



AALBORG UNIVERSITY
DENMARK

Aalborg Universitet

Extraction of Validating Shapes from very large Knowledge Graphs

Rabbani, Kashif; Lissandrini, Matteo; Hose, Katja

Published in:
Proceedings of the VLDB Endowment

DOI (link to publication from Publisher):
[10.14778/3579075.3579078](https://doi.org/10.14778/3579075.3579078)

Creative Commons License
CC BY-NC-ND 4.0

Publication date:
2023

Document Version
Publisher's PDF, also known as Version of record

[Link to publication from Aalborg University](#)

Citation for published version (APA):
Rabbani, K., Lissandrini, M., & Hose, K. (2023). Extraction of Validating Shapes from very large Knowledge Graphs. *Proceedings of the VLDB Endowment*, 16(5), 1023-1032. <https://doi.org/10.14778/3579075.3579078>

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal -

Take down policy

If you believe that this document breaches copyright please contact us at vbn@aub.aau.dk providing details, and we will remove access to the work immediately and investigate your claim.

Extraction of Validating Shapes from very large Knowledge Graphs

Kashif Rabbani

Aalborg University, Denmark
kashifrabbani@cs.aau.dk

Matteo Lissandrini

Aalborg University, Denmark
matteo@cs.aau.dk

Katja Hose

Aalborg University, Denmark
khose@cs.aau.dk

ABSTRACT

Knowledge Graphs (KGs) represent heterogeneous domain knowledge on the Web and within organizations. There exist shapes constraint languages to define *validating shapes* to ensure the quality of the data in KGs. Existing techniques to extract validating shapes often fail to extract complete shapes, are not scalable, and are prone to produce spurious shapes. To address these shortcomings, we propose the **QUALITY SHAPES EXTRACTION (QSE)** approach to extract validating shapes in very large graphs, for which we devise both an exact and an approximate solution. QSE provides information about the reliability of shape constraints by computing their confidence and support within a KG and in doing so allows to identify shapes that are most informative and less likely to be affected by incomplete or incorrect data. To the best of our knowledge, QSE is the first approach to extract a complete set of validating shapes from WikiData. Moreover, QSE provides a 12x reduction in extraction time compared to existing approaches, while managing to filter out up to 93% of the invalid and spurious shapes, resulting in a reduction of up to 2 orders of magnitude in the number of constraints presented to the user, e.g., from 11,916 to 809 on DBpedia.

PVLDB Reference Format:

Kashif Rabbani, Matteo Lissandrini, and Katja Hose. Extraction of Validating Shapes from very large Knowledge Graphs. PVLDB, 16(5): 1023-1032, 2023.

doi:10.14778/3579075.3579078

PVLDB Artifact Availability:

The source code, data, and/or other artifacts have been made available at <https://github.com/dkw-aau/qse>.

1 INTRODUCTION

Knowledge Graphs (KGs), stored as collections of triples of the form $\langle \text{subject}, \text{relation}, \text{object} \rangle$ using the Resource Description Framework (RDF) [10], are in widespread use both within companies [29, 43, 44] and on the Web [46, 49]. Nonetheless, as KGs quickly accrue more data, practical applications impose further demands in terms of quality assessment and validation [34, 38, 52]. Hence, shapes constraint languages, e.g., SHACL [23], and ShEx [35], have been proposed to validate KGs by enforcing constraints represented in the form of *validating shapes*. For instance, we can express that an entity of type Student requires a name, a registration number,

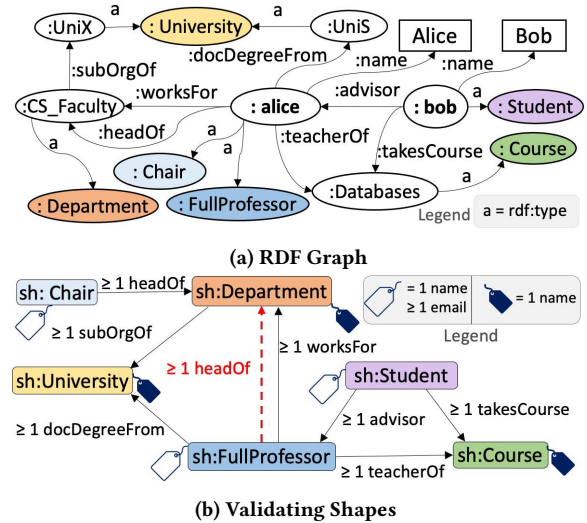


Figure 1: An example RDF Graph and Validating Shapes

and should be enrolled in some courses; and that these attributes should be of type string, integer, and Course, respectively – see Figures 1a and 1b for an example KG and corresponding shapes.

Often, validating shapes are manually specified by domain experts. Yet, when trying to specify validating shapes for already-existing large-scale KGs, data scientists are in need of tools that can speed up this process [38]. Thus, various tools have been proposed to automatically [9, 12, 20, 26] or semi-automatically [5, 32, 36] produce a set of validating shapes for a target KG. Unfortunately, these methods suffer from 3 important limitations: (1) they are not able to produce complete shapes, e.g., they can identify that a student should have a property of type takesCourse but they do not extract the fact that the object should be of type Course; (2) the shapes they produce are easily affected by errors and inconsistencies in the KG, e.g., if some departments, by mistake, are attached the property hasAdvisor, a corresponding *spurious shape* is extracted; and (3) they do not scale to large KGs, e.g., they cannot process the full English WikiData, and they take days to process a subset of it. Therefore, in this work, we present the first techniques for *efficient extraction of validating shapes from very large existing KGs that also ensures robustness against the effects of spuriousness*.

Spuriousness poses important challenges to automatic shape extraction methods. For instance, in DBpedia [4], some of the entities representing musical bands are wrongly assigned to the class dbo:City. As a consequence, when shapes are extracted from its instance data using existing approaches, the resulting node shape for dbo:City specifies that cities are allowed optional properties like dbo:genre and dbo:formerBandMember. Hence, due to the effect

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit <https://creativecommons.org/licenses/by-nc-nd/4.0/> to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing info@vldb.org. Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.

Proceedings of the VLDB Endowment, Vol. 16, No. 5 ISSN 2150-8097.
doi:10.14778/3579075.3579078

of *spuriousness*, existing approaches generate tens of thousands of shapes (our experiments show that standard extraction processes produce more than 2 million property shapes for WikiData [49]). Thus, it becomes unmanageable for domain experts to manually identify valid shapes. SheXer [12] is the only existing approach that attempts to tackle this issue by filtering shapes based on a “trustworthiness” score. Unfortunately, this score does not directly translate into how frequently a shape is satisfied in a dataset, so it is still prone to generate spurious shapes, and it is also hard to tune. Furthermore, SheXer is not able to efficiently process large KGs.

Therefore, to tackle the issue of *spuriousness*, we study and formalize the problem of *support-based shapes extraction* and propose the QUALITY SHAPES EXTRACTION (QSE) approach as a solution to this problem. To tackle the issue of *scalability*, we devise two efficient algorithms, QSE-Exact and QSE-Approximate. Hence, QSE can filter out shapes affected by spurious or erroneous data based on robust and easily understandable measures. QSE allows shapes extraction both from KGs available as files as well as SPARQL endpoints. Moreover, our efficient approximation algorithm enables shape extraction even on a commodity machine by sampling the KG entities via a dynamic multi-tiered reservoir sampling technique.

We perform a thorough experimental evaluation using both synthetic (LUBM [18]) and real (DBpedia [4], YAGO-4 [46], WikiData [49]) KGs demonstrating the benefits of our approach. The shapes produced by our approach are of high quality and instrumental for easily finding errors in real KGs. The results show that QSE-Exact can extract shapes from the entire WikiData’s 2015 dump in 16 minutes and from 2021’s dump (1.9B triples) in 2.5 hours. Similarly, QSE-Approximate can extract shapes from WikiData’s 2021 dump in 90 minutes on a 32GB machine while still achieving 100% precision and 95% recall in the set of shapes produced. Hence, our sampling strategy is accurate and efficient both when extracting shapes from a file as well as when using an endpoint.

2 RDF SHAPES AND THE QSE PROBLEM

In the following, we first introduce the KG data model and the concepts of validating shapes, their support, and confidence, then we define our focus: the QUALITY SHAPES EXTRACTION problem.

2.1 Preliminaries

The standard model for encoding KGs is the Resource Description Framework (RDF [10]), which describes data as a set of $\langle s, p, o \rangle$ triples stating that a subject s is in a relationship with an object o through predicate p . Therefore, we define an RDF graph as follows:

Definition 2.1 (RDF graph). Given pairwise disjoint sets of IRIs \mathcal{I} , blank nodes \mathcal{B} , and literals \mathcal{L} , an RDF Graph $\mathcal{G}:\langle N, E \rangle$ is a graph with a finite set of nodes $N \subset (\mathcal{I} \cup \mathcal{B} \cup \mathcal{L})$ and a finite set of edges $E \subset \{ \langle s, p, o \rangle \in (\mathcal{I} \cup \mathcal{B}) \times \mathcal{I} \times (\mathcal{I} \cup \mathcal{B} \cup \mathcal{L}) \}$.

Moreover, we distinguish two special subsets of the IRIs \mathcal{I} : predicates \mathcal{P} and classes \mathcal{C} . The set of predicates $\mathcal{P} \subset \mathcal{I}$ is the subset of IRIs that appear in the predicate position p in any $\langle s, p, o \rangle \in \mathcal{G}$. Among predicates \mathcal{P} , we identify the type predicate $a \in \mathcal{P}$, which corresponds to IRI `rdf:type` [51] or `wdt:P31` WikiData [49], as the predicate that connects all entities that are instances of a class to the node representing the class itself, i.e., their type. Thus, all the IRIs that are classes in \mathcal{G} form the subset $\mathcal{C}:\{c \in \mathcal{I} \mid \exists s \in \mathcal{I} \text{ s.t. } \langle s, a, c \rangle \in \mathcal{G}\}$.

Given a KG \mathcal{G} , a set of *validating shapes* represents integrity constraints in the form of a shape schema \mathcal{S} over \mathcal{G} . Since the shape schema describes shapes associated with node types and their connections to other attributes and node types, we can also visualize the shape schema \mathcal{S} as a particular type of graph (see Figures 1a and 1b). Therefore, in the following, we refer to two concepts: the *data graph* \mathcal{G} and the *shape graph* derived from \mathcal{S} . The *data graph* is the RDF graph \mathcal{G} to be validated, while the *shape graph* consists of constraints in the form of the shape schema \mathcal{S} against which entities of the data graph are validated. These constraints are defined using node and property shapes. In the following, we adopt the previously defined syntax [42] to refer to the set \mathcal{S} according to the SHACL *core constraint components* [50]. Finally, while validating shapes can also be expressed in ShEx [34], our approach can be trivially extended to output ShEx directly, or it can exploit existing SHACL to ShEx converters [53]. Thus, without loss of generality, we focus on the current standard for SHACL shapes in the following.

Definition 2.2 (Shape Schema). A SHACL shape schema consists of a set of node shapes \mathcal{S} , with $\langle s, \tau_s, \Phi_s \rangle \in \mathcal{S}$, where s is the shape name, $\tau_s \in \mathcal{C}$ is the *target class*, and Φ_s is a set of property shapes of the form $\phi_s:\langle \tau_p, T_p, C_p \rangle$, where $\tau_p \in \mathcal{P}$ is called the target property, $T_p \subset \mathcal{I}$ contains either an IRI defining a *literal type*, e.g., `xsd:string`, or a set of IRIs – called *class type constraint*, and C_p is a pair $(n, m) \in \mathbb{N} \times (\mathbb{N} \cup \{\infty\})$. $n \leq m$ – called *min and max cardinality constraints*.

Therefore, given a node shape $s \in \mathcal{S}$ for the *target class* $\tau_s \in \mathcal{C}$, Φ_s defines which properties each instance of τ_s can or should be associated with. For instance, the shape $\langle \text{sh:Student}, \text{:Student}, \{ \phi_{s_1}, \phi_{s_2} \} \rangle$ from Figure 1b, contains a node shape for target class `:Student` and enforces two property shapes ϕ_1 and ϕ_2 . The property shape ϕ_1 has a target property $\tau_p = \text{:name}$, a literal type constraint $T_p = \text{xsd:string}$, and the cardinality constraints $C_p = (1, 1)$. Similarly, the property shape ϕ_2 has a target property $\tau_p = \text{:takesCourse}$, a class type constraint $T_p = \text{:Course}$, and the cardinality constraint $C_p = (1, \infty)$.

When validating a graph \mathcal{G} against a shape schema \mathcal{S} having a node shape $\langle s, \tau_s, \Phi_s \rangle \in \mathcal{S}$, we verify that each entity $e \in \mathcal{G}$ that is an instance of τ_s satisfies all the constraints Φ_s . Note that we use the term entity and node interchangeably throughout the paper. Thus, we define the semantics of \mathcal{S} as follows:

Definition 2.3 (Validating Shape Semantics). Given a node shape $\langle s, \tau_s, \Phi_s \rangle \in \mathcal{S}$, a graph \mathcal{G} , and an entity e s.t. $\langle e, a, \tau_s \rangle \in \mathcal{G}$, we have that s validates e , and we write $e \models_{\mathcal{G}} \phi$, if for every property shape $\phi_s:\langle \tau_p, T_p, C_p \rangle \in \Phi_s$ the following conditions hold:

- If T_p is a literal type constraint, then for every triple $(e, \tau_p, l) \in \mathcal{G}$, l is a literal of type T_p .
- If T_p is a set of class type constraints $T_p = \{t_1, t_2, \dots, t_n\}$, then for every triple $(e, \tau_p, o) \in \mathcal{G}$, it holds that $\forall t \in T_p$, o is an instance of t (or of a subclass of t) and if $\exists S_t \in \mathcal{S}$, $o \models_{\mathcal{G}} S_t$.
- $n \leq |\{ \langle s, p, o \rangle \in \mathcal{G} : s = e \wedge p = \tau_p \}| \leq m$, where $C_p = (n, m)$.

Here we study the case where \mathcal{G} is given, and we want to extract the set of validating shapes \mathcal{S} that validates every class in \mathcal{C} from \mathcal{G} . This is the *shapes extraction* problem. In this case, existing automatic approaches [38] assume the graph to be correct, then iterate over all entities in it, and extract for each entity e all necessary shapes that validate e . The union of all such shapes is assumed to be the final schema \mathcal{S} . This is useful when we want to validate new data

that will be added *in the future* to the KG so that it will conform to the data already in the graph. Unfortunately, this approach will produce spurious shapes. For instance, in Figure 1, since `:alice` has both type `Full Professor` and `Chair`, when parsing the triple (`:alice`, `:headOf`, `:CS_Faculty`), the property shape `headOf` (the red dotted arrow in Figure 1b) is assigned to both node shapes, instead of assigning it to the `Chair` node shape only.

2.2 Shapes Support and Confidence

To contrast the effect of spuriousness, we want to exploit statistics on how often properties are applied to entities of a given type. Therefore, we introduce the notion of *support* and *confidence* for shape constraints to study the reliability of extracted shapes. These concepts are inspired by the well-known theory developed for the task of frequent patterns mining [19] and the concept of MNI support for graph patterns [7]. The MNI support of a graph pattern is the minimum cardinality of the *set of all nodes of \mathcal{G} that are mapped to a specific pattern node by some isomorphism* across all the nodes of the pattern. In our approach, a property shape corresponds to a node- and edge-labeled graph pattern. Thus, given the shape $s:\langle s, \tau_s, \Phi_s \rangle \in \mathcal{S}$ its support is the number of entities that are of type τ_s , while the support of a property shape $\phi_s:\langle \tau_p, T_p, C_p \rangle \in \Phi_s$ is the cardinality of entities conforming to it.

Definition 2.4 (Support of ϕ_s). Given a shape $\langle s, \tau_s, \Phi_s \rangle \in \mathcal{S}$ with shape constraint $\phi_s:\langle \tau_p, T_p, C_p \rangle \in \Phi_s$, the support of ϕ_s is defined as the number of entities e satisfying ϕ_s , denoted as $e \models \phi_s$, hence:

$$\text{supp}(\phi_s) = |\{e \in \mathcal{I} \mid e \models \phi_s\}| \quad (1)$$

Finally, the confidence of a constraint ϕ_s measures the ratio between how many entities conform to ϕ_s and the total number of entities that are instances of the target class of the shape s .

Definition 2.5 (Confidence of ϕ_s). Given a shape $\langle s, \tau_s, \Phi_s \rangle \in \mathcal{S}$ having shape constraint $\phi_s:\langle \tau_p, T_p, C_p \rangle \in \Phi_s$, the confidence of ϕ_s is defined as the proportion of entities for which $e \models \phi_s$ among the entities that are instances of the target class τ_s of $s \in \mathcal{S}$, hence:

$$\text{conf}(\phi_s) = \frac{\text{supp}(\phi_s)}{|\{e \mid (e, \text{type}, \tau_s) \in \mathcal{G}\}|} \quad (2)$$

As it happens in the case of frequent pattern mining [19], when extracting validating shapes, the support provides insights on how frequently a constraint is matched in the graph, i.e., the number of entities e satisfying a constraint ϕ_s . While similar to the task of itemset mining [6], the confidence can tell us how strong is the association between a node type and a specific constraint, i.e., the proportion of entities e satisfying a constraint ϕ_s among all the entities that are instances of the node type τ_s of $s \in \mathcal{S}$. For instance, the confidence for property shape `headOf` (Figure 1b) in our snapshot of LUBM is 10% for the `Full Professor` node shape and 100% for `Chair`, which indicates a strong association of the `headOf` property shape to latter and a weak association to the former.

2.3 The Quality Shapes Extraction Problem

Given the need to extract shapes from a large existing graph \mathcal{G} while limiting the effect of spuriousness, we formally define the problem of extracting high-quality shapes from KGs as follows:

PROBLEM 1 (QUALITY SHAPES EXTRACTION). *Given an RDF graph \mathcal{G} , a threshold ω for support, and ϵ for confidence, the problem of quality shapes extraction over \mathcal{G} is to find the set of shapes \mathcal{S} such that for all node shapes $\langle s, \tau_s, \Phi_s \rangle \in \mathcal{S}$ it holds that $\text{supp}(s) > \omega$ and for all property shapes $\phi_s:\langle \tau_p, T_p, C_p \rangle \in \Phi_s$, $\text{supp}(\phi_s) > \omega$ and $\text{conf}(\phi_s) > \epsilon$.*

In the following, we provide both, an exact and an approximate solution to the problem of quality shape extraction.

3 QSE-EXACT

Extracting shapes \mathcal{S} from an RDF graph \mathcal{G} requires processing its triples and analyzing the types of nodes involved both as subjects and objects in those triples. At a high level, we need to know for each entity all its types, these will become node shapes, and then for each entity type, identify property shapes, which requires, in turn, knowing the types of the objects as well. Furthermore, we need to keep frequency counts to know how often a specific property connects nodes of two given types compared to how many entities exist of those types. In our solution, this is done in four steps: (1) entity extraction, (2) entity constraints extraction, (3) support and Confidence computation, and (4) shapes extraction. Here we first consider the case where the graph is stored as a complete dump on a single file. Later, we also consider the case for a graph stored within a triplestore [41] for which the KG is not available as a file.

QSE-Exact (file-based). One of the most common ways to store an RDF graph \mathcal{G} on a file F is to represent it as a sequence of triples. Therefore, QSE reads F line by line and processes it as a stream of $\langle s, p, o \rangle$ triples. Algorithm 1 and Figure 2 present the four main steps of QSE to extract shapes for graph \mathcal{G} stored in F . In the entity extraction phase, the algorithm parses each $\langle s, p, o \rangle$ triple containing a type declaration (e.g., `rdf:type` or `wdt:P31` – this can be configured) and for each entity, it stores the set of its entity types and the global count of their frequencies, i.e., the number of instances for each class (Lines 4-8) in maps Ψ_{ETD} (Entity-to-Data) and Ψ_{CEC} (Class-to-Entity-Count), respectively. For example, Figure 2 (phase 1) presents two example entities `:bob` and `:alice` (from the example graph of Figure 1a) having entity types `:Student`, `:FullProfessor`, and `:Chair`, respectively. Figure 2 also presents the structure of the Entity-to-Data Ψ_{ETD} dictionary map to help understand the captured entities and their information. In the second phase, i.e., entity constraints extraction, the algorithm performs a second pass over F (Lines 9-19) to collect the constraints and the meta-data required to compute support and confidence of each candidate property shape. Specifically, it parses all triples except triples containing type declarations (which can be skipped now) to obtain for each predicate the subject and object types from the map Ψ_{ETD} that was populated in the previous step. The type of a literal object is inferred from the value, and for a non-literal object is obtained from Ψ_{ETD} (Lines 11-16). For example, Ψ_{ETD} records that the types of `:alice` are `:FullProfessor` and `:Chair`. Then, the Entity-to-Property-Data map Ψ_{ETPD} is updated to add the candidate property constraints associated with each subject entity (Line 17). Figure 2 (phase 2) shows the meta-data captured for the properties of `:bob` and `:alice`.

In the third phase, i.e., for support and confidence computation, the constraints' information stored in maps (Ψ_{ETD} , Ψ_{CEC}) is used to compute support and confidence for specific constraints. The algorithm iterates over the map Ψ_{ETD} to get the inner map Ψ_{ETPD}

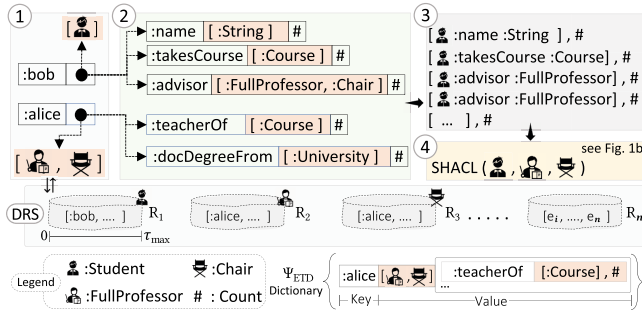


Figure 2: Overview of the four phases of QSE: ① entity extraction, ② entity constraints extraction, ③ support and confidence computation, and ④ shapes extraction. QSE-Approximate uses Dynamic Reservoir Sampling (DRS) in ①.

mapping entities to candidate property shapes $\phi_s: \langle \tau_p, T_p, C_p \rangle \in \Phi_s$, and retrieves the type of each entity using types information stored in Ψ_{ETD} to build triplets of the form $\langle \tau_e, \tau_p, \tau_{p_o} \rangle$ and compute their support and confidence (Line 25). Figure 2 (phase 3) highlights some of these triplets for $\tau_e = \text{Student}$. The value of support and confidence for each distinct triplet is incremented in each iteration and stored in Ψ_{SUPP} and Ψ_{CONF} maps. Additionally, a map Ψ_{PTT} (Property to Types) is populated with distinct properties' frequencies and their object types in order to, later on, establish the corresponding min/max cardinality constraints (Line 26).

Finally, in the shapes extraction phase, the algorithm iterates over the values of the Ψ_{CTP} map and defines the *shape name* of s , the *shape's target definition* τ_s , and the set of *shape constraints* ϕ_s for each candidate class (Lines 27-29). The set of property shapes P for a given *Node Shape* are then extracted from the map $\text{MAP}(\text{PROPERTY}, \text{SET})$ (Lines 30-36). An example shapes graph for our running example is shown in Figure 1. The C_p constraint can possibly have three types of values: `sh:Literal`, `sh:IRI`, and `sh:BlankNode`. In the case of literal types, the literal object types such as `xsd:string`, `xsd:integer`, or `xsd:date` are used. However, in the case of non-literal object types, the constraint `sh:class` is used to declare the type of object to define the type of value for the candidate property. It is possible to have more than one value for the `sh:class` and `sh:datatype` constraints of a candidate property shape, e.g., to state that a property can accept both integers and floats as values, in such cases, we use `sh:or` constraint to encapsulate multiple values. A detailed explanation of each phase is available in the extended version of the paper¹.

QSE-Exact (query-based). To support shapes extraction from a triplestore, we propose QSE-Exact query-based that uses a set of SPARQL queries [40] to extract all the necessary information that we collect across the four phases. In practice, we pose queries to extract all the distinct classes C , then, for each class $c \in C$, its properties $p \in \mathcal{P}$ along with object types are extracted as triplets, and support is computed for each triplet by a count query. This method is based on the standard procedure also implemented in other existing, query-based tools [12, 20].

Cardinality Constraints. QSE supports assigning cardinality constraints (`sh:minCount` and `sh:maxCount`) to C_p to each property shape constraint $\phi_s: \langle \tau_p, T_p, C_p \rangle$. Following the open-world

assumption, all shape constraints are initially assigned a minimum cardinality of 0, making them optional. However, there are cases where we can infer that some properties are mandatory (i.e., should be assigned a min count of 1), and some other properties should appear exactly once for each entity (i.e., should be assigned both a min and a max count equal to 1). Trivially one can assign minimum cardinality 1 to property shapes having confidence 100%, i.e., for those cases in which all entities have that property. In case of incomplete KGs, QSE allows users to provide a different confidence threshold value for adding the min cardinality constraints. To achieve this, we extend the fourth phase and add a min cardinality constraint in property shapes on line 35 based on the min-confidence provided by the user. QSE also keeps track of properties having maximum cardinality equal to 1 in a second phase and assigns `sh:maxCount=1` to those property shapes in the fourth phase of shapes extraction.

Complexity Analysis. The time complexity of QSE-Exact (Algorithm 1) is $O(2 \cdot |F| + |E| \cdot |\Phi_s| + |S| \cdot |\Phi_s|)$. Where $2 \cdot |F|$ refers to the first and second phases having to parse all the triples twice, E is the set of entities (i.e., the set of distinct IRIs that appear as a subject for some triple), S is the set of Node Shapes, and lastly, Φ_s represents a set of all property shape constraints, i.e., $\Phi_s = \{\phi_1, \phi_2, \dots, \phi_n\}$. Therefore, our algorithm scales linearly in the number of edges and nodes in the graph and in the size of the final set of shapes.

Algorithm 1 SHAPES EXTRACTION

```

Input: Graph  $\mathcal{G}$  from File  $F$ ,  $\omega$ : min-support,  $\epsilon$ : min-confidence
Output:  $S(s, \tau_s, \Phi_s)$ 
1:  $E_{data} \leftarrow \{T: \text{SET}_{Types}, \Psi_{ETPD} = \text{Map}(\text{IRI}, P_{data})\}$ 
2:  $P_{data} \leftarrow \{T': \text{SET}_{ObjTypes}, \text{Count}: \text{INT}\}$ 
3:  $\Psi_{ETD} = \text{Map}(\text{IRI}, E_{data}), \Psi_{CEC} = \text{Map}(\text{IRI}, \text{INT}), \Psi_{CTP} = \text{Map}(\text{IRI}, \text{MAP}(\text{IRI}, \text{SET}))$ 
4: for  $t \in \mathcal{G} \wedge t.p \neq \text{Type Predicate}$  do ▷ ① Entity extraction
5:    $\text{entity } e: t.s; \text{entityType } e_t = t.o$  ▷ s: subject, o: object
6:   if  $e \notin \Psi_{ETD}$  then  $\Psi_{ETD}.\text{insert}(e, \dots)$ 
7:    $\Psi_{ETD}.\text{insert}(e, \Psi_{ETD}.\text{get}(e).T.\text{add}(e_t))$  ▷ T: entity types
8:   increment entity count for current  $e_t$  in  $\Psi_{CEC}$ 
9: for  $t \in \mathcal{G} \wedge t.p \neq \text{Type Predicate}$  do ▷ ② Entity constraints extraction
10:   $\text{Set}_{ObjTypes} \leftarrow \emptyset, \text{Set}_{TUPLE} \leftarrow \emptyset$  ▷ init a type and property to type tuple set
11:  if object  $t.o$  is Literal then
12:     $\text{Set}_{ObjTypes}.\text{add}(\text{getLiteralType}(t.o))$ 
13:     $\text{Set}_{TUPLE}.\text{add}(\text{new Tuple}(t.p, \text{getLiteralType}(t.o)))$ 
14:  else ▷ for non-literal objects
15:    for  $\text{obj}_{type} \in \Psi_{ETD}.\text{get}(t.o).T$  do
16:       $\text{Set}_{ObjTypes}.\text{add}(\text{obj}_{type}); \text{Set}_{TUPLE}.\text{add}(\text{new Tuple}(t.p, \text{obj}_{type}))$ 
17:   $\text{addPropertyConstraints}(t.s, \text{Set}_{TUPLE}, \Psi_{ETD})$ 
18:  for  $\text{IRI} \in \Psi_{ETD}.\text{get}(t.s.T)$  do ▷ if  $t.s \in \Psi_{ETD}$ 
19:    update  $\Psi_{CTP}$  with class IRI,  $t.p$ , and object types using  $\text{Set}_{ObjTypes}$ 
▷ ③ Support and Confidence computation
20:  $\Psi_{SUPP} = \text{MAP}(\text{TUPLE}_3, \text{INT}), \Psi_{CONF} = \text{MAP}(\text{TUPLE}_3, \text{INT}), \Psi_{PTT}$ 
21: for  $(e, E_{data}) \in \Psi_{ETD}$  do
22:   for  $(T, \Psi_{ETPD}) \in E_{data}$  do
23:     for  $e_t \in T \wedge (p, p_o, c) \in P_{data}$  do
24:        $\chi \leftarrow \text{createTriplets}(\langle \tau_e, \tau_p, \tau_{p_o} \rangle)$ 
25:        $\text{computeSupportAndConfidence}(\Psi_{SUPP}, \chi, \Psi_{CEC})$ 
26:        $\text{computeMaxCardinality}(\Psi_{PTT}, p, c)$ 
27: for  $(\text{class}, \text{MAP}(\text{Property}, \text{SET}_{ObjTypes})) \in \Psi_{CTP}$  do ▷ ④ Shapes extraction
28:   $\Phi_s \leftarrow \emptyset$  ▷ Property shapes  $\Phi_s = \{\phi_{s_1}, \phi_{s_2}, \dots, \phi_{s_n}\}$  where  $\phi_s: \langle \tau_p, T_p, C_p \rangle$ 
29:   $s = \text{class}.\text{buildShapeName}(), \tau_s = \text{class}$ 
30:  for  $(p, \text{SET}_{ObjTypes}) \in \text{MAP}(\text{Property}, \text{SET})$  do  $\phi_s.\tau_p = p$ 
31:   $p.\omega = \Psi_{SUPP}.\text{get}(p, \text{SET}_{ObjTypes}), p.\epsilon = \Psi_{CONF}.\text{get}(p, \text{SET}_{ObjTypes})$ 
32:  if  $p.\omega > \omega \wedge p.\epsilon > \epsilon$  then
33:     $\text{build}(\text{sh:nodeKind}, \text{sh:maxCount}, \Psi_{PTT})$ 
34:     $\phi_s.C_p.\text{add}(\text{sh:minCount}: 1)$  ▷ if  $p.\epsilon > \epsilon'$ 
35:     $\Phi_s.\text{add}(\phi_s)$ 
36:   $S.\text{add}(s, \tau_s, \Phi_s)$  ▷ if  $s.\omega > \omega \wedge \phi_s!0$ 

```

¹<https://relweb.cs.aau.dk/qse/>

4 QSE-APPROXIMATE

QSE-Exact keeps type and property information for each entity in memory while extracting shapes. As a result, its memory requirements are prohibitively large when dealing with large KGs. Therefore, we propose QSE-Approximate to enable shape extraction from very large KGs with reduced memory requirements. Our goal is to *solve the scalability issue in shapes extraction approaches by using only the resources available to a commodity machine*. QSE-Approximate is based on a multi-tiered dynamic reservoir-sampling (DRS) algorithm that we introduce. We maintain as many reservoirs as types in the graph, and we dynamically resize each reservoir as new triples are parsed. Moreover, the replacement of nodes in the reservoir is performed based on the number of node types across reservoirs. The resulting algorithm replaces the first phase of QSE. After sampling, the information about the sampled entities is used in the same way as before in the remaining phases of Algorithm 1. Hence, we maintain information only for a small representative sample of entities in memory but enough to detect all shapes.

Algorithm 2 receives as input a graph file F , sampling percentage (Sampling%), and maximum size of the reservoir per class (τ_{max}). After initialization, triples t of F are parsed (Line 3) and filtered based on whether they contain a type declaration. From these, we extract the entities to populate the Entity-to-Data map Ψ_{ETD} (Lines 4-24), while non-type triples are parsed on Line 24 to keep count of distinct properties in the Property-Count map Ψ_{PC} . For instance, `:alice` is an entity of type `:FullProfessor` and `:Chair` in Ψ_{ETD} shown in Figure 2. QSE-Approximate maintains a reservoir for each distinct entity type e_t , e.g., maintaining a distinct reservoir of entities of type `:Student` (R_1), `:FullProfessor` (R_2), and `:Chair` (R_3) shown in Figure 2, using a map of sampled entities per class (Ψ_{SEPC}). The reservoir capacity map (Ψ_{RCPC}) stores the current max capacities for the reservoir for each e_t . If e_t does not exist in Ψ_{SEPC} and Ψ_{RCPC} , i.e., if it has not a reservoir, one is created (lines 6-7). Then, e is inserted in the reservoir for e_t (Lines 8-11), e.g., `:alice` is inserted into both reservoirs R_2 and R_3 shown in Figure 2. If the reservoir has reached its current capacity limit, we may have to replace an entity in the reservoir with the current one. Hence, neighbor-based dynamic reservoir sampling is performed (Lines 13-18), i.e., a random number r is generated between zero and the current number of type declarations read from F . If r falls within the reservoir size, then a node in the reservoir is replaced with e . To select which node to replace, we identify as \tilde{n} the target node at index r , and with \bar{n} and $\bar{\bar{n}}$ its neighbors at indexes $r-1$ and $r+1$, respectively. Among these, the node having minimum scope (i.e., the minimum number of types that are known at this point in time) is selected to be replaced by the current e (Line 17). Additionally, the algorithm keeps track of actual Class-to-Entity-Count in Ψ_{CEC} (Line 19), i.e., the exact count of how many entities of each type we have seen. Once the reservoir for e_t is updated, the sampling ratio for this type is computed, i.e., the proportion of entities kept so far with type e_t over the total number of entities of that type seen up to now. Given the current and target sampling ratio (Sampling%) provided as input, the algorithm evaluates whether to resize the reservoir for e_t , if it has not already reached the limit τ_{max} (Lines 21-23).

While performing shapes pruning using counts over sampled entities, QSE-Approximate requires to estimate actual support $\bar{\omega}_\phi$

Algorithm 2 QSE-APPROXIMATE RESERVOIR SAMPLING

Input: Graph \mathcal{G} from File F , maximum entity threshold τ_{max} , Sampling%
Output: Ψ_{ETD}, Ψ_{CEC}

```

1: init maps  $\Psi_{ETD}, \Psi_{SEPC}, \Psi_{RCPC}, \Psi_{CEC}, \Psi_{PC}$ 
2:  $\tau_{min} = 1$  (minimum entity threshold); lineCounter = 0
3: for  $t \in \mathcal{G}$  do ▷ parse s,p,o of the triple t
4:   if  $t.p = \text{Type Predicate}$  then
5:     entity  $e : t.s$ ; entityType  $e_t = t.o$  ▷ s: subject, o: object
6:      $\Psi_{SEPC}.putIfAbsent(e_t, [ ])$  ▷ if  $e_t \notin \Psi_{SEPC}$ 
7:      $\Psi_{RCPC}.putIfAbsent(e_t, \tau_{min})$  ▷ if  $e_t \notin \Psi_{RCPC}$ 
8:     if  $|\Psi_{SEPC}.get(e_t)| < \Psi_{RCPC}.get(e_t)$  then ▷ Add entity  $e$  in reservoir
9:       if  $\Psi_{ETD}.get(e).T$  is  $\emptyset$  then  $\Psi_{ETD}.insert(e, \dots)$  ▷ T : entity types
10:       $\Psi_{ETD}.insert(e, \Psi_{ETD}.get(e).T.add(t.o))$ 
11:       $\Psi_{SEPC}.get(e_t).insert(e)$ 
12:   else ▷ Replace random entity in reservoir with current entity  $e$ 
13:      $r = \text{generateRandomNumber}(0, \text{lineCounter})$ 
14:     if  $r < |\Psi_{SEPC}.get(e_t)|$  then
15:        $\tilde{n}, \bar{n}, \bar{\bar{n}} = \Psi_{SEPC}.get(e_t).nodeAtIndex(r - 1, r, r + 1)$ 
16:        $n = \text{getNodeWithMinimumScope}(\tilde{n}, \bar{n}, \bar{\bar{n}})$ 
17:       replace node at index  $n$  with current  $e$  &  $e_t$  in  $\Psi_{ETD}$ 
18:        $\Psi_{SEPC}.get(e_t).add(e)$ 
19:   increment entity count for current  $e_t$  in  $\Psi_{CEC}$ 
20:   ▷ Resize reservoir
21:   ratio =  $(\Psi_{SEPC}.get(e_t).size() / \Psi_{CEC}.get(e_t)) \times 100$ 
22:   capacity = Sampling%  $\times \Psi_{SEPC}.get(e_t).size()$ 
23:   if capacity <  $\tau_{max} \wedge \text{ratio} \leq \text{Sampling\%}$  then  $\Psi_{RCPC}.insert(e_t, \text{capacity})$ 
24:   else  $\rightarrow$  increment property count for current  $t.p$  in  $\Psi_{PC}$ 
25:   lineCounter + +
```

and confidence $\bar{\varepsilon}_\phi$ of a property shape ϕ from the current values ω and ε computed from the sampled data. Thus, it estimates with $\bar{\omega}_\phi = \omega_\phi / \min(|P_r^*|/|P|, |T_r|/|T|)$ the effective support for a property shape ϕ , where ω_ϕ is the support computed for ϕ in the current sample, P represents all triples in \mathcal{G} having property τ_p , P_r^* represents triples having property τ_p across all entities in all reservoirs, T represents all entities of type e_t in \mathcal{G} , and T_r represents all entities of type e_t in the reservoir. Similarly, the confidence $\bar{\varepsilon}_\phi$ of a property shape is estimated by replacing denominator in eq. (2) with $|T_r|$.

QSE-Approximate (query-based). We apply the same sampling technique in the query-based shapes extraction approach where in Algorithm 2 entities and their meta-data are retrieved via SPARQL queries, resulting in *query-based* QSE-Approximate.

Space Analysis. The space requirement of QSE-Approximate depends on the values of target Sampling%, the maximum reservoir size τ_{max} , and the number of entity types $|T|$ in \mathcal{G} . In the worst case, it requires $O(2 \cdot |T| \cdot \tau_{max})$, therefore while \mathcal{G} can contain hundreds of millions of entities, we can still easily estimate how many distinct types are in the graph and select τ_{max} to fit the available memory.

5 EVALUATION

In the following, we evaluate our QSE solutions and their effectiveness in tackling the problem of *spuriousness* along with a comparison to existing state-of-the-art approaches.

Datasets. We selected a synthetic dataset, LUBM-500 [18], and three real-world datasets: DBpedia [4] downloaded on 01.10.2020; YAGO-4 [46], for which we use the subset containing instances from the English Wikipedia, downloaded on 01.12.2020; and WikiData [49], in two variants, i.e., a dump from 2015 [54] (Wdt15), used in the original evaluation of SheXer [12], and the truthy dump from September 2021 (Wdt21) filtered by removing non-English strings. Table 1 provides a comparison of their contents.

Experimental Setup. We have implemented QSE algorithms in JAVA-11. All experiments are performed on a single machine

with Ubuntu 18.04, having 16 cores and 256 GB RAM. We have used GraphDB [16] 9.9.0 to experiment with *query-based* variants of QSE with a maximum memory usage limit of 16 GB. The source code of QSE is available as open-source [37] along with experimental settings and datasets. We have also published the extracted SHACL shapes of all our datasets on Zenodo [39]. For SheXer, we cloned its original code from GitHub and used the same settings as the original paper, i.e., default tuned parameters for the sheXing process and customized tuned parameters to output shapes equivalent to QSE.

Metrics. We measure the *running time*, and maximum *memory usage* (defined using Java `-Xmx`) during QSE shapes extraction process, and *Shape Statistics* of the output shapes.

QSE-Exact. We use QSE-Exact to extract shapes from LUBM (L), DBpedia (D), YAGO-4 (Y), and WikiData (W). The statistics of the shapes extracted from these datasets using QSE-Exact (file-based) are shown in Table 2. It shows the count of Node Shapes (NS), Property Shapes (PS), and Property Shape Constraints (PSc), i.e., literal and non-literal node types constraints. We refer to these statistics as *default shape statistics*. We initially considered SheXer [12], ShapeDesigner [5], and SHACLGEN [20] as state-of-the-art approaches [38] to compare against QSE. Among these, both ShapeDesigner and SHACLGEN load the whole graph into a triplestore similar to our QSE-Exact (query-based). Yet, their current implementations cannot handle large KGs with more than a few million triples and do not manage to extract shapes of KGs having more than some hundreds of classes. In our experiments, either they crashed because they tried to load the graph into an in-memory triplestore or required multiple hours to generate shapes for large KGs such as YAGO-4 (with 8,897 classes). Therefore, in the following, we focus our comparison on SheXer, which supports both the file-based and the query-based methods. Table 3 shows the running time and memory consumption to extract shapes for all datasets using File (F) and Query-based (Q) variants of SheXer, QSE-Exact, and QSE-Approximate. Among the *file-based* approaches, QSE-Exact is 1 order of magnitude faster than SheXer for all datasets. It consumes up to 50% less memory than SheXer to extract shapes from D, L, Y,

Table 1: Size and characteristics of the datasets

| | DBpedia | LUBM | YAGO-4 | Wdt15 | Wdt21 |
|-----------------|---------|-------|--------|--------|---------|
| # of triples | 52 M | 91 M | 210 M | 290 M | 1.926 B |
| # of objects | 19 M | 12 M | 126 M | 64 M | 617 M |
| # of subjects | 15 M | 10 M | 5 M | 40 M | 196 M |
| # of literals | 15 M | 5.5 M | 111 M | 40 M | 904 M |
| # of instances | 5 M | 1 M | 17 M | 3 M | 91 M |
| # of classes | 427 | 22 | 8,902 | 13,227 | 82,693 |
| # of properties | 1,323 | 20 | 153 | 4,906 | 9,017 |
| Size in GBs | 6.6 | 15.66 | 28.59 | 42 | 234 |

Table 2: Shapes Statistics using QSE-Exact.

| | NS | PS | Non-Literal PSc | Literal PSc |
|---------|--------|------------------|-----------------|-----------------|
| | COUNT | COUNT/AVG | COUNT/AVG | COUNT/AVG |
| LUBM | 23 | 164 / 7.1 | 323 / 3.0 | 57 / 1.0 |
| DBpedia | 426 | 11,916 / 27.9 | 38,454 / 6.9 | 5,335 / 1.0 |
| YAGO-4 | 8,897 | 76,765 / 8.6 | 315,413 / 14.5 | 50,708 / 1.0 |
| Wdt15 | 13,227 | 202,085 / 15.2 | 114,890 / 3.0 | 106,599 / 1.0 |
| Wdt21 | 82,651 | 2,051,538 / 24.8 | 3,765,953 / 5.6 | 1,113,856 / 1.0 |

Table 3: Running Time (T) in minutes (m) and hours (h) along with Memory (M) consumption in GB.

| | | DBpedia | | LUBM | | YAGO-4 | | Wdt15 | | Wdt21 | |
|---|------------|-------------|----|-------------|-----|------------------|----|--------------|-----|------------------|------------------|
| | | T | M | T | M | T | M | T | M | T | M |
| F | SheXer | 26 m | 18 | 58 m | 33 | 1.9 h | 24 | 3.2 h | 59 | - | Out _M |
| | QSE-Exact | <u>3 m</u> | 16 | <u>8 m</u> | 16 | <u>23 m</u> | 16 | <u>16 m</u> | 50 | <u>2.5 h</u> | 235 |
| | QSE-Approx | 1 m | 10 | 2 m | 10 | 13 m | 10 | 13 m | 16 | 1.3 h | 32 |
| Q | SheXer | 9 h | 65 | 15 h | 140 | Out _T | - | 13 h | 180 | Out _T | - |
| | QSE-Exact | <u>34 m</u> | 16 | <u>47 m</u> | 16 | <u>2.4 h</u> | 16 | <u>1.2 h</u> | 16 | Out _T | - |
| | QSE-Approx | 16 m | 6 | 3 m | 7 | 39 m | 16 | 49 m | 16 | 5.7 h | 64 |

and Wdt15, whereas SheXer goes out of memory (Out_M) for Wdt21. Similarly, among the *query-based* approaches, QSE-Exact is 1 order of magnitude faster and consumes less than 50% memory to extract shapes from D, Y, L, and Wdt15. SheXer timed out (Out_T – 24 hours) for Y and Wdt21, while QSE-Exact timed out for Wdt21 only.

Taming spuriousness. To deal with the issue of spuriousness, we analyze the shapes extracted and kept after pruning. QSE performs *support-based shapes extraction* by producing only the shapes with support and confidence greater than or equal to a threshold specified by the user. For instance, given a minimum support threshold of 100 and minimum confidence value 25%, for every PS, QSE prunes all the PSc that do not appear with at least 100 entities or if not at least for 25% of entities for that type. We remind that the pruning of PSc has a cascading effect that also affects the pruning of PS, and the pruning of PS can, in turn, cause the pruning of NS. To study the impact of various confidence and support thresholds on the number of PSc, PS, and NS, we analyze the effect of pruning by specifying various values for confidence and support. Figure 3 shows the result of pruning PSc (3a,b), PS and NS (3c,d) for confidence >(25, 50, 75,90)% and support ($\geq 1, >100$) on DBpedia and Wdt21. Experimental results on LUBM, YAGO-4, and Wdt15 are comparable to the results presented for DBpedia and Wdt21, and are reported in the extended version of the paper¹. In general, as expected, the results show that the higher we set the threshold for support and confidence, the higher the percentage of PSc and PS to be pruned. Precisely, DBpedia contains 11K PS, 38K non-literal, and 5K literal PSc (Table 2), when QSE performs pruning with confidence >25% and support ≥ 1 , it prunes out 99% PSc and PS (Figure 3a,b). Similarly for Wdt21, QSE prunes 85% non-literal and 97% literal constraints, and 66% PS for confidence >25% and support ≥ 1 (Figure 3b). In comparison to the default shape statistics (Table 2), increasing confidence to >50%, 75%, and 90%, pruning

Table 4: QSE-Approximate: Effect of Sampling_s (S%) and reservoir size (τ_{max}) on Precision (P), Recall (R), and Relative Error (Δ) with min. support 1 and confidence 25% on Wdt21

| S% | τ_{max} | Property Shapes (PS) | | | | Time (Min) | Mem (GB) |
|------|--------------|----------------------|---------|-------------|----------|------------|----------|
| | | Real | Sample | P / R | Δ | | |
| 10% | 20 | 698,825 | 470,562 | 1.00 / 0.61 | 228,263 | 81 | 16 |
| | 200 | 698,825 | 497,035 | 0.92 / 0.65 | 201,790 | 81 | 16 |
| 50% | 500 | 698,825 | 548,381 | 0.96 / 0.79 | 150,444 | 82 | 24 |
| | 5000 | 698,825 | 605,785 | 0.96 / 0.83 | 93,040 | 95 | 24 |
| 100% | 500 | 698,825 | 617,349 | 1.00 / 0.88 | 81,476 | 87 | 32 |
| | 5000 | 698,825 | 645,810 | 1.00 / 0.92 | 53,015 | 98 | 32 |

resulted in a drastic decrease in the number of PSc and PS. In DBpedia, the majority of non-literal PSc are pruned out, and in Wdt21, the majority of literal constraints are pruned out. Pruning of NS is lower compared to PS and PSc for all combinations of support and confidence, showing that almost all types are associated at least with some very common PSc, e.g., the fact to have a :name.

QSE-Approximate. The QSE-Approximate approach reduces the memory requirements of the exact approach by allowing users to specify the *sampling percentage* (Sampling%, S% for short) and maximum limit of the *reservoir size* (τ_{max}), i.e., the maximum number of entities to be sampled per class, to reduce the number of entities to keep in memory. Table 3 shows that among the *file-based* approaches, QSE-Approximate is the most efficient approach compared to QSE-Exact and SheXer. For example, to extract shapes from Wdt21, QSE-Approximate (with $\tau_{max} = 1000$ and S%=100%) required almost half the time with 1 order of magnitude less memory than QSE-Exact, while SheXer could not complete the computation. Similarly, among *query-based* approaches, QSE-Approximate proved to be the only approach to extract shapes from the Wdt21 endpoint in 5.7 hours with 64 GB memory consumption. In contrast, QSE-Exact and SheXer timed out (24 hours). Analogously to Wdt21, QSE-Approximate remains 1 order of magnitude faster with 50% less memory consumption than SheXer (for both query and file-based variants) to extract shapes from D, L, Y, and Wdt15. Overall, these results show that our proposals have solved scalability issues in shape extraction approaches regardless of the type of input data source (file or endpoint). The choice of using a query-based or file-based version depends on the given setting. For instance, querying an endpoint to extract shapes can impose excessive stress on a production DBMS serving other applications. On the other hand, the file-based approach is less resource-intensive and can be used if the user can afford the cost of dumping the graph into a file.

QSE Sampling Parameters. We further evaluate the quality of the output of QSE-Approximate using multiple combinations of values for S% and τ_{max} on Wdt21 with a fixed confidence and support threshold. This analysis helps the user to choose the best values for S% and τ_{max} parameters given some memory constraints. We show the results in Table 4, where the values shown in columns *Real* and *Sampled* are extracted by QSE-Exact and QSE-Approximate, respectively. Here we skip listing values for NS as they are not affected by the values of S%, τ_{max} , confidence, and support. The results show that S%=10 and τ_{max} up to 200 provide a 92% precision for PS extracted using QSE-Approximate and pruned with support ≥ 1 and confidence $>25\%$. This requires only 16 GB RAM and 81 minutes. If a machine up to 24 GB RAM is available, then S%=50% and $\tau_{max}=5K$ provide 96% precision with $\Delta = 93K$ in 95 minutes. Similarly, on a machine having 32 GB RAM, S%=100% and $\tau_{max}=5K$ provide 100% precision with $\Delta = 53K$ in 98 minutes. The non-perfect precision translates into some shapes being produced despite their support and confidence is slightly lower than required. We also see that, for very small values of τ_{max} we achieve a lower recall, meaning that some shapes that should have been produced are instead wrongly pruned. We note though that min support 1 and confidence 25% are still quite low values and the shapes produced are thus more affected by spuriousness. Nonetheless, on a standard commodity machine with 32GB we see we can easily achieve perfect precision (100%) and very high recall (92%).

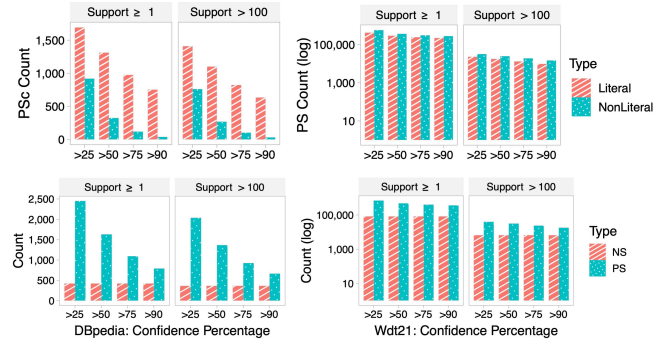


Figure 3: QSE-Exact on DBpedia and Wdt21

We further study the effect of pruning on shapes extracted from Wdt21 using QSE-Approximate with confidence $>25\%$ and $>75\%$ having support 1, 10, and 100 (shown in Table 5). We see that with support ≥ 1 and confidence $>25\%$, QSE-approximate is able to get almost all the PS extracted by QSE-Exact for Wdt21 (Figure 3d) having 89% recall and 100% precision. Additionally, upon increasing the support to 10 and 100, we notice a constant recall of around 88-99% and a slight reduction in precision, i.e., 98% and 96% with decreasing relative error (i.e., Δ). Similarly, we notice the same trend with confidence=75%. Therefore, while we very rarely overestimate the support and confidence of the shapes produced, we underestimate some of these values, although still in a few cases only.

Practical Implications of QSE. We show the practical utility of QSE by evaluating the correctness of extracted shapes and their effect when used to validate the KG. We extracted shapes from DBpedia using QSE with confidence $>25\%$ and support >100 . Then, we randomly selected 10 shapes and manually inspected them to evaluate their correctness, i.e., whether these shapes describe valid constraints. This allows us to measure precision and recall based on the pruning parameters. The results of this analysis showed that QSE extracts shapes with 100% precision in terms of correct shapes constraints that should be part of the final set of shapes (qualified as quality shapes) by removing spurious shape constraints. Further, we used these 10 shapes, extracted by QSE, to validate DBpedia using a SHACL validator and found 20,916 missing triples and 155 erroneous triples. The detailed results of this analysis are contained in the extended version¹. Overall, this experiment shows that by using our technique the user is provided with a refined set of valid shapes that can effectively identify errors in the KG.

Table 5: Output quality of QSE-Approximate on Wdt21 with S% = 100% and $\tau_{max} = 500$ as # of real and sampled NS, PS, and corresponding Precision (P), Recall (R), and Relative Error Δ .

| Conf | Supp | Node Shapes (NS) | | | | Property Shapes (PS) | | | |
|----------|----------|------------------|--------|-----------|----------|----------------------|---------|-------------|----------|
| | | Real | Sample | P / R | Δ | Real | Sample | P / R | Δ |
| $> 25\%$ | ≥ 1 | 82,651 | 82,651 | 1.0 / 1.0 | 0 | 698,825 | 620,622 | 1.00 / 0.89 | 78,203 |
| | 10 | 23,640 | 23,640 | 1.0 / 1.0 | 0 | 158,283 | 141,040 | 0.99 / 0.88 | 17,243 |
| | 100 | 6,596 | 6,596 | 1.0 / 1.0 | 0 | 39,877 | 36,362 | 0.96 / 0.88 | 3,515 |
| $> 75\%$ | ≥ 1 | 82,651 | 82,651 | 1.0 / 1.0 | 0 | 405,344 | 362,717 | 1.00 / 0.89 | 42,627 |
| | 10 | 23,640 | 23,640 | 1.0 / 1.0 | 0 | 91,947 | 83,329 | 0.99 / 0.90 | 8,618 |
| | 100 | 6,596 | 6,596 | 1.0 / 1.0 | 0 | 23,944 | 22,193 | 0.97 / 0.90 | 1,751 |

Constraints Coverage. Comparing the constraints supported by QSE and existing approaches (i.e., SheXer [12], SHACLGEN [20], and ShapeDesigner [5]), we report that QSE is able to extract the widest range of constraints (i.e., 15 out of 16 specific core constraints). Amongst those that are usually not supported, we support `sh:in`, `sh:Literal`, `sh:class`, `sh:not`, and `sh:node`. We currently do not support `sh:inverse` but we plan to support it in the future. More details are available in the extended version of our paper¹.

Optimal Pruning Thresholds. For each class in the KG, QSE computes its frequency. Thus, this information can be used as a reference for the support and confidence thresholds. Further, QSE also supports the extraction of shapes for specific classes only. Therefore, the user can make use of frequency information and set class-specific pruning thresholds.

6 RELATED WORK

KG Data Validation. Integrity constraints for KGs were initially defined with the RDF schema vocabulary [11] and then with the OWL language [27, 28, 47]. Later, the SPARQL Inferencing Notation (SPIN) [22] was proposed. SHACL [23] (a W3C standard since 2017) is known as the next generation of SPIN. Similar to SHACL, ShEX [35] is a constraint language that is built on regular bag expressions inspired by schema languages for XML. While ShEx is not a standard, it is used within the WikiData project [48]. Even though SHACL and ShEx are not completely equivalent [15], their core mechanism revolves around the same concept of enforcing for each node to satisfy specific constraints on the combination of its types and predicates [12]. In this work, we support the extraction of validating shapes that can be represented in both languages.

Shape Extraction. Given the abundance of large-scale KGs, various applications have been created to assist the process of extracting information about its implicit or explicit schema [21]. Among these, shapes construction or extraction approaches, i.e., to generate a set of shapes given information from an existing KG, are used in order to obtain validating schema to ensure the quality of a KG’s content. We have classified existing approaches in Table 6 based on their features, i.e., support for shapes extraction from data or ontologies, support for automatic extraction of shapes, support for shapes extraction from a SPARQL triplestore, and whether they extract SHACL, ShEx, or both types of validating shapes. In our recent community survey [38] on extraction and adoption of validating shapes, we show that *there is a growing need among practitioners for techniques for efficient extraction of validating shapes from very*

Table 6: State-of-the-art to extract validating shapes [38]

| Approach | Extracted from | | Auto-matic | Triple-store | Type |
|----------------------|----------------|----------|------------|--------------|------------|
| | data | ontology | | | |
| Shape Induction [26] | ✓ | ✗ | ✓ | ✓ | SHACL,ShEx |
| SheXer [12] | ✓ | ✗ | ✓ | ✓ | SHACL,ShEx |
| Spahiu et al. [45] | ✓ | ✗ | ✓ | ✓ | SHACL |
| ShapeDesigner. [5] | ✓ | ✗ | ✓ | ✓ | SHACL,ShEx |
| SHACLGEN [20] | ✓ | ✓ | ✓ | ✓ | SHACL |
| TopBraid [36] | ✓ | ✓ | ✓ | ✓ | SHACL |
| Pandit et al. [32] | ✗ | ✓ | ✗ | ✓ | SHACL |
| Astrea [9] | ✗ | ✓ | ✓ | ✗ | SHACL |
| SHACLearner [30] | ✓ | ✗ | ✓ | ✗ | SHACL |
| Groz et al. [17] | ✓ | ✗ | ✓ | ✗ | ShEx |

large existing KGs. Note that there exist approaches for schema extraction from property graphs as well [24]. Such approaches are not directly applicable to RDF KGs since their schema is more complex, moreover they focus on identifying sub-types based on node labels (which do not exist in RDF data, since types are nodes in the graph), and finally are not designed to handle the issue of spuriousness. Once shapes are extracted, they can be used to validate KGs using validation approaches like MagicShapes [2] and Trav-SHACL [13].

Rules, Patterns, and Summaries. There exist various approaches for rule discovery in graphs [25]. These systems [1, 14, 31] derive rules from large KGs using structural information by exploring the frequently occurring graph patterns. In contrast to validating shapes, rules are mainly used to derive new facts from an incomplete KG or identify specific sets of wrong connections. Frequent subgraph mining (FPM) approaches, instead, are designed to find frequently recurring structures in a large graph. In FPM, the occurrence of subgraphs (the number of times a subgraph appears) cannot be taken as the support of subgraphs since it does not satisfy the non-monotonic property [7]. The most practical measurement for measuring this support is, instead, the minimum image-based support (MNI [7]). Our proposed definition of support for shape constraints is inspired by the concept of MNI support and its use in FPM [19]. Yet, different than FPM, we do not extract patterns of arbitrary shape and size, thus we are able to provide better performance guarantees as we solve a simpler problem. Finally, our approach is also related to the techniques of graph summarization [8] and can be seen as a special form of structural summarization [33]. Additionally, QSE provides a scalable solution for understanding the content of large KGs (by extracting their shapes) like ABSTAT-HD [3], which is based on exploring semantic profiles of large KGs.

7 CONCLUSION

In this paper, we propose an automatic shape extraction approach that addresses the two common limitations in other existing techniques, i.e., scalability and spuriousness. We addressed these limitations by introducing the **QUALITY SHAPES EXTRACTION (QSE)** problem. We devised an exact and approximate solution for QSE to enable the efficient extraction of shapes on commodity machines. Our method is based on the well-understood concepts of support and confidence, hence it allows a data scientist to focus on the shapes providing the highest reliability first when addressing issues of data quality. By setting even low pruning thresholds, QSE can prune up to 93% of the shapes that a trivial extraction would produce (i.e., a reduction of 2 orders of magnitude), shapes that hence have little support from the data and are thus likely spurious. Furthermore, we show that our approximate technique introduces only negligible loss in the quality and completeness of the produced shapes. In the future, we will extend the scope of constraints covered by QSE and a solution to automatically learn the optimal configurations for pruning thresholds for QSE.

ACKNOWLEDGMENTS

This research was partially funded by the Danish Council for Independent Research (DFF) under grant agreement no. DFF-8048-00051B, the EU’s H2020 research and innovation programme under grant agreement No 838216, and the Poul Due Jensen Foundation.

REFERENCES

- [1] Naser Ahmadi, Thi-Thuy-Duyen Truong, Le-Hong-Mai Dao, Stefano Ortona, and Paolo Papotti. 2020. Rulehub: A public corpus of rules for knowledge graphs. *Journal of Data and Information Quality (JDQ)* 12, 4 (2020), 1–22.
- [2] Shqiponja Ahmetaj, Bianca Löhnert, Magdalena Ortiz, and Mantas Šimkus. 2022. Magic shapes for SHACL validation. *Proceedings of the VLDB Endowment* 15, 10 (2022), 2284–2296.
- [3] Renzo Arturo Alva Principe, Andrea Maurino, Matteo Palmonari, Michele Ciavotta, and Blerina Spahiu. 2022. ABSTAT-HD: a scalable tool for profiling very large knowledge graphs. *The VLDB Journal* 31, 5 (2022), 851–876.
- [4] Sören Auer, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, and Zachary G. Ives. 2007. DBpedia: A Nucleus for a Web of Open Data. In *The Semantic Web, 6th International Semantic Web Conference (Lecture Notes in Computer Science)*, Vol. 4825. Springer, Busan, Korea, 722–735.
- [5] Iovka Boneva, Jérémie Dusart, Daniel Fernández-Álvarez, and José Emilio Labra Gayo. 2019. Shape Designer for ShEx and SHACL constraints. In *Proceedings of the ISWC 2019 Satellite Tracks (CEUR Workshop Proceedings)*, Vol. 2456. CEUR-WS.org, Auckland, New Zealand, 269–272.
- [6] Christian Borgelt. 2012. Frequent item set mining. *Wiley interdisciplinary reviews: data mining and knowledge discovery* 2, 6 (2012), 437–456.
- [7] Björn Bringmann and Siegfried Nijssen. 2008. What Is Frequent in a Single Graph?. In *PAKDD (Lecture Notes in Computer Science)*, Vol. 5012. Springer, Osaka, Japan, 858–863.
- [8] Šejla Čebirić, François Goasdoué, Haridimos Kondylakis, Dimitris Kotzinos, Ioana Manolescu, Georgia Troullinou, and Mussab Zneika. 2019. Summarizing semantic graphs: a survey. *The VLDB journal* 28, 3 (2019), 295–327.
- [9] Andrea Cimmino, Alba Fernández-Izquierdo, and Raúl García-Castro. 2020. Astrea: Automatic Generation of SHACL Shapes from Ontologies. In *ESWC (Lecture Notes in Computer Science)*, Vol. 12123. Springer, Heraklion, Crete, Greece, 497.
- [10] WWW Consortium. 2014. RDF 1.1. <https://w3.org/RDF/>. Accessed 17th January, 2023.
- [11] WWW Consortium. 2014. RDF Schema 1.1. <https://www.w3.org/TR/rdf-schema/>. Accessed 17th January, 2023.
- [12] Daniel Fernandez-Álvarez, Jose Emilio Labra-Gayo, and Daniel Gayo-Avello. 2022. Automatic extraction of shapes using sheXer. *Knowledge-Based Systems* 238 (2022), 107975.
- [13] Mónica Figuera, Philipp D. Rohde, and Maria-Esther Vidal. 2021. Trav-SHACL: Efficiently Validating Networks of SHACL Constraints. In *WWW '21: The Web Conference 2021, Virtual Event / Ljubljana, Slovenia, April 19-23, 2021. ACM / IW3C2, Ljubljana, Slovenia, 3337–3348*. <https://doi.org/10.1145/3442381.3449877>
- [14] Luis Galárraga, Christina Teflioudi, Katja Hose, and Fabian M Suchanek. 2015. Fast rule mining in ontological knowledge bases with AMIE+. *The VLDB Journal* 24, 6 (2015), 707–730.
- [15] Jose Emilio Labra Gayo, Eric Prud'Hommeaux, Iovka Boneva, and Dimitris Kontokostas. 2017. Validating RDF data. *Synthesis Lectures on Semantic Web: Theory and Technology* 7, 1 (2017), 1–328.
- [16] GraphDB. 2023. GraphDB. <https://graphdb.ontotext.com>. Accessed 17th January, 2023.
- [17] Benoît Groz, Aurélien Lemay, Slawek Staworko, and Piotr Wiecezorek. 2022. Inference of Shape Graphs for Graph Databases. In *25th International Conference on Database Theory, ICDT 2022, March 29 to April 1, 2022, Edinburgh, UK (Virtual Conference) (LIPICs)*, Vol. 220. Schloss Dagstuhl - Leibniz-Zentrum für Informatik, Edinburgh, UK, 14:1–14:20. <https://doi.org/10.4230/LIPICs.ICDT.2022.14>
- [18] Y. Guo, Z. Pan, and J. Heflin. 2005. LUBM: A benchmark for OWL knowledge base systems. *Journal of Web Semantics* 3 (2005), 158–182.
- [19] Jiawei Han, Hong Cheng, Dong Xin, and Xifeng Yan. 2007. Frequent pattern mining: current status and future directions. *Data mining and knowledge discovery* 15, 1 (2007), 55–86.
- [20] Alexis Keely. 2023. SHACLGEN. <https://pypi.org/project/shaclgen/>. Accessed 17th January, 2023.
- [21] Kenza Kellou-Menouer, Nikolaos Kardoulakis, Georgia Troullinou, Zoubida Kedad, Dimitris Plexousakis, and Haridimos Kondylakis. 2022. A survey on semantic schema discovery. *VLDB J.* 31, 4 (2022), 675–710. <https://doi.org/10.1007/s00778-021-00717-x>
- [22] Holger Knublauch, James A Hendler, and Kingsley Idehen. 2011. SPIN-overview and motivation. *W3C Member Submission* 22 (2011), W3C.
- [23] Holger Knublauch and Dimitris Kontokostas. 2017. Shapes constraint language (SHACL). *W3C Candidate Recommendation* 11, 8 (2017).
- [24] Haná Lbath, Angela Bonifati, and Russ Harmer. 2021. Schema Inference for Property Graphs. In *Proceedings of the 24th International Conference on Extending Database Technology, EDBT 2021, Nicosia, Cyprus, March 23 - 26, 2021. OpenProceedings.org, Nicosia, Cyprus, 499–504*. <https://doi.org/10.5441/002/edbt.2021.58>
- [25] Michael Loster, Davide Mottin, Paolo Papotti, Jan Ehmler, Benjamin Feldmann, and Felix Naumann. 2021. Few-Shot Knowledge Validation using Rules. In *WWW '21: The Web Conference 2021, Virtual Event / Ljubljana, Slovenia, April 19-23, 2021. ACM / IW3C2, Ljubljana, Slovenia, 3314–3324*. <https://doi.org/10.1145/3442381.3450040>
- [26] Nandana Mihindukulasooriya, Mohammad Rifat Ahmmad Rashid, Giuseppe Rizzo, Raúl García-Castro, Óscar Corcho, and Marco Torchiano. 2018. RDF shape induction using knowledge base profiling. In *Proceedings of the 33rd Annual ACM Symposium on Applied Computing, SAC, ACM, Pau, France, 1952–1959*.
- [27] Boris Motik, Ian Horrocks, and Ulrike Sattler. 2007. Adding Integrity Constraints to OWL. In *Proceedings of the OWLED 2007 Workshop on OWL: Experiences and Directions (CEUR Workshop Proceedings)*, Vol. 258. CEUR-WS.org, Austria.
- [28] Boris Motik, Ian Horrocks, and Ulrike Sattler. 2009. Bridging the gap between OWL and relational databases. *Journal of Web Semantics* 7, 2 (2009), 74–89.
- [29] Natasha Noy, Yuqing Gao, Anshu Jain, Anant Narayanan, Alan Patterson, and Jamie Taylor. 2019. Industry-scale knowledge graphs: lessons and challenges. *Commun. ACM* 62, 8 (2019), 36–43.
- [30] Pouya Ghiasnezhad Omran, Kerry Taylor, Sergio José Rodríguez Méndez, and Armin Haller. 2020. Towards SHACL Learning from Knowledge Graphs. In *Proceedings of the ISWC 2020 Demos and Industry Tracks (CEUR Workshop Proceedings)*, Vol. 2721. CEUR-WS.org, Globally online, 94–99.
- [31] Stefano Ortona, Venkata Vamsikrishna Meduri, and Paolo Papotti. 2018. Robust Discovery of Positive and Negative Rules in Knowledge Bases. In *34th IEEE International Conference on Data Engineering, ICDE 2018, Paris, France, April 16-19, 2018. IEEE Computer Society, Paris, France, 1168–1179*. <https://doi.org/10.1109/ICDE.2018.00108>
- [32] Harshvardhan J. Pandit, Declan O'Sullivan, and Dave Lewis. 2018. Using Ontology Design Patterns To Define SHACL Shapes. In *Proceedings of the 9th Workshop on Ontology Design and Patterns (CEUR Workshop Proceedings)*, Vol. 2195. CEUR-WS.org, Monterey, USA, 67–71.
- [33] Minh-Duc Pham, Linnea Passing, Orri Erling, and Peter A. Boncz. 2015. Deriving an Emergent Relational Schema from RDF Data. In *Proceedings of the 24th International Conference on World Wide Web, WWW 2015, Florence, Italy, May 18-22, 2015. ACM, Florence, Italy, 864–874*. <https://doi.org/10.1145/2736277.2741121>
- [34] Eric Prud'hommeaux, José Emilio Labra Gayo, and Harold R. Solbrig. 2014. Shape expressions: an RDF validation and transformation language. In *Proceedings of the 10th International Conference on Semantic Systems, SEMANTICS 2014, Leipzig, Germany, September 4-5, 2014. ACM, Leipzig, Germany, 32–40*. <https://doi.org/10.1145/2660517.2660523>
- [35] Eric Prud'hommeaux, José Emilio Labra Gayo, and Harold R. Solbrig. 2014. Shape expressions: an RDF validation and transformation language. In *Proceedings of the 10th International Conference on Semantic Systems, SEMANTICS. ACM, Leipzig, Germany, 32–40*.
- [36] Top Quadrant. 2023. TopBraid. <https://www.topquadrant.com/products/topbraid-composer/>. Accessed 17th January, 2023.
- [37] Kashif Rabbani. 2023. Quality Shape Extraction - resources and source code. <https://github.com/dkw-aa/qse>. Accessed 17th January, 2023.
- [38] Kashif Rabbani, Matteo Lissandrini, and Katja Hose. 2022. SHACL and ShEx in the Wild: A Community Survey on Validating Shapes Generation and Adoption. In *Proceedings of the ACM Web Conference 2022. ACM, Online, Lyon, France, 260–263*. <https://www2022.thewebconf.org/PaperFiles/65.pdf>
- [39] Kashif Rabbani, Matteo Lissandrini, and Katja Hose. 2022. SHACL shapes for DBpedia, LUBM, YAGO-4, and WikiData published on Zenodo. <https://doi.org/10.5281/zenodo.5958985>.
- [40] Kashif Rabbani, Matteo Lissandrini, and Katja Hose. 2023. SPARQL queries used for shape extraction in QSE-Exact (query-based) approach. <https://github.com/dkw-aa/qse/tree/main/src/main/resources/queries>. Accessed 17th January, 2023.
- [41] Tomer Sagi, Matteo Lissandrini, Torben Bach Pedersen, and Katja Hose. 2022. A design space for RDF data representations. *VLDB J.* 31, 2 (2022), 347–373. <https://doi.org/10.1007/s00778-021-00725-x>
- [42] Ognjen Savkovic, Evgeny Kharlamov, and Steffen Lamparter. 2019. Validation of SHACL Constraints over KGs with OWL 2 QL Ontologies via Rewriting. In *The Semantic Web - 16th International Conference, ESWC 2019, Portorož (Lecture Notes in Computer Science)*, Vol. 11503. Springer, Slovenia, 314–329.
- [43] Stefan Schmid, Cory Henson, and Tuan Tran. 2019. Using Knowledge Graphs to Search an Enterprise Data Lake. In *The Semantic Web: ESWC 2019 Satellite Events - ESWC (Lecture Notes in Computer Science)*, Vol. 11762. Springer, Portorož, Slovenia, 262–266. https://doi.org/10.1007/978-3-030-32327-1_46
- [44] Juan Sequeda and Ora Lassila. 2021. Designing and Building Enterprise Knowledge Graphs. *Synthesis Lectures on Data, Semantics, and Knowledge* 11, 1 (2021), 1–165.
- [45] Blerina Spahiu, Andrea Maurino, and Matteo Palmonari. 2018. Towards Improving the Quality of Knowledge Graphs with Data-driven Ontology Patterns and SHACL. In *Emerging Topics in Semantic Technologies - ISWC 2018 Satellite Events (Studies on the Semantic Web)*, Vol. 36. IOS Press, Satellite, USA, 103–117.
- [46] Thomas Pellissier Tanon, Gerhard Weikum, and Fabian M. Suchanek. 2020. YAGO 4: A Reason-able Knowledge Base. In *The Semantic Web - 17th International Conference, ESWC (Lecture Notes in Computer Science)*, Vol. 12123. Springer, Heraklion, Crete, Greece, 583–596.
- [47] Jiao Tao, Evren Sirin, Jie Bao, and Deborah L. McGuinness. 2010. Integrity Constraints in OWL. In *Proceedings of the Twenty-Fourth AAAI Conference on Artificial Intelligence, AAAI. AAAI Press, Atlanta, Georgia, USA*.

- [48] Katherine Thornton, Harold Solbrig, Gregory S Stupp, Jose Emilio Labra Gayo, Daniel Mietchen, Eric Prud’Hommeaux, and Andra Waagmeester. 2019. Using shape expressions (ShEx) to share RDF data models and to guide curation with rigorous validation. In *ESWC*. Springer, Cham, 606–620.
- [49] Denny Vrandečić and Markus Krötzsch. 2014. Wikidata: a free collaborative knowledgebase. *Commun. ACM* 57, 10 (2014), 78–85.
- [50] W3C. 2023. SHACL- core constraint components. <https://www.w3.org/TR/shacl/#core-components>. Accessed 17th January, 2023.
- [51] W3C. 2023. W3C: RDF Type. <http://www.w3.org/1999/02/22-rdf-syntax-ns#type>. Accessed 17th January, 2023.
- [52] WESO. 2023. RDFShape. <http://rdfshape.weso.es>. Accessed 17th January, 2023.
- [53] WesoShaclConvert. 2023. SHACL to ShEx converter. <https://rdfshape.weso.es/shaclConvert>. Accessed 17th January, 2023.
- [54] WikiData. 2023. WikiData-2015. <https://archive.org/details/wikidata-json-20150518>. Accessed 17th January, 2023.