

# Using Pregel to create Knowledge Graphs subsets described by non-recursive Shape Expressions

Ángel Iglesias Préstamo<sup>1</sup>[0009-0004-0686-4341] and Jose Emilio Labra Gayo<sup>1</sup>[0000-0001-8907-5348]

WESO Lab - University of Oviedo, Spain

**Abstract.** Knowledge Graphs have been successfully adopted in recent years, existing general-purpose ones, like Wikidata, as well as domain-specific ones, like UniProt. Their increasing size poses new challenges with regards to their practical usage. As an example, Wikidata has been growing the size of its contents and their data since its inception making it difficult to download and process its data. Although the structure of Wikidata items is flexible, it tends to be heterogeneous: the shape of an entity representing a human is distinct from that of a mountain. Recently, Wikidata adopted Entity Schemas to facilitate the definition of different schemas using Shape Expressions, a language that can be used to describe and validate RDF data. In this paper, we present an approach to obtain subsets of knowledge graphs based on Shape Expressions that use an implementation of the Pregel algorithm implemented in Rust. We have applied our approach to obtain subsets of Wikidata and UniProt and present some of these experiments' results.

**Keywords:** Knowledge Graphs · Graph algorithms · RDF · Linked Data · RDF Validation · Shape Expressions · Subsets · Pregel

## 1 Introduction

Knowledge graphs have emerged as powerful tools for representing and organizing vast amounts of information in a structured manner. As their applications continue to expand across various domains, the need for efficient and scalable processing of these graphs becomes increasingly critical.

Creating subsets of knowledge graphs is a common approach to tackle the challenges posed by their size and complexity. Such subsets are essential not only to reduce computational overhead but also to focus on specific aspects of the data.

In this paper, we explore the synergy between two essential concepts in the field of graph processing: Shape Expressions (ShEx) [10] and the Pregel model [12]. Shape Expressions allow to describe and validate knowledge graphs based on the Resource Description Framework (RDF). These expressions have gained significant adoption in prominent projects like Wikidata. On the other

hand, Pregel is a distributed graph processing model designed for efficiently handling large-scale graphs across multiple machines.

Motivated by the need for handling massive graphs in a scalable way, we propose the concept of creating subsets of knowledge graphs using Shape Expressions. By selecting relevant portions of the graph, we can focus computational efforts on specific areas of interest, leading to enhanced efficiency and reduced processing times.

Furthermore, we delve into the capabilities of the Pregel algorithm and its potential for distributed graph processing. We emphasize that the scalability of graph computation can be achieved not only by increasing the number of machines but also by optimizing the use of multi-threading solutions to leverage a single machine’s capabilities. Hence, our solution aims for distributing the problem across multiple threads of a single-node machine. This is, a multi-threaded Pregel. The idea is not only to provide a solution that can run on any hardware efficiently but also to explore the capabilities of Rust for enabling some performance gains regarding single-node computation.

The main contributions of this paper are the following:

1. We present an approach for subset generation of Knowledge Graphs based on Shape Expressions using the Pregel algorithm.
2. We have implemented the previous approach in Rust.
3. We have applied it to generate subsets of Wikidata and UniProt and present some optimizations and results.

The structure of this document is as follows: Section 2 presents the key concepts required for describing the foundations of the problem to be solved. Next, Section 3 explains the most important algorithms for creating Knowledge Graph subsets. After that, in Section 4, the novel approach introduced by this paper is described. Section 5 depicts the experiment for analyzing how the Pregel-based Schema validating algorithm behaves. Next, Section 6 establishes the alternatives and work related to what is presented in the document. And lastly, Section 7 contains the conclusions and future work.

## 2 Background

### 2.1 Knowledge graphs

**Definition 1 (Knowledge Graph [3,5]).** *A Knowledge Graph is a graph-structured data model that captures knowledge in a specific domain, having nodes that represent entities and edges modeling relationships between those.*

Definition 1 is a general and open description of a Knowledge Graph. There are several data models for representing Knowledge Graphs, including Directed edge-labeled and Property Graphs [5], to name a few. In this paper, we will focus on RDF-based Knowledge Graphs, a standardized data model based on directed edge-labeled graphs [5].

**RDF-based Knowledge Graphs** The Resource Description Framework (RDF) is a standard model for data interchange on the Web. It is a W3C Recommendation for representing information based on a directed edge-labeled graph, where labels are the resource identifiers. The idea behind the RDF model is to make statements about things in the form of subject-predicate-object triples. The subject denotes the resource itself, while the predicate expresses traits or aspects of it and expresses a relationship between the subject and the object, another resource. This linking system forms a graph data structure, which is the core of the RDF model. If the dataset represents Knowledge of a specific domain, the Graph will be an RDF-based Knowledge Graph. There are several serialization formats for RDF-based Knowledge Graphs, including Turtle, N-Triples, and JSON-LD. Its formal definition is as follows:

**Definition 2 (RDF-based Knowledge Graph [3]).** *Given a set of IRIs  $\mathcal{I}$ , a set of blank nodes  $\mathcal{B}$  and a set of literals  $\mathcal{L}$ , and RDF-based Knowledge Graph is defined as a triple-based graph  $\mathcal{G} = \langle \mathcal{S}, \mathcal{P}, \mathcal{O}, \rho \rangle$  where  $\mathcal{S} \subseteq \mathcal{I} \cup \mathcal{B}$ ,  $\mathcal{P} \subseteq \mathcal{I}$ ,  $\mathcal{O} \subseteq \mathcal{I} \cup \mathcal{B} \cup \mathcal{L}$ , and  $\rho \subseteq \mathcal{S} \times \mathcal{P} \times \mathcal{O}$ .*

*Example 1 (RDF-based Knowledge Graph of Alan Turing).* Alan Turing (23 June 1912 – 7 June 1954) was employed by the government of the United Kingdom in the course of WWII. During that time he invented the computer for deciphering Enigma-machine-encrypted secret messages, namely, the Bombe computer. Additional information about relevant places where he lived is also annotated, including his birthplace, and the place where he died.

$$\begin{aligned} \mathcal{I} &= \{ \text{alanTuring, wilmslow, GCHQ, unitedKingdom, warringtonLodge, bombe} \\ &\quad \text{town, computer, dateOfBirth, placeOfBirth, employer, placeOfDeath,} \\ &\quad \text{country, manufacturer, instanceOf} \} \\ \mathcal{B} &= \{ \emptyset \} \\ \mathcal{L} &= \{ 23 \text{ June } 1912 \} \\ \rho &= \{ (\text{alanTuring, instanceOf, Human}), \\ &\quad (\text{alanTuring, dateOfBirth, } 23 \text{ June } 1912), \\ &\quad (\text{alanTuring, placeOfBirth, warringtonLodge}), \\ &\quad (\text{alanTuring, placeOfDeath, wilmslow}), \\ &\quad (\text{alanTuring, employer, GCHQ}), \\ &\quad (\text{bombe, discoverer, alanTuring}), \\ &\quad (\text{bombe, manufacturer, GCHQ}), \\ &\quad (\text{bombe, instanceOf, computer}), \\ &\quad (\text{wilmslow, country, unitedKingdom}), \\ &\quad (\text{wilmslow, instanceOf, town}), \\ &\quad (\text{warringtonLodge, country, unitedKingdom}) \} \end{aligned}$$

Through tools like RDFShape it is possible to visualize the RDF-based Knowledge Graph of Alan Turing.

## 2.2 ShEx

Shape Expressions (ShEx) were designed as a high-level, domain-specific language for describing RDF graph structures. The syntax of ShEx is inspired by Turtle and SPARQL, while the semantics are motivated by RelaxNG and XML Schema. In this manner, Shape Expressions specify the requirements that RDF data graphs must fulfill to be considered conformant, they allow systems to establish contracts for sharing information; through a common schema, systems agree that a certain element should appear. This pattern behaves similarly to interfaces in the object-oriented paradigm. Shapes can be specified using a JSON-LD syntax or a human-friendly concise one called ShExC.

*Example 2.* The following ShEx schema describes the **Person** Shape Expression, which is used to validate the RDF-based Knowledge Graph of Alan Turing (see example 1). This example can be found in RDFShape but using ShExC serialization.

$$\begin{aligned}
 \mathcal{L} &= \{ \textit{Person}, \textit{Place}, \textit{Country}, \textit{Organization}, \textit{Date} \} \\
 \delta(\textit{Person}) &= \{ \_ \xrightarrow{\textit{placeOfBirth}} @\textit{Place}, \\
 &\quad \_ \xrightarrow{\textit{dateOfBirth}} @\textit{Date}, \\
 &\quad \_ \xrightarrow{\textit{employer}} @\textit{Organization} \} \\
 \delta(\textit{Place}) &= \{ \_ \xrightarrow{\textit{country}} @\textit{Country} \} \\
 \delta(\textit{Country}) &= \{ \} \\
 \delta(\textit{Organization}) &= \{ \} \\
 \delta(\textit{Date}) &\in \text{xsd:date}
 \end{aligned}$$

## 2.3 Pregel

Pregel (*Parallel, Graph, and Google*) is a data flow paradigm and system created by Google to handle large-scale graphs. Even though the original instance remains proprietary at Google, it was adopted by many graph-processing systems, including Apache Spark. For a better understanding of Pregel, the idea is to *think like a vertex*[11]; this way, the state of a given node, will only depend on those of its neighbors, those nodes connected to it by an outgoing edge (see definition 4). Hence, by *thinking like a vertex*, the problem is divided into several sub-problems: instead of dealing with a huge graph, smaller graphs are the way to go: a vertex and its neighbors. This is especially crucial for huge graphs, with millions of entities and hundreds of relationships for each node.

The series of steps that the Pregel framework follows to process a graph are depicted in Figure 1. The execution starts by sending the initial messages to the vertices at iteration 0. Then, the first – actual – *superstep* begins. In our current implementation, this loop will last until the current iteration is greater than the threshold set at the creation of the Pregel instance. At each iteration, the

vertices will send messages to their neighbors, provided the given direction, and subsequently, they may receive messages sent from other nodes. Moving forward, an aggregation function is applied, and the vertices are updated accordingly. Finally, the iteration counter is incremented, and the next iteration starts.

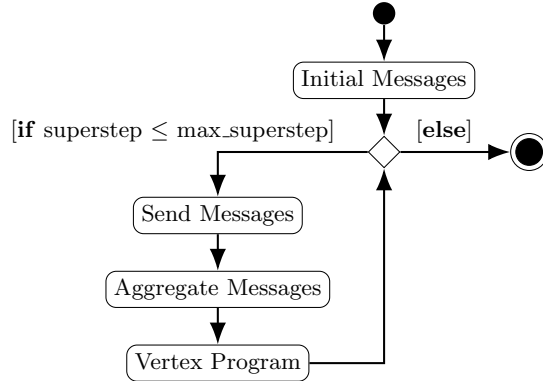


Fig. 1: Pregel model as implemented in `pregel-rs`

### 3 Subsetting approaches

#### 3.1 Knowledge Graph subsets, a formal definition

Given the definition 1 of a Knowledge Graph, the formal definition of a Knowledge Graph subset is seen in definition 3.

**Definition 3 (RDF-based Knowledge Graph subset [3]).** *Given a Knowledge Graph  $\mathcal{G} = \langle \mathcal{S}, \mathcal{P}, \mathcal{O}, \rho \rangle$ , a RDF sub-graph is defined as  $\mathcal{G}' = \langle \mathcal{S}', \mathcal{P}', \mathcal{O}', \rho' \rangle$  such that:  $\mathcal{S}' \subseteq \mathcal{S}$ ,  $\mathcal{P}' \subseteq \mathcal{P}$ ,  $\mathcal{O}' \subseteq \mathcal{O}$  and  $\rho' \subseteq \rho$ .*

*Example 3 (Example of an RDF-based Knowledge Graph subset).* Given the RDF-based Knowledge Graph  $\mathcal{G} = \langle \mathcal{S}, \mathcal{P}, \mathcal{O}, \rho \rangle$  from example 1, the subset of  $\mathcal{G}'$  that contains only the information about Alan Turing's birthplace is defined as follows:

$$\begin{aligned}
 \mathcal{I}' &= \{ \text{alanTuring}, \text{warringtonLodge}, \text{dateOfBirth}, \text{placeOfBirth} \} \\
 \mathcal{B}' &= \{ \emptyset \} \\
 \mathcal{L}' &= \{ 23 \text{ June } 1912 \} \\
 \rho' &= \{ (\text{alanTuring}, \text{dateOfBirth}, 23 \text{ June } 1912), \\
 &\quad (\text{alanTuring}, \text{placeOfBirth}, \text{warringtonLodge}) \}
 \end{aligned}$$

### 3.2 ShEx-based matching generated subsets

ShEx-based matching comprises using a ShEx schema  $\mathcal{S}$  as input, including any nodes whose neighborhood matches any of the shapes from  $\mathcal{S}$  in the produced subset [3]. This approach is used by the `PSchema` algorithm. Having the neighborhood of a node  $s \in \mathcal{S}$  defined as follows:

**Definition 4 (Neighborhood of a node in a Knowledge graph).** *The neighbors of an item  $s \in \mathcal{S}$  in a RDF-based Knowledge graph  $\mathcal{G} = \langle \mathcal{S}, \mathcal{P}, \mathcal{O}, \rho \rangle$  are defined as  $neighbors(s) = \{(s, p, o) : \exists v \in \mathcal{S} \wedge v = (s, p, o)\}$ .*

*Example 4 (Neighborhood of Alan Turing (Q7251)).* Given the RDF-based Knowledge graph  $\mathcal{G} = \langle \mathcal{S}, \mathcal{P}, \mathcal{O}, \rho \rangle$  from example 1, the neighborhood of Alan Turing (Q7251)  $\in \mathcal{S}$  is defined as follows:

$$neighbors(alanTuring) = \{ (alanTuring, instanceOf, Human), \\ (alanTuring, dateOfBirth, 23\ June\ 1912), \\ (alanTuring, placeOfBirth, warringtonLodge), \\ (alanTuring, placeOfDeath, wilmslow), \\ (alanTuring, employer, GCHQ) \}$$

*Example 5 (Example of a ShEx-based matching subgraph).* Given the RDF-based Knowledge Graph  $\mathcal{G} = \langle \mathcal{S}, \mathcal{P}, \mathcal{O}, \rho \rangle$  from example 1 and the ShEx schema  $\mathcal{S}$  defined in example 2, the *ShEx-based matching subgraph* of  $\mathcal{G}$  from  $\mathcal{S}$  is the RDF-based Knowledge graph  $\mathcal{G}'$ , which defined as follows:

$$\begin{aligned} \mathcal{T}' &= \{ alanTuring, wilmslow, GCHQ, unitedKingdom, warringtonLodge, \\ &\quad dateOfBirth, placeOfBirth, employer, country \} \\ \mathcal{B}' &= \{ \emptyset \} \\ \mathcal{L}' &= \{ 23\ June\ 1912 \} \\ \rho &= \{ (alanTuring, dateOfBirth, 23\ June\ 1912), \\ &\quad (alanTuring, placeOfBirth, warringtonLodge), \\ &\quad (alanTuring, employer, GCHQ), \\ &\quad (wilmslow, country, unitedKingdom) \\ &\quad (warringtonLodge, country, unitedKingdom) \} \end{aligned}$$

## 4 Pregel-based Schema validating algorithm

In this section, both the support data structure and the subsetting algorithm are described, including the different steps followed in the Pregel implementation.

#### 4.1 Shape Expression tree

The Shape Expression tree is a hierarchical data structure that represents Shape Expressions in a tree format. Each node in the tree corresponds to a Shape Expression, with the root node being the Shape Expression being validated. Nodes can reference other Shape Expressions, which become its children in the tree.

**Definition 5 (Shape Expression tree).** *Given a Shape Expression  $\mathcal{S}$ , the Shape Expression tree  $\mathcal{T}$  is defined as follows:*

- *If  $\mathcal{S}$  is a terminal **Shape**; that is, it does not reference any other **Shape**, then  $\mathcal{T}$  is a leaf node.*
- *If  $\mathcal{S}$  is a non-terminal **Shape**; that is, it references other **Shapes**, then  $\mathcal{T}$  is an internal node, and its children are the **Shapes** referenced by  $\mathcal{S}$ . Which will be the root nodes of their respective Shape Expression trees.*

*Example 6.* Given the Shape Expression  $\mathcal{S}$  defined in example 2, the Shape Expression tree  $\mathcal{T}$  obtained from  $\mathcal{S}$  was created using the RDFShape and is depicted in Figure 2. **Person** is the root node of  $\mathcal{T}$ , a non-terminal **Shape** that references **Organization**, **Date**, and **Place**. Thus, the children of the root are the **Shapes** referenced by **Person**, which are the root nodes of their respective Shape Expression trees. In the case of the first child, **Organization** is a terminal **Shape**, and thus, it is a leaf node. The same applies to **Date**. However, **Place** is a non-terminal **Shape**, and thus, it is an internal node. Its children are the **Shapes** referenced by it. This representation is recursive, and thus, the **Shapes** referenced by **Place** are the root nodes of their respective ShEx trees.

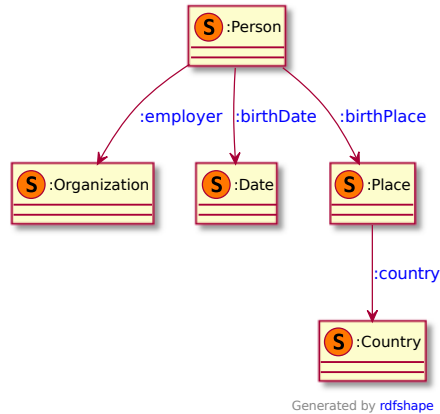


Fig. 2: Example of a Shape Expression tree for the **Person** Shape Expression

The currently supported ShEx language does not support recursion; however, it is planned to implement a solution based on the idea of *unfolding* the recursive schema.

## 4.2 Subsetting algorithm using Pregel and ShEx

PSchema is a Pregel-based algorithm that creates subsets of RDF-based Knowledge Graphs using Shape Expressions. The algorithm's core idea is to validate the nodes of the Shape Expression tree  $\mathcal{T}$  in a bottom-up manner, proceeding from the leaves to the root. The validation is performed in a *reverse level-order traversal* of the tree. The algorithm comprises two main phases: initialization and validation. During the initialization phase, the initial messages are generated and sent to the vertices, while also setting up the *superstep* counter and threshold. This phase establishes the baseline for subsequent steps. In the validation phase, referred to as the *local computation*, the **Shapes** of the tree  $\mathcal{T}$  are validated. The vertices are updated based on the messages they receive from their neighbors. The aforementioned Pregel fork is publicly accessible on Github <sup>1</sup>. For a formal description of the procedure, refer to Algorithm 1.

---

**Algorithm 1:** The PSchema algorithm as implemented in Rust

---

**Input parameters:**

[  $g$ : Graph[ $\mathcal{V}$ ,  $\mathcal{E}$ ]  
 [  $l$ :  $\mathcal{L}$

**Output:**

[ sub: Graph[ $\mathcal{V}$ ,  $\mathcal{E}$ ]

maxIters = see Lemma 1

initialMsgs = None

return Pregel( $g$ , maxIters, initialMsgs, sendMsg, aggMsgs, vProg)

def sendMsg( $l$ :  $\mathcal{L}$ ,  $g$ : Graph[ $\mathcal{V}$ ,  $\mathcal{E}$ ]) = msgs where **foreach**  $l \in \mathcal{L}$

[ msgs.insert  $\left( \begin{array}{l} \left( \begin{array}{ll} \text{validate}(l, g) \text{ if } l = \text{TripleConstraint} & \text{see Algorithm 2} \\ \text{validate}(l, g) \text{ if } l = \text{ShapeReference} & \text{see Algorithm 3} \\ \text{validate}(l, g) \text{ if } l = \text{ShapeAnd} & \text{see Algorithm 4} \\ \text{validate}(l, g) \text{ if } l = \text{ShapeOr} & \text{see Algorithm 5} \\ \text{validate}(l, g) \text{ if } l = \text{Cardinality} & \text{see Algorithm 6} \\ \text{None} & \text{otherwise} \end{array} \right) \end{array} \right)$

def aggMsgs(msgs:  $\mathcal{M}$ ) = msgs where

[ msgs.insert  $\left( \begin{array}{ll} msg & \text{if } msg \neq \text{None} \\ \emptyset & \text{otherwise} \end{array} \right)$

def vProg( $l$ :  $\mathcal{L}$ ,  $g$ : Graph[ $\mathcal{V}$ ,  $\mathcal{E}$ ], msgs:  $\mathcal{M}$ ) = labels.concatenate(msgs)

---

To continue, a formal description of the *validating* algorithms for each of the currently implemented Shapes is provided. The input parameters are simplified. In the actual implementation of the algorithm, it can access the whole Graph and its state.

<sup>1</sup> <https://github.com/weso/pregel-rs>



---

**Algorithm 2:** Validate method for the TripleConstraint Shape

---

**Input parameters:**

```

┌  $l: \mathcal{L}$ 
└  $(-, p, o): (s \in \mathcal{S}, p \in \mathcal{P}, o \in \mathcal{O})$ 

```

**Output:**

```

└  $msg: \mathcal{M}$ 

```

**match**  $l.object$ 

```

┌ case  $Value(v)$ 
└   ┌ if  $p == l.predicate \wedge o == v$  then
      └   └ return  $l$ 
┌ case  $Any$ 
└   ┌ if  $p == l.predicate$  then
      └   └ return  $l$ 

```

---



---

**Algorithm 3:** Validate method for the ShapeReference Shape

---

**Input parameters:**

```

┌  $l: \mathcal{L}$ 
└  $(-, p, o): (s \in \mathcal{S}, p \in \mathcal{P}, o \in \mathcal{O})$ 

```

**Output:**

```

└  $msg: \mathcal{M}$ 

```

**if**  $p == l.predicate \wedge o == l.reference.object$  **then**

```

└   └ return  $l$ 

```

---



---

**Algorithm 4:** Validate method for the ShapeAnd Shape

---

**Input parameters:**

```

┌  $l: \mathcal{L}$ 
└  $(-, p, o): (s \in \mathcal{S}, p \in \mathcal{P}, o \in \mathcal{O})$ 

```

**Output:**

```

└  $msg: \mathcal{M}$ 

```

 $ans \leftarrow \text{true}$ **forall**  $l \in l.shapes$  **do**

```

└   └  $ans \leftarrow ans \wedge \text{validate}(l, g)$ 

```

**if**  $ans$  **then**

```

└   └ return  $l$ 

```

---

---

**Algorithm 5:** Validate method for the ShapeOr Shape
 

---

**Input parameters:**  
 |  $l: \mathcal{L}$   
 |  $(-, p, o): (s \in \mathcal{S}, p \in \mathcal{P}, o \in \mathcal{O})$

**Output:**  
 |  $msg: \mathcal{M}$

$ans \leftarrow \text{false}$   
**forall**  $l \in l.shapes$  **do**  
 |  $ans \leftarrow ans \vee \text{validate}(l, g)$

**if**  $ans$  **then**  
 | **return**  $l$

---



---

**Algorithm 6:** Validate method for the Cardinality Shape
 

---

**Input parameters:**  
 |  $l: \mathcal{L}$   
 |  $(-, p, o): (s \in \mathcal{S}, p \in \mathcal{P}, o \in \mathcal{O})$   
 |  $prevMsg: \mathcal{M}$

**Output:**  
 |  $msg: \mathcal{M}$

$count \leftarrow prevMsg.count()$   
**match**  $l.min$   
 | **case**  $Inclusive(min)$   
 | | **if**  $count \leq min$  **then**  
 | | |  $min \leftarrow \text{true}$   
 | **case**  $Exclusive(min)$   
 | | **if**  $count < min$  **then**  
 | | |  $min \leftarrow \text{true}$

**match**  $l.max$   
 | **case**  $Inclusive(max)$   
 | | **if**  $count \geq max$  **then**  
 | | |  $max \leftarrow \text{true}$   
 | **case**  $Exclusive(max)$   
 | | **if**  $count > max$  **then**  
 | | |  $max \leftarrow \text{true}$

**if**  $min \wedge max$  **then**  
 | **return**  $l$

---

**Lemma 1 (Convergence of the PSchema algorithm).** *Given a Shape Expression tree  $\mathcal{T}$  and a Knowledge Graph  $\mathcal{G}$ , let  $h$  denote the height of  $\mathcal{T}$ ; then the PSchema algorithm is going to converge in  $h$  supersteps. This is, the algorithm is going to validate all the Shapes of  $\mathcal{T}$  in  $h$  supersteps.*

*Example 7 (Example of the subset generated by the PSchema algorithm).* Given the RDF-based Knowledge Graph  $\mathcal{G} = \langle \mathcal{S}, \mathcal{P}, \mathcal{O}, \rho \rangle$  from example 1 and the ShEx schema  $\mathcal{S}$  defined in example 2, the ShEx-based matching subgraph of  $\mathcal{G}$  from  $\mathcal{S}$  is the RDF-based Knowledge graph  $\mathcal{G}'$  from example 5, which is represented in Turtle syntax as follows, refer to RDFShape for more information:

---

```

1 PREFIX :      <http://example.org/>
2 PREFIX xsd:   <http://www.w3.org/2001/XMLSchema#>
3
4 :alan        :placeOfBirth    :warrington          ;
5              :dateOfBirth     "1912-06-23"^^xsd:date ;
6              :employer        :GCHQ                      .
7
8 :warrington  :country          :uk                    .
9
10 :wilmslow   :country          :uk                    .

```

---

### 4.3 Optimizations

**Columnar storage** [1] is a data warehousing technique used in databases and data processing systems. Unlike traditional row-based storage, where data is stored in rows, columnar storage organizes data in columns, grouping values of the same attribute. The columnar storage format offers several advantages. Firstly, it enables better data compression, as similar data types are stored together, reducing the storage footprint. This leads to faster data retrieval and reduced I/O operations when querying specific columns. Columnar stores are advantageous for analytical workloads that involve aggregations or filtering on specific attributes, as it only needs to read the relevant columns. Moreover, columnar storage enhances query performance by leveraging vectorized processing. Modern CPUs can perform operations on entire sets of data (vectors) more efficiently than on individual elements, further improving query speeds.

**Caching** Dictionary encoding [15] is a data compression technique where unique values in a column are assigned numerical identifiers (dictionary indices) and stored in a separate data structure. The actual data in the column is replaced with these compact numerical representations. This technique significantly reduces the storage footprint, especially when columns contain repetitive or categorical data with limited distinct values, as predicates in real-life scenarios.

When combined, columnar storage and dictionary encoding offer several caching benefits. The reduced data size allows more data to be cached within the

same memory capacity, maximizing cache utilization. Additionally, the focused access ensures that only the necessary data is retrieved, further enhancing cache efficiency. In data processing scenarios with repeated patterns and aggregations, caching with columnar storage and dictionary encoding can lead to performance improvements. The cache can hold a large amount of relevant data, minimizing the need for costly disk accesses and accelerating query response times, ultimately resulting in a more efficient and responsive data processing system.

## 5 Experiments and results

Two different Knowledge Graphs are going to be used to test the algorithm, namely, Uniprot<sup>2</sup> and Wikidata. The former is a database that contains information about proteins [2], while the latter is a general-purpose knowledge base having a dump created the 21<sup>st</sup> August 2017. As the serialization format of *Uniprot* is RDF/XML, the `riot` utility from *Apache Jena* is used to convert from RDF/XML to N-Triples. Refer to the examples in the GitHub repository<sup>3</sup>.

As it is seen, the results obtained are similar regarding the size of the subsets. That is, the optimizations have no impact on the *validity* of the tool; this is, the subsets are correct for the **Shape** defined. For this, two Shapes were created during the Japan BioHackathon 2023 [8], namely, `protein` and `subcellular_location`. When comparing the optimized version against its counterpart, the time consumption is reduced by 38% and 35%; while the memory consumption is decreased by 43% and 38%, respectively. Hence, the optimizations are effective both time and memory-wise. The next experiment will focus on how the algorithm behaves when its parameters are modified.

Shape Expression	Initial triples	Resulting triples	Time (s)	Memory (GB)
<code>protein</code>	7,346,129	226,241	23.35	6.74
<code>subcellular_location</code>	7,346,129	1,084,151	57.56	6.04

(a) Execution of the **PSchema** algorithm with no optimization enabled

Shape Expression	Initial triples	Resulting triples	Time (s)	Memory (GB)
<code>protein</code>	7,346,129	226,241	14.58	3.87
<code>subcellular_location</code>	7,346,129	1,084,151	37.76	3.75

(b) Execution of the **PSchema** algorithm with all the optimizations enabled

Fig. 3: Time and memory consumption to create **Uniprot**'s *subsets*

In the second experiment, the number of **Wikidata** entities, the depth, and the width of the **ShEx** tree were modified. The results are depicted in Figure 4. It was observed that the execution time followed a linear trend in all scenarios.

<sup>2</sup> [https://ftp.uniprot.org/pub/databases/uniprot/current\\_release/rdf/](https://ftp.uniprot.org/pub/databases/uniprot/current_release/rdf/)

<sup>3</sup> [https://github.com/angelip2303/pschema-rs/tree/main/examples/from\\_uniprot](https://github.com/angelip2303/pschema-rs/tree/main/examples/from_uniprot)

This indicates that as the number of `Wikidata` entities increased, the execution time increased at a consistent rate. Additionally, the depth and width of the `ShEx` tree influenced the execution time similarly, displaying a linear correlation.

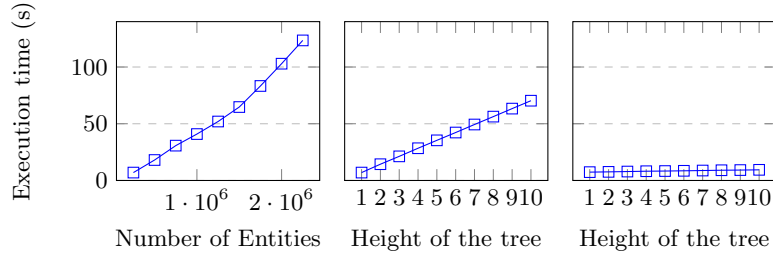


Fig. 4: Time to create the *subsets* of Wikidata with `pschema-rs`

## 6 Related work

### 6.1 Knowledge Graph descriptions

Several Knowledge Graph descriptions have been proposed, with many introduced in [3,5]. Among them, the most relevant ones for this paper are Property graphs, RDF graphs, and Wikibase graphs.

Shape Expressions are used to create the subsets, which were first introduced in 2014. While SHACL (Shapes Constraint Language) is the W3C recommendation<sup>4</sup> since 2019, the Wikidata community has been using Shape Expressions [14] in the entity schemas namespace. This preference stems from the fact that Shape Expressions better adapt to describing data models compared to SHACL which is more focused on constraint violations. A comparison between the two can be found in the following book [9].

### 6.2 Knowledge Graph subsets

Although it is possible to create subsets of the RDF Knowledge Graph through SPARQL construct queries, there are limitations to this approach. Notably, the lack of support for recursion. While proposals to extend SPARQL with recursion have been made [13], such extensions are not widely supported by existing processors. In light of these limitations, a new method using Shape Expressions for creating Knowledge Graph subsets is described in [3]. `PSchema` follows a similar approach to that presented in [16]. However, `SP-Tree` uses SPARQL to query the Knowledge Graph, while `PSchema` uses Shape Expressions and Rust. As such,

<sup>4</sup> <https://www.w3.org/TR/2017/REC-shacl-20170720/>

the **PSchema** algorithm is more flexible than **SP-Tree**, and several optimizations can be applied to it.

The creation of Knowledge graph subsets has gained attention, starting from the 12th International SWAT4HCLS Conference<sup>5</sup>. It has since been selected as a topic of interest in the Elixir Europe Biohackathon 2020<sup>6</sup> and the SWAT4HCLS 2021 Hackathon, which resulted in several publications collecting different approaches [7,6]. This paper is inspired on one of those approaches that was based on Apache Spark but using a new implementation of the Pregel algorithm in Rust.

## 7 Conclusions and future work

In this paper, a novel approach for creating subsets of RDF-based Knowledge Graphs using Shape Expressions and the Pregel framework is presented. The **PSchema** algorithm is described, including the different steps followed in the Pregel implementation. Moreover, the support data structure and the optimizations applied are also described. Two Shapes were created during the Japan BioHackathon 2023 [8] for testing the tool and its validity regarding the optimizations applied. The subsets resulting subsets have the same size in both scenarios. When comparing the optimized version against its counterpart, the time consumption is reduced by 38% and 35%; while the memory consumption is decreased by 43% and 38%, respectively. Hence, the optimizations are effective. The next experiment focused on how the algorithm behaves when its parameters are modified. It was observed that the execution time followed a linear trend in all cases.

In the future, **PSchema** could be extended to support more complex **ShEx** features like recursive Shapes, and an early-prune strategy to reduce the cost of the local computation. The algorithm should receive the **ShEx** schema as an input, rather than programmatically creating descired Shape instance. It is also planned to give support for **WShEx** [4], a **ShEx**-inspired language for describing Wikidata entities, where qualifiers about statements and references can be used for validating purposes. This would allow the algorithm to be used in a wider range of scenarios.

To conclude, **PSchema**, being a Pregel-based Knowledge Graph validating algorithm, allows the processing of large-scale Knowledge Graphs. This is especially relevant in the *Bioinformatics*, where the integration of data from multiple sources is needed. What's more, inference algorithms can be applied to the subsets generated, which is not possible in larger Graphs due to their sizes.

---

<sup>5</sup> [https://www.wikidata.org/wiki/Wikidata:WikiProject\\_Schemas/Subsetting](https://www.wikidata.org/wiki/Wikidata:WikiProject_Schemas/Subsetting)

<sup>6</sup> <https://github.com/elixir-europe/BioHackathon-projects-2020/tree/master/projects/35>

## References

1. Abadi, D.J., Madden, S.R., Hachem, N.: Column-stores vs. row-stores: How different are they really? In: Proceedings of the 2008 ACM SIGMOD International Conference on Management of Data. p. 967–980. SIGMOD '08, Association for Computing Machinery, New York, NY, USA (2008). <https://doi.org/10.1145/1376616.1376712>, <https://doi.org/10.1145/1376616.1376712>
2. Consortium, T.U.: UniProt: the Universal Protein Knowledgebase in 2023. *Nucleic Acids Research* **51**(D1), D523–D531 (11 2022). <https://doi.org/10.1093/nar/gkac1052>, <https://doi.org/10.1093/nar/gkac1052>
3. Gayo, J.E.L.: Creating knowledge graphs subsets using shape expressions (2021). <https://doi.org/10.z8550/ARXIV.2110.11709>, <https://arxiv.org/abs/2110.11709>
4. Gayo, J.E.L.: Wshex: A language to describe and validate wikibase entities (2022)
5. Hogan, A., Blomqvist, E., Cochez, M., d'Amato, C., de Melo, G., Gutiérrez, C., Kirrane, S., Labra Gayo, J.E., Navigli, R., Neumaier, S., Ngonga Ngomo, A.C., Polleres, A., Rashid, S.M., Rula, A., Schmelzeisen, L., Sequeda, J.F., Staab, S., Zimmermann, A.: Knowledge Graphs. No. 22 in Synthesis Lectures on Data, Semantics, and Knowledge, Springer (2021). <https://doi.org/10.2200/S01125ED1V01Y202109DSK022>, <https://kgbook.org/>
6. Hosseini Beghaeiraveri, S.A., Labra-Gayo, J.E., Waagmeester, A., Ammar, A., Gonzalez, C., Slenter, D., Sabah Ul-Hasan Willighagen, E., McNeill, F., Gray, A.: Wikidata subsetting: approaches, tools, and evaluation. *Semantic Web Journal* (Jun 2023), <https://www.semantic-web-journal.net/content/wikidata-subsetting-approaches-tools-and-evaluation-0>
7. Labra-Gayo, J.E., Hevia, A.G., Álvarez, D.F., Ammar, A., Brickley, D., Gray, A.J.G., Prud'hommeaux, E., Slenter, D., Solbrig, H., Beghaeiraveri, S.A.H., et al.: Knowledge graphs and wikidata subsetting (Apr 2021). <https://doi.org/10.37044/osf.io/wu9et>, [biohackrxiv.org/wu9et](https://biohackrxiv.org/wu9et)
8. Labra-Gayo, J.E., Waagmeester, A., Yamamoto, Y., Iglesias-Préstamo, A., Katayama, T., Liener, T., Unni, D., Bolleman, J., Aoki-Kinoshita, K.F., Yokochi, M., et al.: RDF Data integration using Shape Expressions (Jul 2023), <https://biohackrxiv.org/md73k>
9. Labra Gayo, J.E., Prud'hommeaux, E., Boneva, I., Kontokostas, D.: Validating RDF Data. No. 1 in Synthesis Lectures on the Semantic Web: Theory and Technology, Morgan & Claypool Publishers LLC (sep 2017). <https://doi.org/10.2200/s00786ed1v01y201707wbe016>, <https://doi.org/10.2200/s00786ed1v01y201707wbe016>
10. Malewicz, G., Austern, M.H., Bik, A.J., Dehnert, J.C., Horn, I., Leiser, N., Czajkowski, G.: Pregel: a system for large-scale graph processing. In: Proceedings of the 2010 international conference on Management of data. pp. 135–146. New York, NY, USA (2010), <http://doi.acm.org/10.1145/1807167.1807184>
11. McCune, R.R., Weninger, T., Madey, G.: Thinking like a vertex: A survey of vertex-centric frameworks for large-scale distributed graph processing. *ACM Comput. Surv.* **48**(2) (oct 2015). <https://doi.org/10.1145/2818185>, <https://doi.org/10.1145/2818185>
12. Prud'hommeaux, E., Labra Gayo, J.E., Solbrig, H.: Shape expressions: an RDF validation and transformation language. In: Proceedings of the 10th International Conference on Semantic Systems, SEMANTICS 2014. pp. 32–40. ACM (2014)

13. Reutter, J.L., Soto, A., Vrgoč, D.: Recursion in sparql. In: Arenas, M., Corcho, O., Simperl, E., Strohmaier, M., d'Aquin, M., Srinivas, K., Groth, P., Dumontier, M., Heflin, J., Thirunarayan, K., Thirunarayan, K., Staab, S. (eds.) *The Semantic Web - ISWC 2015*. pp. 19–35. Springer International Publishing, Cham (2015)
14. Thornton, K., Solbrig, H., Stupp, G.S., Labra Gayo, J.E., Mietchen, D., Prud'hommeaux, E., Waagmeester, A.: Using shape expressions (shex) to share rdf data models and to guide curation with rigorous validation. In: Hitzler, P., Fernández, M., Janowicz, K., Zaveri, A., Gray, A.J., Lopez, V., Haller, A., Hammar, K. (eds.) *The Semantic Web*. pp. 606–620. Springer International Publishing, Cham (2019)
15. Witten, I.H., Moffat, A., Bell, T.C.: *Managing Gigabytes: Compressing and Indexing Documents and Images*. Morgan Kaufmann Series in Multimedia Information and Systems, Morgan Kaufmann, San Francisco, CA, 2 edn. (1999)
16. Xu, Q., Wang, X., Li, J., Zhang, Q., Chai, L.: Distributed subgraph matching on big knowledge graphs using pregel. *IEEE Access* **7**, 116453–116464 (2019). <https://doi.org/10.1109/ACCESS.2019.2936465>