \in \mathcal{P}^G$ s.t. $\mu_1 \sim \mu_2$. Thus, we can deduce $|\mu, \text{JOIN}(\mathcal{P}, P(x, y))^G| = |\mu, \mathcal{P}^G|$, which implies $\text{JOIN}(\mathcal{P}, P(x, y)) \equiv_{\mathcal{S}} \mathcal{P}$ by Definition 3.

Clause 2: For the pattern $\text{OPT}_{\mathcal{F}}(\mathcal{P}, C'(y))$, constraint $(\geq_n P.\beta) \in \phi_s$ (resp., $(\leq_0 P.\neg\beta) \in \phi_s$) together with the condition $P(x, y) \blacktriangleleft \mathcal{P}$ guarantees that for any mapping in the first part \mathcal{P}^G , we can find exactly one mapping in the second part $C'(y)^G$. Thus, any solution that satisfies \mathcal{P} and $C'(y)$ will also be part of the solutions that satisfy only \mathcal{P} w.r.t. G and vice-versa, which means that $|\mu, \text{OPT}_{\mathcal{F}}(\mathcal{P}, C'(y))^G| = |\mu, \text{FILTER}_{\mathcal{F}}(\text{JOIN}(\mathcal{P}, C'(y)))^G|$, i.e., $\text{OPT}_{\mathcal{F}}(\mathcal{P}, C'(y)) \equiv_{\mathcal{S}} \text{FILTER}_{\mathcal{F}}(\text{JOIN}(\mathcal{P}, C'(y)))$ by Definition 3.

Further, the restriction of solutions to $\mathcal{F}^\mu = \top$ is identical for both $\text{JOIN}(\mathcal{P}, C'(y))^G$ and \mathcal{P}^G since $\text{var}(\mathcal{F}) \subseteq \text{var}(\mathcal{P})$ can be deduced from the fact $\text{var}(C'(y)) \subseteq \text{var}(\mathcal{P})$. Then, using clause 1, i.e., $\text{JOIN}(\mathcal{P}, C'(y)) \equiv_{\mathcal{S}} \mathcal{P}$, we can infer that $\text{FILTER}_{\mathcal{F}}(\text{JOIN}(\mathcal{P}, C'(y))) \equiv_{\mathcal{S}} \text{FILTER}_{\mathcal{F}}(\mathcal{P})$. Thus,

$$\text{OPT}_{\mathcal{F}}(\mathcal{P}, C'(y)) \equiv_{\mathcal{S}} \text{FILTER}_{\mathcal{F}}(\text{JOIN}(\mathcal{P}, C'(y))) \equiv_{\mathcal{S}} \text{FILTER}_{\mathcal{F}}(\mathcal{P}).$$

Example 7. Let G be a graph that satisfies the shapes $\langle \exists \text{hasName}, \tau_{\exists \text{hasName}}, \phi_{\exists \text{hasName}} \rangle$ and $\langle \text{Program}, \tau_{\text{Program}}, \phi_{\text{Program}} \rangle$, where $(\leq_0 : \text{enrolledIn}.\neg, \text{Program}) \in \phi_{\exists \text{hasName}}$. Consider the query,

$$\text{PROJ}_{xz}(\text{OPT}_{\top}(\text{hasName}(x, y) \text{ enrolledIn}(x, z), \text{Program}(z)))$$

over the G . According to the constraint $\phi_{\exists \text{hasName}}$, all enrollments are only in the ‘Program’. Hence, using clause 1 and 2 of Lemma 5, we can obtain the following equivalent query:

$$\text{PROJ}_{xz}(\text{hasName}(x, y) \text{ enrolledIn}(x, z)).$$

Given a shape $\langle s, \tau_s, \phi_s \rangle \in \mathcal{S}$, let \mathcal{K} be a graph pattern s.t.,

$$\mathcal{K} = \begin{cases} C(x)R(x, z), & \text{if } \tau_s = C, \\ R(x, z), & \text{if } \tau_s = \exists R, \\ R^-(x, z), & \text{if } \tau_s = \exists R^-. \end{cases}$$

Lemma 6. Let $\langle s, \tau_s, \phi_s \rangle \in \mathcal{S}$ with $(R = P) \in \phi_s$, and \mathcal{P} a graph pattern s.t. $\mathcal{K} \blacktriangleleft \mathcal{P}$. If $y \notin \text{var}(\mathcal{P})$, then

1. $\text{OPT}_{\mathcal{F}}(\mathcal{P}, P(x, y)) \equiv_{\mathcal{S}} \text{FILTER}_{\mathcal{F}}(\text{JOIN}(\mathcal{P}, P(x, y)))$
2. $\text{JOIN}(\mathcal{P}, P(x, y)) \cong_{\mathcal{S}, y} \mathcal{P}$

Proof. Assuming $\mathcal{K} \blacktriangleleft \mathcal{P}$ w.r.t. the target τ_s of the shape $\langle s, \tau_s, \phi_s \rangle \in \mathcal{S}$, we have $\text{var}(\mathcal{P}) \cap \text{var}(P(x, y)) = \{x\}$ since $R(x, z) \blacktriangleleft \mathcal{K}$. Let G be a graph s.t. $G \models \mathcal{S}$.

Clause 1: For $\text{OPT}_{\mathcal{F}}(\mathcal{P}, P(x, y))$ s.t. $y \notin \text{var}(\mathcal{P})$, the constraint $(R = P) \in \phi_s$ together with the condition $R(x, z) \blacktriangleleft \mathcal{P}$ guarantees that for any mapping in the first part \mathcal{P}^G ,

we can find exactly one mapping in the second part $P(x, y)^G$. Thus, any solution that satisfies \mathcal{P} and $P(x, y)$ will also be part of the solutions that satisfy only \mathcal{P} w.r.t. G and vice-versa, which means that $|\mu, \text{OPT}_{\mathcal{F}}(\mathcal{P}, P(x, y))^G| = |\mu, \text{FILTER}_{\mathcal{F}}(\text{JOIN}(\mathcal{P}, P(x, y)))^G|$ if $y \notin \text{var}(\mathcal{P})$, i.e.,

$$\text{OPT}_{\mathcal{F}}(\mathcal{P}, P(x, y)) \equiv_S \text{FILTER}_{\mathcal{F}}(\text{JOIN}(\mathcal{P}, P(x, y))) \text{ by Definition 3.}$$

Clause 2: Consider the left-hand side $\text{JOIN}(\mathcal{P}, P(x, y))$ s.t. $y \notin \text{var}(\mathcal{P})$. Then, the constraint $(R = P) \in \phi_s$ together with the condition $R(x, z) \blacktriangleleft \mathcal{P}$ guarantees that there exists exactly one mapping $\mu_2 \in P(x, y)^G$ for $\forall \mu_1 \in \mathcal{P}^G$ s.t. $\mu_1 \sim \mu_2$. From this, we can deduce that $\text{sup}(\text{JOIN}(\mathcal{P}, P(x, y))^G) = \text{sup}(\mathcal{P}^G)$, which in turn implies $\text{JOIN}(\mathcal{P}, P(x, y)) \cong_{S,y} \mathcal{P}$ by Definition 5, if $y \notin \text{var}(\mathcal{P})$.

Corollary 5 follows from Lemma 6.

Corollary 5. Let $\langle s, \tau_s, \phi_s \rangle \in S$ with $(R = P) \in \phi_s$, and \mathcal{P} a graph pattern s.t. $\mathcal{K} \blacktriangleleft \mathcal{P}$. If $y \notin \text{var}(\mathcal{P} \cup \mathcal{F})$, then

1. $\text{FILTER}_{\mathcal{F}}(\text{JOIN}(\mathcal{P}, P(x, y))) \cong_{S,y} \text{FILTER}_{\mathcal{F}}(\mathcal{P})$
2. $\text{OPT}_{\mathcal{F}}(\mathcal{P}, P(x, y)) \cong_{S,y} \text{FILTER}_{\mathcal{F}}(\mathcal{P})$

Example 8. Consider a SPARQL query,

$$\text{DIST}(\text{PROJ}_{xy}(\text{JOIN}(\text{Employee}(x) \text{ insuredBy}(x, y), \text{employedBy}(x, z))))$$

over G s.t. $G \models \langle \text{Employee}, \tau_{\text{Employee}}, \phi_{\text{Employee}} \rangle$ with $(\text{insuredBy} = \text{employedBy}) \in \phi_{\text{Employee}}$. Then, using clause 2 of Lemma 6, we can obtain the following equivalent query:

$$\text{DIST}(\text{PROJ}_{xy}(\text{Employee}(x) \text{ insuredBy}(x, y))).$$

Similarly, assume an optional query scenario where variable z only occurs within the pattern ‘ $\text{employedBy}(x, z)$ ’, such as

$$\text{DIST}(\text{PROJ}_{xy}(\text{OPT}_{\top}(\text{Employee}(x) \text{ insuredBy}(x, y), \text{employedBy}(x, z)))).$$

Using Corollary 5.2, which applies $\cong_{S,y}$ equivalence to reduce the entire ‘ OPT ’-pattern to the first part of the Opt -pattern with the filter applied, we can deduce the following equivalent query:

$$\text{DIST}(\text{PROJ}_{xy}(\text{FILTER}_{\top}(\text{Employee}(x) \text{ insuredBy}(x, y)))).$$

In Lemma 7, let $\phi' = \begin{cases} \top, & \text{if } \mathcal{P} = T, \\ \bigwedge_{i=1}^n (\leq_1 P_i, \top), & \text{if } \mathcal{P} = (T P_1(x, z_1) \dots P_i(x, z_i) \dots P_n(x, z_n)). \end{cases}$

Lemma 7. Let $\langle s, \tau_s, \phi_s \rangle \in S$ with $\{(\triangleright_{\tau_s} P), \phi'\} \subseteq \phi_s$ and \mathcal{P} a graph pattern s.t. $\mathcal{P} = T$ or $\mathcal{P} = (T P_1(x, z_1) \dots P_i(x, z_i) \dots P_n(x, z_n))$ for $n \geq i \geq 1$. Then,

$$\text{DIST}(\text{PROJ}_X(\text{JOIN}(\mathcal{P}, P(x, y)))) \equiv_S \text{PROJ}_X(\text{JOIN}(\mathcal{P}, P(x, y))).$$

Proof. For the pattern $\text{PROJ}_X(\text{JOIN}(\mathcal{P}, P(x, y)))$, we can observe that $\text{var}(\mathcal{P}) \cap \text{var}(P(x, y)) = \{x\}$ because: (a) $\mathcal{P} = T$ or $T \blacktriangleleft \mathcal{P}$, and (b) T is either $C(x)$ or $R(x, z)$ or $R^-(x, z)$ based on the target τ_s of a shape $\langle s, \tau_s, \phi_s \rangle \in \mathcal{S}$. Using the definition of Join and Proj, the solution of $\text{PROJ}_X(\text{JOIN}(\mathcal{P}, P(x, y)))$ is all possible combinations of solutions that satisfy \mathcal{P} and $P(x, y)$, and restricted to variable set X . Let G be a graph s.t. $G \models \mathcal{S}$. Then, constraint $(\triangleright_{\tau_s} P) \in \phi_s$ ensures that all P -successors binding to y are unique (i.e., distinct) among the target nodes τ_s over G . Therefore, for each mapping in \mathcal{P}^G , we can only find at most one corresponding mapping in $P(x, y)^G$ when $\mathcal{P} = T$. This implies $|\mu, \text{PROJ}_X(\text{JOIN}(\mathcal{P}, P(x, y)))^G| = 1$ for $\forall \mu \in \text{sup}(\text{PROJ}_X(\text{JOIN}(\mathcal{P}, P(x, y)))^G)$. Hence, $\text{Dist}(\text{PROJ}_X(\text{JOIN}(\mathcal{P}, P(x, y)))) \equiv_S \text{PROJ}_X(\text{JOIN}(\mathcal{P}, P(x, y)))$.

In the case of $\mathcal{P} = \{T P_1(x, z_1) \dots P_i(x, z_i) \dots P_n(x, z_n)\}$ s.t. $n \geq 1$, the constraint component $\bigwedge_{i=1}^n (\leq_1 P_i, \top) \in \phi_s$ guarantees the absence of duplicate mappings $\mu \in \text{PROJ}_X(\text{JOIN}(\mathcal{P}, P(x, y)))^G$ that could have originated from the join operations with the (inner) patterns $P_1(x, z_1) \dots P_i(x, z_i) \dots P_n(x, z_n)$. Thus, equivalence \equiv_S holds.

Example 9. Consider the SPARQL query,

$$\text{DIST}(\text{PROJ}_{yz}(\text{Employee}(x) \text{ hasName}(x, y) \text{ hasID}(x, z)))$$

over a graph $\mathcal{G} \models \langle \text{Employee}, \tau_{\text{Employee}}, \phi_{\text{Employee}} \rangle$ such that $\{(\triangleright_{\text{Employee}} \text{ hasID}), (\leq_1 \text{ hasName}, \top)\} \subseteq \phi_{\text{Employee}}$ (resp., $\mathcal{G} \models \langle \exists \text{hasName}, \tau_{\exists \text{hasName}}, \phi_{\exists \text{hasName}} \rangle$ s.t. $\{(\triangleright_{\text{Employee}} \text{ hasID}), (\leq_1 \text{ hasName}, \top)\} \subseteq \phi_{\exists \text{hasName}}$). Then, ‘Dist’ modifier of the query can be removed by using Lemma 7,

$$\text{PROJ}_{yz}(\text{Employee}(x) \text{ hasName}(x, y) \text{ hasID}(x, z)).$$

Definition 6. Let Q be a SPARQL query, and let \mathcal{P} and \mathcal{U} be two graph patterns. Then, we write $\mathcal{U} \trianglelefteq Q$ if $\text{Dist}(\text{PROJ}_X(\mathcal{P})) \trianglelefteq Q$ and $\mathcal{U} \trianglelefteq \mathcal{P}$.

Theorem 1 specifies query rewriting, outlining the necessary conditions for applying and propagating the established equivalences from Lemmas 1 to 7 to a larger query.

Theorem 1. Let Q be a SPARQL query and \mathcal{S} a SHACL document. Let \mathcal{U} and \mathcal{V} be two graph patterns. Then,

1. $Q \equiv_S Q^{\mathcal{U} \rightarrow \mathcal{V}}$ if $\mathcal{U} \equiv_S \mathcal{V}$
2. $\text{PROJ}_X(Q) \equiv_S \text{PROJ}_X(Q^{\mathcal{U} \rightarrow \mathcal{V}})$ if $\mathcal{U} \trianglelefteq Q$, $\mathcal{U} \equiv_{S,y} \mathcal{V}$ and $y \notin \text{var}(\text{PROJ}_X(Q) \setminus \mathcal{U})$
3. $Q \equiv_S Q^{\mathcal{U} \rightarrow \mathcal{V}}$ if $\mathcal{U} \trianglelefteq Q$, $\mathcal{U} \cong_{S,y} \mathcal{V}$ and $y \notin \text{var}(Q \setminus \mathcal{U})$

Example 10. Consider a query that asks for all distinct combinations of student ids and their corresponding advisors (if available), for those who are currently enrolled in a university program, such as

$$\text{DIST}(\text{PROJ}_{zk}(\text{OPT}_{\top}(\text{Student}(x) \text{ enrolledIn}(x, y) \text{ Program}(y) \text{ hasID}(x, z), \text{ hasAdvisor}(x, k))))$$

over a graph G s.t. $G \models \langle \text{Student}, \tau_{\text{Student}}, \phi_{\text{Student}} \rangle$ with $\{(\geq_1 \text{ enrolledIn}, \text{Program}), (\triangleright_{\text{Student}} \text{ hasID}), (=_1 \text{ hasAdvisor}, \top)\} \subseteq \phi_{\text{Student}}$. Then, using clause 1

of Lemma 1 (i.e., \equiv_S), we can reduce ‘OPT’ to ‘JOIN’ pattern as follows,

$$\text{DIST}(\text{PROJ}_{z,k}(\text{JOIN}(\text{Student}(x) \text{ enrolledIn}(x,y) \text{ Program}(y) \text{ hasID}(x,z) \text{ hasAdvisor}(x,k)))) .$$

Similarly, using clause 2 of Lemma 3 (i.e., $\cong_{S,y}$), we can remove the join clause ‘ $\text{enrolledIn}(x,y) \text{ Program}(y)$ ’ from the query,

$$\text{DIST}(\text{PROJ}_{z,k}(\text{JOIN}(\text{Student}(x) \text{ hasID}(x,z), \text{ hasAdvisor}(x,k)))) .$$

Finally, ‘DIST’ modifier can be removed by using Lemma 7 (i.e., \equiv_S),

$$\text{PROJ}_{z,k}(\text{JOIN}(\text{Student}(x) \text{ hasID}(x,z), \text{ hasAdvisor}(x,k))) .$$

Theorems 2 and 3 state the *confluence* [8, Chapter 6] property, which guarantees that the query rewriting, defined by employing equivalences from Lemmas 1–7, is deterministic and results in the same optimized query, regardless of the order in which the rewriting rules (i.e., replacing left-hand-side of equivalences with the right-hand-side) are applied. In this context, the query rewriting defined by Lemmas 1–6 (resp., 1–7) freely employs equivalence \equiv_S for rewriting graph patterns anywhere within a query, permits the use of equivalence $\cong_{S,y}$ when y is confined within the joint (or optional) triples pattern to be removed (or transformed into a join), and exclusively employs equivalence $\cong_{S,y}$ for rewriting distinct patterns, as described in Theorem 1.

Theorem 2. *Query rewriting defined by Lemmas 1 to 6 is a confluent reduction.*

Proof. The main point of the proof is to establish that all overlapping rewriting rules (i.e., application of equivalences) defined by Lemmas 1 to 6 always yield the same uniquely simplified pattern when applied to a graph pattern. Taking $\mathcal{P}' = P(x,y)$ into account, observe that satisfying the constraint $\phi = (\geq_n P.\beta)$ with $\beta = s'$ in Lemma 3.1 also satisfies $(\geq_n P.\top)$ in Lemma 1.1, and applying either lemma to an optional pattern of the form $\text{OPT}_{\mathcal{F}}(\mathcal{P}, P(x,y))$ where $y \notin \text{var}(\mathcal{P})$ results in the same pattern, namely $\text{FILTER}_{\mathcal{F}}(\text{JOIN}(\mathcal{P}, P(x,y)))$. Similarly, fulfilling $(=_1 P.\beta)$ s.t. $\beta = s'$ in Lemma 4 satisfies $(=_1 P.\top)$ in Lemma 2, and applying either lemma leads to the same optimized join pattern. This property also applies to the implied constraint case, e.g., $(=_1 P.\top) \rightarrow (\geq_1 P.\top)$. Likewise, satisfying ϕ in Lemma 3.2 satisfies the condition $(\geq_n P.\top)$ in Lemma 1.2, and applying either lemma produces the same optimized joint pattern.

Due to the implication $(=_1 P.\top) \rightarrow (\geq_1 P.\top)$, when the condition of Lemma 2 is fulfilled, the rewriting of Lemma 1.2 also applies. As Lemma 2 can be applied to all patterns (Clause 1 of Theorem 2), applying both rewritings to a distinct query pattern leads to the same optimized pattern. A similar situation arises for Lemmas 2.2 and 4 due to the implication $(=_1 P.\beta) \rightarrow (\geq_1 P.\beta)$, and they also lead to the same optimized pattern.

Theorem 3. *Query rewriting defined by Lemmas 1 to 7 is a confluent reduction iff*

$$\phi' = \begin{cases} \top, & \text{if } \mathcal{P} = T, \\ \wedge_{i=1}^n (_1 P_i. \top), & \text{if } \mathcal{P} = (T P_1(x, z_1) \dots P_i(x, z_i) \dots P_n(x, z_n)) \text{ in Lemma 7.} \end{cases}$$

The rewriting employing equivalence supported by Lemma 7 transforms a distinct pattern into a distinct-free pattern. However, certain rewritings, such as Lemmas 1.2, 3.2 and 6.2, are exclusive to distinct patterns. Hence, the introduction of constraint $\wedge_{i=1}^n (=_1 P_i, \top)$ in Theorem 3 becomes necessary to ensure the application of complementary rewritings, particularly based on equivalences supported by Lemmas 2 and 4, when the aforementioned rewritings cannot be directly applied after Lemma 7.

Example 11. Consider a SPARQL query,

$$\text{DIST}(\text{PROJ}_{xy}(\text{employeeID}(x, y) \text{ hiredBy}(x, k) \text{ insuredBy}(x, z)))$$

over a graph G s.t. $G \models \langle \exists \text{employeeID}, \tau_{\exists \text{employeeID}}, \phi_{\exists \text{employeeID}} \rangle$ with $\{(\triangleright_{\exists \text{employeeID}} \text{employeeID}), (=_1 \text{insuredBy}, \top), (\text{hiredBy} = \text{insuredBy})\} \subseteq \phi_{\exists \text{employeeID}}$. Subsequently, the constraint $(=_1 \text{hiredBy}, \top)$ can be inferred from the $(=_1 \text{insuredBy}, \top)$ and $(\text{hiredBy} = \text{insuredBy})$. Based on Theorem 1, the query is subject to optimization using the Lemmas 1.2 (i.e., $\cong_{S,y}$), 2 (i.e., $\equiv_{S,y}$), 6.2 (i.e., $\cong_{S,y}$) and 7 (i.e., \equiv_S). Then, there are seven possible rewriting combinations that can be applied using the lemmas mentioned above: (Lemmas (1.2, 1.2, 7), Lemmas (1.2, 2, 7), Lemmas (2, 1.2, 7), Lemmas (2, 2, 7), Lemmas (6.2, 1.2, 7), Lemmas (6.2, 2, 7), and Lemmas (7, 2, 2)). All these combinations result in a uniquely optimized query, which is

$$\text{PROJ}_{xy}(\text{employeeID}(x, y)).$$

Note that Lemmas 1.2 and 6.2 cannot be applied after Lemma 7. This information is further useful in determining the cheapest order in complex query optimization, which is a topic related to our future goals.

4 Discussion

We have presented a set of query rewriting rules that exploit the relation between the SPARQL queries containing graph patterns and SHACL constraints describing these graph patterns in the RDF. By treating each lemma as a set of rules under distinct constraints, we emphasized the significance of each constraint-enabled rule and made it easier to track and analyze their order and individual contributions to the overall robust query rewriting process, see Example 10. They also offer a clearer understanding of the different orders and strategies that can be employed while rewriting complex queries, see Example 11.

Constraints, known for restricting data to useful relations [23], can be used to rewrite queries into more efficient equivalents [50]. Leveraging constraints for optimizing queries has become a well-established practice in knowledge graphs [13, 20, 56], deductive databases [14], relational databases [3, 25] and other fields [48]. Regarding SHACL, Rabbani et al. [47] were the first to propose shapes as global statistics for estimating cardinality and optimizing SPARQL query plans for the RDF graph under query. Further, Abbas et al. [1] studied the containment of restricted classes of SPARQL queries under set semantics, while considering ShEx constraints [46]. They also leveraged information obtained from ShEx schemas to optimize the execution of SPARQL

queries [2]. In their work [2] (resp., [1]), a notion of well-formed ShEx schemas (resp., query patterns), especially suitable for query optimization (resp., containment), was first introduced, then subsequently used to reorganise execution orders (resp., to determine subsumption) of triple patterns in a SPARQL query. In contrast to previous studies [1, 2, 47], we have prioritized optimizing a much larger class of queries, and subsequently, introduced an explicit set of core rewriting rules for optimizing SPARQL queries under bag semantics when SHACL constraints are satisfied, and proved their correctness. We have further shown that our proposed query rewriting rules are confluent and always lead to a unique simplified form. We plan to investigate the potential benefits of integrating these rewriting rules into a query processing engine to enhance query execution in the near future. Additionally, we intend to extend our results to cover the broader expressivity and more general cases of SPARQL queries.

5 Conclusion

In this paper, we have proposed a set of optimizations for SPARQL query constructed from core operators of SPARQL 1.1 under the SHACL constraints. We believe that the proposed optimizations could play a crucial role in simplifying and answering SPARQL queries over large-scale RDF triples when SHACL descriptions of the underlying RDF dataset are available. In future, we aim to investigate an optimal normal form of SPARQL 1.1 [24] queries studied in [12, 30] (which includes nesting, assignment, construct [35, 44] and aggregation as in online analytical processing cube and window-based queries) with respect to (full) SHACL constraints [32] formalized in [19, 41, 42] that should allow us to introduce a confluent set of query rewriting [50, Sect. 5] rules.

Acknowledgements. This work is supported by the Norwegian Research Council via the SIR-IUS SFI (237898). We thank Egor Kostylev for many fruitful discussions.

References

1. Abbas, A., Genevès, P., Roisin, C., Layaïda, N.: SPARQL query containment with ShEx constraints. In: Kirikova, M., Nørvåg, K., Papadopoulos, G.A. (eds.) ADBIS 2017. LNCS, vol. 10509, pp. 343–356. Springer, Cham (2017). https://doi.org/10.1007/978-3-319-66917-5_23
2. Abbas, A., Genevès, P., Roisin, C., Layaïda, N.: Selectivity estimation for SPARQL triple patterns with shape expressions. In: Mikkonen, T., Klamma, R., Hernández, J. (eds.) ICWE 2018. LNCS, vol. 10845, pp. 195–209. Springer, Cham (2018). https://doi.org/10.1007/978-3-319-91662-0_15
3. Abiteboul, S., Hull, R., Vianu, V.: Foundations of Databases, vol. 8. Addison-Wesley, Reading (1995)
4. Ahmetaj, S., David, R., Polleres, A., Šimkus, M.: Repairing SHACL constraint violations using answer set programming. In: Sattler, U., et al. (eds.) ISWC 2022. LNCS, vol. 13489, pp. 375–391. Springer, Cham (2022). https://doi.org/10.1007/978-3-031-19433-7_22
5. Angles, R., Gutierrez, C.: The multiset semantics of SPARQL patterns. In: Groth, P., et al. (eds.) ISWC 2016. LNCS, vol. 9981, pp. 20–36. Springer, Cham (2016). https://doi.org/10.1007/978-3-319-46523-4_2

6. Arias, M., Fernández, J.D., Martínez-Prieto, M.A., de la Fuente, P.: An empirical study of real-world SPARQL queries. *arXiv preprint*<https://arxiv.org/pdf/1103.5043.pdf> (2011)
7. Atre, M.: For the DISTINCT clause of SPARQL queries. In: Proceedings of the 25th International Conference Companion on World Wide Web, pp. 7–8 (2016)
8. Baader, F., Nipkow, T.: Term Rewriting and All That. Cambridge University Press, Cambridge (1998)
9. Bilidas, D., Koubarakis, M.: Efficient duplicate elimination in SPARQL to SQL translation. In: Description Logics (2018)
10. Bogaerts, B., Jakubowski, M., Van den Bussche, J.: Expressiveness of SHACL features and extensions for full equality and disjointness tests. *arXiv preprint*<https://arxiv.org/pdf/2212.03553.pdf> (2022)
11. Broekstra, J., Kampman, A., Van Harmelen, F.: Sesame: an architecture for storing and querying RDF data and schema information. In: Spinning the Semantic Web, pp. 197–222 (2001)
12. Buil-Aranda, C., Arenas, M., Corcho, O., Polleres, A.: Federating queries in SPARQL 1.1: syntax, semantics and evaluation. *J. Web Semant.* **18**(1), 1–17 (2013)
13. Calvanese, D., et al.: Ontop: answering SPARQL queries over relational databases. *Semantic Web* **8**(3), 471–487 (2017)
14. Chakravarthy, U.S., Grant, J., Minker, J.: Foundations of semantic query optimization for deductive databases. In: Foundations of Deductive Databases and Logic Programming, pp. 243–273. Elsevier (1988)
15. Cheng, S., Hartig, O.: OPT+: a monotonic alternative to OPTIONAL in SPARQL. *J. Web Eng.* (2019)
16. W3C Consortium. Apache Jena (2016). <https://jena.apache.org>
17. Corman, J., Florenzano, F., Reutter, J.L., Savković, O.: SHACL2SPARQL: validating a SPARQL endpoint against recursive SHACL constraints. In: ISWC (Satellites), pp. 165–168 (2019)
18. Corman, J., Florenzano, F., Reutter, J.L., Savković, O.: Validating SHACL constraints over a SPARQL endpoint. In: Ghidini, C., et al. (eds.) ISWC 2019. LNCS, vol. 11778, pp. 145–163. Springer, Cham (2019). https://doi.org/10.1007/978-3-030-30793-6_9
19. Corman, J., Reutter, J.L., Savković, O.: Semantics and validation of recursive SHACL. In: Vrandečić, D., et al. (eds.) ISWC 2018. LNCS, vol. 11136, pp. 318–336. Springer, Cham (2018). https://doi.org/10.1007/978-3-030-00671-6_19
20. Di Pinto, F., et al.: Optimizing query rewriting in ontology-based data access. In: Proceedings of the 16th International Conference on Extending Database Technology, pp. 561–572 (2013)
21. Erling, O.: Implementing a SPARQL compliant RDF triple store using a SQLORDBMS. OpenLink Software Virtuoso (2001)
22. Inc Franz. AllegroGraph (2017). *Compatible Semantic Technologies*, <https://allegrograph.com>
23. Grefen, P.W.P.J., Apers, P.M.G.: Integrity control in relational database systems—an overview. *Data Knowl. Eng.* **10**(2), 187–223 (1993)
24. Harris, S., Seaborne, A., Prud'hommeaux, E.: SPARQL 1.1 query language. W3C Recommendation **21**(10), 778 (2013)
25. Jarke, M., Koch, J.: Query optimization in database systems. *ACM Comput. Surv. (CsUR)* **16**(2), 111–152 (1984)
26. Jian, X., Wang, Y., Lei, X., Zheng, L., Chen, L.: SPARQL rewriting: towards desired results. In: Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data, pp. 1979–1993 (2020)
27. Johnson, D.S., Klug, A.: Testing containment of conjunctive queries under functional and inclusion dependencies. In: Proceedings of the 1st ACM SIGACT-SIGMOD Symposium on Principles of Database Systems, pp. 164–169 (1982)

28. Jung, H., Kim, W.: Automated conversion from natural language query to SPARQL query. *J. Intell. Inf. Syst.* **55**(3), 501–520 (2020). <https://doi.org/10.1007/s10844-019-00589-2>
29. Kaminski, M., Kostylev, E.V.: Beyond well-designed SPARQL. In: 19th International Conference on Database Theory (ICDT 2016). Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik (2016)
30. Kaminski, M., Kostylev, E.V., Grau, B.C.: Query nesting, assignment, and aggregation in SPARQL 1.1. *ACM Trans. Database Syst. (TODS)* **42**(3), 1–46 (2017)
31. Kaufmann, E., Bernstein, A., Zumstein, R.: Querix: a natural language interface to query ontologies based on clarification dialogs. In: 5th International Semantic Web Conference (ISWC 2006), pp. 980–981. Citeseer (2006)
32. Knublauch, H., Kontokostas, D.: Shapes constraint language (SHACL). *W3C Candidate Recommend.* **11**(8) (2017)
33. Kontchakov, R., Rezk, M., Rodríguez-Muro, M., Xiao, G., Zakharyaschev, M.: Answering SPARQL queries over databases under OWL 2 QL entailment regime. In: Mika, P., et al. (eds.) ISWC 2014. LNCS, vol. 8796, pp. 552–567. Springer, Cham (2014). https://doi.org/10.1007/978-3-319-11964-9_35
34. Kostylev, E.V., Reutter, J.L., Romero, M., Vrgoč, D.: SPARQL with property paths. In: Arenas, M., et al. (eds.) ISWC 2015. LNCS, vol. 9366, pp. 3–18. Springer, Cham (2015). https://doi.org/10.1007/978-3-319-25007-6_1
35. Kostylev, E.V., Reutter, J.L., Ugarte, M.: CONSTRUCT queries in SPARQL. In: 18th International Conference on Database Theory (ICDT 2015). Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik (2015)
36. Le, W., Duan, S., Kementsietsidis, A., Li, F., Wang, M.: Rewriting queries on SPARQL views. In: Proceedings of the 20th International Conference on World Wide Web, pp. 655–664 (2011)
37. Manola, F., Miller, E., McBride, B., et al.: RDF primer. *W3C Recommend.* **10**(1–107), 6 (2004)
38. Mohamed, A., Abuoda, G., Ghanem, A., Kaoudi, Z., Aboulnaga, A.: RDFFrames: knowledge graph access for machine learning tools. *VLDB J.* **31**(2), 321–346 (2022)
39. Pandit, H.J., O’Sullivan, D., Lewis, D.: Test-driven approach towards GDPR compliance. In: Acosta, M., Cudré-Mauroux, P., Maleshkova, M., Pellegrini, T., Sack, H., Sure-Vetter, Y. (eds.) SEMANTiCS 2019. LNCS, vol. 11702, pp. 19–33. Springer, Cham (2019). https://doi.org/10.1007/978-3-030-33220-4_2
40. Paretí, P., Konstantinidis, G.: A review of SHACL: from data validation to schema reasoning for RDF graphs. *Reasoning Web International Summer School*, pp. 115–144 (2021)
41. Paretí, P., Konstantinidis, G., Mogavero, F.: Satisfiability and containment of recursive SHACL. *J. Web Semantics* **74**, 100721 (2022)
42. Paretí, P., Konstantinidis, G., Mogavero, F., Norman, T.J.: SHACL satisfiability and containment. In: Pan, J.Z., et al. (eds.) ISWC 2020. LNCS, vol. 12506, pp. 474–493. Springer, Cham (2020). https://doi.org/10.1007/978-3-030-62419-4_27
43. Pérez, J., Arenas, M., Gutierrez, C.: Semantics and complexity of SPARQL. *ACM Trans. Database Syst. (TODS)* **34**(3), 1–45 (2009)
44. Polleres, A., Reutter, J., Kostylev, E.: Nested constructs vs. sub-selects in SPARQL. In: CEUR Workshop Proceedings (2016)
45. Prud’Hommeaux, E., Seaborne, A.: SPARQL query language for RDF. *W3C Recommendation*, January 15, 2008 (2011)
46. Prud’hommeaux, E., Gayo, J.E.L., Solbrig, H.: Shape expressions: an RDF validation and transformation language. In: Proceedings of the 10th International Conference on Semantic Systems, pp. 32–40 (2014)
47. Rabbani, K., Lissandrini, M., Hose, K.: Optimizing SPARQL queries using shape statistics (2021)

48. Reiter, R.: Nonmonotonic reasoning. In: Exploring Artificial Intelligence, pp. 439–481. Elsevier (1988)
49. Rodriguez-Muro, M., Rezk, M.: Efficient SPARQL-to-SQL with R2RML mappings. *J. Web Semantics* **33**, 141–169 (2015)
50. Schmidt, M., Meier, M., Lausen, G.: Foundations of SPARQL query optimization. In: Proceedings of the 13th International Conference on Database Theory, pp. 4–33 (2010)
51. Sequeda, J.F., Miranker, D.P.: Ultrawrap: SPARQL execution on relational data. *J. Web Semantics* **22**, 19–39 (2013)
52. Taelman, R., Van Herwegen, J., Vander Sande, M., Verborgh, R.: Comunica: a modular SPARQL query engine for the web. In: Vrandečić, D., et al. (eds.) ISWC 2018. LNCS, vol. 11137, pp. 239–255. Springer, Cham (2018). https://doi.org/10.1007/978-3-030-00668-6_15
53. Thapa, R.B., Giese, M.: A source-to-target constraint rewriting for direct mapping. In: Hotho, A., et al. (eds.) ISWC 2021. LNCS, vol. 12922, pp. 21–38. Springer, Cham (2021). https://doi.org/10.1007/978-3-030-88361-4_2
54. Thapa, R.B., Giese, M.: Mapping relational database constraints to SHACL. In: The Semantic Web-ISWC 2022: 21st International Semantic Web Conference, Virtual Event, October 23–27, 2022, Proceedings, pp. 214–230 (2022)
55. Thapa, R.B., Giese, M.: Optimizing SPARQL queries with SHACL (extended version). Research Report 504, Dept. of Informatics, University of Oslo, July 2023
56. Xiao, G., Kontchakov, R., Cogrel, B., Calvanese, D., Botoeva, E.: Efficient handling of SPARQL OPTIONAL for OBDA. In: Vrandečić, D., et al. (eds.) ISWC 2018. LNCS, vol. 11136, pp. 354–373. Springer, Cham (2018). https://doi.org/10.1007/978-3-030-00671-6_21
57. Yin, X., Gromann, D., Rudolph, S.: Neural machine translating from natural language to SPARQL. *Futur. Gener. Comput. Syst.* **117**, 510–519 (2021)