



HAL
open science

An Algorithm Based on Grammatical Evolution for Discovering SHACL Constraints

Rémi Felin, Pierre Monnin, Catherine Faron, Andrea G. B. Tettamanzi

► **To cite this version:**

Rémi Felin, Pierre Monnin, Catherine Faron, Andrea G. B. Tettamanzi. An Algorithm Based on Grammatical Evolution for Discovering SHACL Constraints. EuroGP 2024 - 27th European Conference on Genetic Programming, Apr 2024, Aberystwyth, United Kingdom. pp.176-191, 10.1007/978-3-031-56957-9_11 . hal-04446252

HAL Id: hal-04446252

<https://inria.hal.science/hal-04446252v1>

Submitted on 8 Feb 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

An Algorithm Based on Grammatical Evolution for Discovering SHACL Constraints

Rémi Felin¹[0000-0003-2532-7555], Pierre Monnin¹[0000-0002-2017-8426],
Catherine Faron¹[0000-0001-5959-5561], and Andrea G. B.
Tettamanzi¹[0000-0002-8877-4654]

Université Côte d’Azur, Inria, CNRS, I3S, Sophia-Antipolis, France
{name.surname}@inria.fr

Abstract. The continuous evolution of heterogeneous RDF data has led to an increase of inconsistencies on the Web of data (i.e. missing data and errors) that we assume to be inherent to RDF data graphs. To improve their quality, the W3C recommendation SHACL allows to express various constraints that RDF data must conform to and detect nodes violating them. However, acquiring representative and meaningful SHACL constraints from complex and very large RDF data graphs is very challenging and tedious. Consequently, several recent works focus on the automatic generation of these constraints. We propose an approach based on grammatical evolution (GE) for extracting representative SHACL constraints by mining an RDF data graph. This approach uses a probabilistic SHACL validation framework to consider the inherent errors in RDF data. The results highlight the relevance of this approach in discovering SHACL shapes inspired by association rule patterns from a real-world RDF data graph.

Keywords: Grammatical Evolution · Shape Mining · Web of Data

1 Introduction

Over the years, the Web has witnessed an increasing publication of heterogeneous data graphs forming the “Web of Data”, aimed at being consumed by humans and artificial agents. These graphs are represented using the RDF (Resource Description Format) standard and queried using the SPARQL standard. RDF relies on a triple model: an RDF triple $\langle s, p, o \rangle$ expresses a relation p between a subject s and an object o . An **RDF data graph** is a set of interconnected RDF triples which terms are IRIs, literals and blank nodes (anonymous resources): s is a resource (a IRI or a blank node), p is a IRI, and o is any RDF term. RDF data can be serialized with several syntaxes among which Turtle is the simplest and most readable one. It uses qualified names with *prefixes* associated to namespaces to simplify the notation of IRIs, e.g. `foaf:name` is parsed into `<http://xmlns.com/foaf/0.1/name>`. The increase of available and heterogeneous RDF data is the result of global initiatives for producing open data,

e.g. Linked Open Data¹. It is well known that this evolution has revealed issues of errors and incompleteness in real-world RDF data graphs and we consider it essential to recognize the principle that these inconsistencies are inherent to RDF data graphs.

To improve the quality of RDF data, the W3C recommended SHACL [11] as a standard to represent constraints on RDF data graphs. **SHACL shapes** are instances of `sh:NodeShape`² that allow targeting a specific set of nodes in an RDF data graph and assessing them against a set of SHACL constraints, i.e. searching possible nodes in the RDF data graph that do not conform to the shape. Overall, the evaluation process considers a shapes graph (i.e., a set of SHACL shapes) and evaluate them against an RDF graph. SHACL thus addresses the requirements for RDF data quality control, contributing to reducing the inherent inconsistencies in RDF data graphs.

We are interested in the extraction of SHACL shapes that express domain constraints from RDF data graphs. Given that SHACL is a relatively new language, real-world data graphs have a minimal set of associated SHACL shapes. This issue, extensively discussed by Rabbati et al. [28], has led the community to explore various research directions. We can distinguish between approaches aiming to extract SHACL shapes from an ontology (i.e., the “schema” associated with RDF graphs) [4, 27], approaches aiming to extract SHACL shapes by mining regularities from RDF facts [3, 8, 25], and approaches combining both [1]. In general, one of the most significant challenges in extracting SHACL shapes is scaling these methods to handle substantial data graphs. Only Fernandez-Álvarez et al. [8] address this problem by optimizing machine memory consumption. Ontology-based approaches are limited to the coverage degree of the ontologies regarding the RDF data graph, which impacts the type of constraints that can be extracted. Some approaches for extracting SHACL shapes from RDF data take errors and incompleteness into account, using a threshold of tolerance [8] or quality measures [25], but they have limitations regarding the types of SHACL shapes they can extract. For instance, the method proposed in [25] is limited to the extraction of a specific type of rules that can be later translated into SHACL shapes. Boneva et al. [3] propose a semi-automatic method to discover SHACL shapes (via a user interface), while Pandit et al. [27] suggest a manual construction of SHACL shapes using Ontology Design Patterns.

We aim to address the limitations regarding the kind of SHACL shapes that can be extracted from an RDF data graph, which motivated the following research question: *How to automatically discover SHACL shapes from RDF data?* We believe that a generative approach for the automatic construction of SHACL candidate shapes using RDF data is one of those that can achieve this ideal. To this end, we propose a mining method based on **Grammatical Evolution (GE)**, which is a particular type of genetic programming. In this approach, it is feasible to automatically generate variable-length expressions in any language [26] using well-defined grammars, composed of production rules. Gram-

¹ <https://lod-cloud.net/>

² Or `<http://www.w3.org/ns/shacl#NodeShape>`

mational Evolution has been the subject of ongoing work in recent years [30]. Recent work discusses some limitations concerning the grammar design [6], their complexity and a “poor” initialisation of individuals [9]. These limitations have been the subject of contributions intended to propose general guidelines for grammar design [23], automated techniques for finding optimal parameters [2], or new techniques for improving the population initialisation [24]. However, the two major and recurring problems with this method are *redundancy* between individuals and low *locality* [13, 16], i.e. “how well neighbouring genotypes correspond to neighbouring phenotypes” [29]. Lourenço et al. have proposed an extension to this approach called Structured Grammatical Evolution (SGE) [12, 14, 15], which enables one-to-one mapping between genes and non-terminals belonging to the grammar, with effective responses to these two problems. Other approaches aim to introduce a probabilistic approach which is based on a probabilistic selection of production rules during the genotype building phase [10, 17]. Mégane et al. extend their approach with a probabilistic SGE [18, 19] to improve their results and resolve these two issues. In these works, some well-known benchmarks are used to assess the effectiveness of the proposed models, like *Santa Fe Trail* and *Boston Housing*, but their application to tasks related to RDF data mining has not been demonstrated yet. Only Nguyen and Tettamanzi have proposed an adaptation of GE for extracting OWL disjointness axioms [21] and complex disjointness axioms [22], with some promising results.

In this paper, we propose an algorithm based on Grammatical Evolution for generating candidate SHACL shapes using a BNF grammar and RDF data as input. The algorithm exploits a probabilistic framework for SHACL validation, which is required given the heterogeneity and incompleteness inherent in open RDF data. The approach is validated by applying it to the extraction of SHACL shapes from a set of RDF data relating to the scientific domain. The remainder of the paper is organized as follows: in Section 2, we present the design of BNF grammars describing candidate SHACL shapes; in Section 3, we present the probabilistic SHACL validation (3.1) and how we use it to define an acceptance measure and a fitness function (3.2) in order to evaluate candidate shapes; in Section 4 we present a recombination operator responding to the redundancy problem and variation operators used to ensure a broad exploration of the solution space; in Section 5, we present the experiments carried out on a real-world RDF dataset. The results of our experiments are presented in Section 5.2, and we conclude with a discussion of future research in Section 6.

2 BNF Grammars of SHACL Shapes

In order to produce and exploit well-formed SHACL shapes as individuals in an evolutionary process, we defined a BNF grammar compliant with the SHACL W3C recommendation [11]. Fig. 1 presents a subset of this BNF grammar to produce shapes targeting nodes of a specified class (`sh:targetClass`) and constraining them to be linked through the predicate `rdf:type` to another specified class. It should be noted that the two classes may be the same. The grammar

provides the phenotypic and genotypic characterisation for each individual using a set of *static rules* and *dynamic rules*. The dynamic rules system, proposed by Nguyen and Tettamanzi [21, 22], allows the mapping between rules and RDF data. They wrote BNF grammars to build candidate OWL axioms and exploit them with an evolutionary algorithm based on GE. The static rules are the *immutable* components of the phenotypic character whereas the dynamic rules are the *problem instance-dependent* components of the phenotypic character, where each rule has one or many possible values depending on the considered RDF dataset, and each value is identified by a genotype. However, their dynamic rules were hard-coded in their system. In contrast, we extended the dynamic rules design by directly enabling the user to write embedded SPARQL queries as value of one or more production rules in the grammar in order to perform the mapping with the desired granularity.

```

1 <Shape>      := "a " <NodeShape>
2 <NodeShape> := "sh:NodeShape; " <ShapeBody>
3 <ShapeBody> := "sh:targetClass " <Class> "; " <ShapeProp>
4 <ShapeProp> := "sh:property [ " + <PropBody> " ] ."
5 <PropBody>  := "sh:path rdf:type; sh:hasValue " <Class> " ;"
6 <Class>     := "SPARQL ?x rdf:type ?Class"

```

Fig. 1: An extract of the BNF grammar for SHACL shapes

To illustrate, in Fig. 1, the `<Class>` non-terminal (used in lines 3 and 5) is defined by a dynamic rule to extract all possible classes from the RDF dataset using a SPARQL query. The keyword `SPARQL` in line 6 is used to specify the query graph pattern to be matched on RDF data; the result of this query is the set of nodes in the RDF data graph \mathcal{C} bound to variable `?Class` in the query graph pattern: $\mathcal{C} = \{c_i, i \in [1, n]\}$. The final step is to replace the initial value of the rule `<Class>`, i.e. the SPARQL query "SPARQL ?x a ?Class", by the SPARQL results \mathcal{C} , i.e. $c_1 \mid c_2 \mid \dots \mid c_n$. The whole process is presented in Fig. 2.

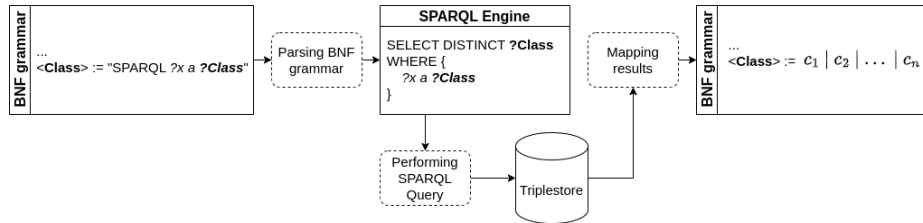


Fig. 2: Dynamic rules process based on the BNF grammar presented in Fig. 1

Using the BNF grammar presented in Fig. 1, the genotype of an individual is a pair of codons $[i, j]$, which are decoded into two classes from the dataset using

a classic genotype-phenotype mapping, and produce the following phenotype structure:

```
"a sh:NodeShape ; sh:targetClass c_i ; sh:property [ sh:path rdf:type ;
sh:hasValue c_j " ; ] ."
```

It is noteworthy that the proposed grammar can be extended to produce a wider array of SHACL shapes using a variable-length template, e.g. replacing the rule `<ShapeProp>` from Fig. 1 by `<ShapeProp> := <Prop> <ShapeProp> | <Prop>` where `<Prop> := "sh:property [" + <PropBody> "] ."`. Such an extended grammar would produce SHACL shapes specifying one or more constraints (depending on the chosen length).

3 Probabilistic SHACL Validation as a Fitness Function

3.1 Preliminaries

We rely on the probabilistic framework for SHACL validation proposed by Felin et al. [7] to assess SHACL shapes considering inherent inconsistencies from RDF data. It extends the standard evaluation of RDF data against SHACL shapes by considering a *physiological* error proportion p in the RDF data, i.e. a possible acceptable violation rate. *Physiological* errors in RDF data graph are inherent errors whose origins can be diverse, e.g. from collaborative building of large RDF data graphs (e.g. Wikidata) or automatically constructed RDF data graphs (e.g. DBpedia). In the rest of this section, we summarize the principles of this model.

Considering an RDF data graph v , the *support* of a shape s , v_s , is the set of RDF triples in v targeted by s (and therefore tested during the validation). The *reference cardinality* (`refCard`) of s is the cardinality of its support: $|v_s|$. The *confirmations* and *violations* of s , respectively v_s^+ and v_s^- , are the sets of triples that, respectively, are consistent with s and violate s : $v_s^+ \cap v_s^- = \emptyset$, and $v_s = v_s^+ \cup v_s^-$. The probabilistic model for SHACL validation relies on a binomial distribution $X \sim B(|v_s|, p)$ where p is the physiological error proportion. The *likelihood* to observe a number of violations $|v_s^-|$ in the support of a shape s considering $X \sim B(|v_s|, p)$ is defined as follows:

$$L_{|v_s^-|} = P(X = |v_s^-|) = \binom{|v_s|}{|v_s^-|} \cdot p^{|v_s^-|} \cdot (1-p)^{|v_s^+|} \quad (1)$$

The *acceptance* of a shape s depends on the proportion of violations for s , i.e. $\hat{p}_s = \frac{|v_s^-|}{|v_s|}$: s is consistent with v if \hat{p}_s is smaller than the theoretical violation proportion p :

$$\hat{p}_s \leq p \implies KG \models s \quad (2)$$

In the case where $\hat{p}_s > p$, a Chi-Square Goodness of Fit test is performed using the *test statistic* X_s^2 defined as follows:

$$X_s^2 = \sum_{i=1}^k \frac{(n_i - T_i)^2}{T_i} \sim \chi_{k-1; \alpha}^2 \quad (3)$$

where k is the total number of groups, i.e. $k = 2$, n_i is the observed number of individuals and T_i is the theoretical number of individuals. The acceptance of the null hypothesis H_0 implies the acceptance of s :

$$X_s^2 \leq \chi_{k-1;\alpha}^2 \implies KG \models s \quad (4)$$

3.2 Acceptability and Fitness Score

In order to iteratively produce a final population of SHACL shapes expressing some domain constraints that are implicit in an RDF dataset, we propose a fitness function based on an acceptability measure of a shape combined with the probabilistic framework presented.

The acceptability of a SHACL shape s , $A(s)$, regarding an RDF dataset, depends on the error rate when validating this RDF dataset against s , which depends on the theoretical error proportion p . The acceptability measure $A(s) \in [0, 1]$ of a SHACL shape is defined by:

$$A(s) = \begin{cases} 1 & \text{if } \hat{p}_s \leq p \text{ or } X_s^2 \leq \chi_{k-1;\alpha}^2 \quad (\text{Equation (2) and 4}), \\ \frac{L_{|v_s^-|}}{P(X=|v_s| \times p)} & \text{otherwise} \quad (\text{Equation (1)}). \end{cases} \quad (5)$$

In the computation of $A(s)$, for the hypothesis testing, we consider a margin error $\alpha = 0.05$. Therefore, the *critical value* is $\chi_{k-1;0.05}^2 = 3.84$. When $A(s) = 1$, s is acceptable and is (probably) selected as one of the most fit individuals, in the sense that it captures some domain knowledge extracted from the RDF data.

In the case where the null hypothesis is rejected, $A(s) \neq 1$ and so s is not acceptable but it may be considered in the grammatical evolution algorithm for crossover or mutation operations. For this purpose, $A(s)$ is equal to the likelihood of s normalized by the maximal value of the probability mass function for a binomial distribution $X \sim B(|v_s|, p)$. It ensures a better distribution of $A(s)$ values between 0 and 1 in contrast to the lonely likelihood value $L_{|v_s^-|}$ and therefore avoids excessively penalising individuals who are "close" to being acceptable but for whom the likelihood is very low.

The fitness function of a SHACL shape s , $F(s)$, regarding an RDF dataset, combines its acceptability $A(s)$ and the cardinality of its confirmations $|v_s^+|$:

$$\forall s \in P, \quad F(s) = |v_s^+| \times A(s) \quad (6)$$

4 Variation and Recombination Operators

In this paper, we adapt the main components of the GE variation operators to discover SHACL shapes over RDF data, considering the problem of *redundancy* and *low locality* presented as the main issues of GE by the community.

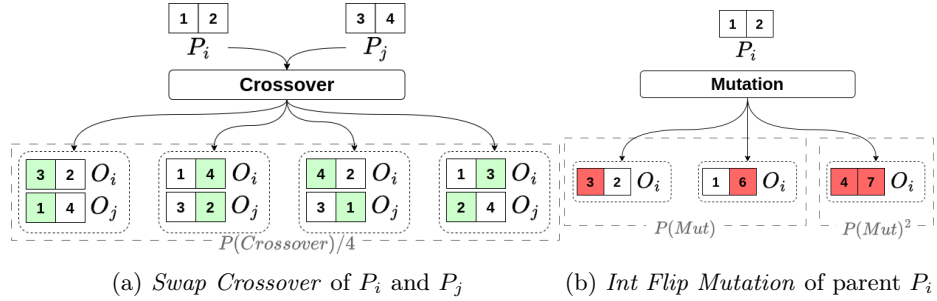


Fig. 3: Representation of GE operators and their probabilities of occurrence

The *redundancy* is observed when many genotypes map the same phenotype expression [13]. Based on this fact, we adapt the recombination phase to filter every offspring using a **phenotypic comparison**: Algorithm 1 presents the recombination of selected individuals \mathcal{S} among the whole population \mathcal{P} (\mathcal{E} represents the elite individuals). Line 11 describes the conditions for integrating an offspring i into the replacement population \mathcal{R} : i is integrated into \mathcal{R} if the phenotypic expression of i is not already observed among the elitist individuals \mathcal{E} and the replacement individuals \mathcal{R} . As a consequence, we avoid the reflection of the redundancy in the final population: $\forall i \in \mathcal{P}, \nexists j \in \mathcal{P} : i \equiv j$.

Algorithm 1: Recombination of a population \mathcal{P}

Data: elite individuals \mathcal{E} and selected individuals \mathcal{S}
Result: replacement population \mathcal{R}

```

1  $\mathcal{R} \leftarrow \{\}$ ;
2 while  $|\mathcal{R}| \neq |\mathcal{P}| - |\mathcal{E}|$  do
3    $\mathcal{C} \leftarrow \{\}$ ;
4    $p_1 \leftarrow \mathcal{S}[\text{random}() * |\mathcal{S}|]$ ;
5    $p_2 \leftarrow \mathcal{S}[\text{random}() * |\mathcal{S}|]$ ;            $\triangleright \text{random}() * |\mathcal{S}|$  as integer;
6   if  $p_1 \neq p_2$  then
7      $\mathcal{C} \leftarrow p_1 \cup p_2$ ;
8      $\mathcal{C} \leftarrow \text{crossover}(\mathcal{C})$ ;            $\triangleright$  Figure 3a;
9      $\mathcal{C} \leftarrow \text{mutation}(\mathcal{C})$ ;        $\triangleright$  Figure 3b;
10    for  $i \in \mathcal{C}$  do
11      if  $i \notin \mathcal{E} \cup \mathcal{R}$  and  $|\mathcal{R}| \neq |\mathcal{P}| - |\mathcal{E}|$     $\triangleright$  Phenotypic comparison
12        then
13           $\mathcal{R} \leftarrow \mathcal{R} \cup i$ ;
14    end
15 end
16 return  $\mathcal{R}$ 

```

In this context, the *locality* issue is dependent on the neighbourhood between the selected rules from parent to offspring. Consequently, some results from the variation operators presented in Fig. 3 can lead to a low locality, i.e. a very different offspring, but also to a fairly strong locality. Considering the grammar presented in Fig. 1, a modification of the first codon (impacting the value of the `sh:targetClass` c_i) significantly changes the meaning of the phenotypic trait, as SHACL validation is performed on the nodes instantiating c_i : a new production rule c'_i replacing c_i leads to a locality as low as the proximity (e.g. common instances) between them. The modification of the last codon has a lower impact on SHACL validation since the targeted nodes are the same, resulting in a fairly strong locality (even if the meaning of the phenotype is different).

5 Experiences

5.1 Shape Mining over the Covid-on-the-Web Dataset

Setup *Covid-on-the-Web*³ [20] is an RDF dataset produced from the *COVID-19 Open Research Dataset (CORD-19)*. It describes articles and named entities identified in these articles and linked to *Wikidata* entities. We consider a subset containing 18.79% of the articles and 0.01% of the named entities. This dataset contains 226,647 RDF triples, 20,912 distinct articles and 6,331 distinct named entities.

We consider the mining of SHACL shapes representing **association rules** between *Wikidata entities*⁴, i.e., rules of the form $\mathcal{X} \rightarrow \mathcal{Y}$. We use the BNF grammar presented in Fig. 1 to generate candidate shapes. Each candidate involves a first Wikidata entity, i.e., the *antecedent* called \mathcal{X} , and targets nodes n (i.e., scientific articles) typed by this entity using the `sh:targetClass` property. The proposed constraint verifies if these nodes are also typed by a second Wikidata entity, i.e., the *consequent* \mathcal{Y} , using the `sh:hasValue` constraint applied on the `rdf:type` property.

Concerning the aim of these experiments, we considered the diversity and the discovery of a large set of acceptable shapes as the most important aspects. While it is obvious that a resource-intensive parameter setting (high population size, high effort, etc.) would provide the best results at the cost of a long computation time, we have focused on a reasonable set of parameters in order to minimize the computation time invested.

We used an implementation of the presented algorithm combined with the probabilistic SHACL validation engine implemented in the *Corese* semantic Web factory [5]. We considered a theoretical error proportion $p = 0.5$ (i.e., *physiological error*), according to the experiment’s results on this dataset reported by Felin et al. [7]: this value p maximises the mean value of the likelihood measure L (see Equation (1)).

³ <https://github.com/Wimmics/CovidOnTheWeb>

⁴ https://www.wikidata.org/wiki/Wikidata:Main_Page

The experiments have been performed on a server equipped with an Intel(R) Xeon(R) CPU E5-2637 v2 processor at 3.50GHz clock speed, with 172 GB of RAM, 1 TB of disk space running under the Ubuntu 20.06.4 LTS 64-bit operating system.

Recall In order to assess the ability of our approach to find acceptable candidates in the solution space of our problem, we define the recall R of our algorithm. The recall provides the rate of distinct solutions found by our algorithm among the total number of solutions \mathcal{A} . Let Ω be the set of all possible (acceptable and non-acceptable) pairs $(\mathcal{X}, \mathcal{Y})$ of distinct named entities extracted from the *Covid-on-the-web dataset*. We may compute its cardinality using the number of distinct entities in the dataset: $|\Omega| = 6,331 \times 6,330 = \mathbf{40,075,230}$.

To estimate the number of acceptable shapes in Ω , we sample a random subset $\Omega' \subseteq \Omega$, representative of the solution space. To ensure the representativity of Ω' , we determine its minimal size using the *Cochran* formula: $|\Omega'| = \frac{z^2 \times p \times (1-p)}{m^2} = \frac{2.58^2 * 0.5^2}{0.02^2} \approx \mathbf{4,161}$, where z is the standard normal z-table with a confidence level of 99% (so $z \approx 2.58$), m is the tolerated margin of error (2%); p' the probability that the candidate shape is acceptable (unknown in this context so $p' = 0.5$).

Consequently, we generated 4,161 distinct and random shapes which have been evaluated over the *Covid-on-the-Web* subgraph with probabilistic SHACL validation: the results show that only 2 shapes in Ω' are accepted, *i.e.*, 0.05% of the total. This allows us to estimate the total number of acceptable shapes \mathcal{A} : $|\mathcal{A}| = |\Omega| \times 0.0005 = \mathbf{20,037.6}$.

Finally, the recall of our algorithm $R(x)$, *i.e.*, how well the x acceptable shapes cover the solution space, is defined by:

$$\mathbf{R}(x) = \frac{x}{|\mathcal{A}|} \times 100 \quad (7)$$

5.2 Results

$|\mathcal{P}|/E$ choice We assessed our approach with manually defined small population sizes ($|\mathcal{P}|$) and quite low *effort* values E and we analysed the effects of the ratio $|\mathcal{P}|/E$. This corresponds to verifying if our algorithm can find credible and surprising candidate shapes using a minimum investment of CPU time. Consequently, we performed 10 executions of our algorithm using the different parameter settings presented in Table 1 (90 in total) and analyzed the final whole population \mathcal{P} and the final elitist subset \mathcal{E} ($\mathcal{E} \subseteq \mathcal{P}$). Each configuration has been assessed regarding the following metrics: the average fitness value $\bar{\mathbf{F}}$; the average rate of accepted shapes $\% \mathbf{A}$; the average likelihood value $\bar{\mathbf{L}}$; the average CPU time (in ms) for evaluating an individual $\bar{\mathbf{T}}$ and the average recall $\bar{\mathbf{R}}$.

According to the results presented in Table 2, a gradual increase of the effort E tends to enhance the global quality of candidate shapes into \mathcal{P} : all the metrics related to the individual quality (\bar{F} , $\%A$, \bar{L} and \bar{R}) are significantly better

Table 1: Used parameters to analyse the impact of $|\mathcal{P}|/E$ choice.

Parameters	Value(s)
GE	
$ \mathcal{P} $	{100; 200; 500}
Effort (E)	{5,000; 10,000; 20,000}
% Selection (\mathcal{E})	20%
% Selection (\mathcal{R})	40%
Selection type	Tournament
% Tournament	25%
Crossover type - P	Swap (Fig. 3a) - 75%
Mutation type - P	Int Flip (Fig. 3b) - 5%
Probabilistic SHACL	
Confidence level α	5%
Theoretical inconsistencies proportion p	50%

Table 2: Results obtained using the parameters presented in Table 1. Best result for each metric is in bold and second best underlined. Highlighted columns are the best.

	$ \mathcal{P} = 100$			$ \mathcal{P} = 200$			$ \mathcal{P} = 500$			
	E = 5,000	E = 10,000	E = 20,000	E = 5,000	E = 10,000	E = 20,000	E = 5,000	E = 10,000	E = 20,000	
From \mathcal{P}	\overline{F}	0.79 ± 1.4	0.64 ± 1.05	1.52 ± 1.7	0.53 ± 0.81	1.24 ± 1.83	<u>1.51 ± 1.14</u>	0.9 ± 0.97	1.07 ± 1	1.3 ± 0.87
	%A	1.9 ± 2.77	4 ± 2.71	10.1 ± 4.18	2 ± 0.75	3.3 ± 2.15	<u>8.05 ± 3.11</u>	1.56 ± 0.76	2.36 ± 1.05	4.98 ± 2.33
	\overline{L}	2.55 ± 1.77	4.86 ± 1.66	6.74 ± 2.07	2.42 ± 0.97	4.41 ± 1.85	<u>6.18 ± 1.13</u>	1.25 ± 0.29	2.07 ± 0.59	4.47 ± 1.29
	\overline{R}	18 ± 2.73	<u>17.98 ± 2.96</u>	21.94 ± 3.73	16.84 ± 2.4	19.22 ± 3.03	<u>21.14 ± 3.65</u>	19.21 ± 2.28	18.39 ± 3.71	19.31 ± 3.71
From \mathcal{E}	\overline{F}	3.91 ± 7	2.77 ± 5.28	7.41 ± 8.48	2.65 ± 4.06	6.11 ± 9.15	7.41 ± 5.74	4.41 ± 4.9	5.29 ± 4.99	<u>6.36 ± 4.39</u>
	%A	9.5 ± 13.83	18.5 ± 12.92	49.5 ± 19.64	9.5 ± 4.22	15.75 ± 11.31	<u>39.5 ± 15.27</u>	6.7 ± 3.13	11.2 ± 4.94	24.1 ± 10.99
	\overline{L}	9.36 ± 6.74	16.18 ± 4.47	21.51 ± 4.88	8.21 ± 3.23	14.96 ± 5.84	<u>20.73 ± 2.25</u>	4.03 ± 1.1	7.22 ± 1.86	15.19 ± 3.67
	\overline{R}	10.6 ± 2.45	9.35 ± 1.43	7.68 ± 1.26	9.22 ± 1.28	7.53 ± 1.11	6.65 ± 0.65	10.64 ± 2.83	7.81 ± 0.94	<u>7.21 ± 2.66</u>

regardless of the population size $|\mathcal{P}|$. This is clearer regarding the elitist part of the population \mathcal{E} .

Globally, it appears that the smallest $|\mathcal{P}|$ with a high effort provide the best results: the results obtained with $(|\mathcal{P}| = 100; E = 20,000)$ and $(|\mathcal{P}| = 200; E = 20,000)$ are very similar, except the proportion of acceptable shapes in \mathcal{E} (respectively 49.5% and 39.5%). For each metric, we performed a *Mann-Whitney-Wilcoxon* (Table 3) test to highlight any differences between each obtained result: it appears that only the average recall \overline{R} values (from \mathcal{P} and \mathcal{E}) are significantly different between them (< 0.05) which suggests that the choice $(|\mathcal{P}| = 200; E = 20,000)$ is the best one for this measure.

Selection (\mathcal{R}) pressure We assume that the smallest population combined with the highest effort, *i.e.*, $(|\mathcal{P}| = 100; E = 20,000)$, is the best choice for analysing selective pressure and learning about its impact on metrics. Conse-

Table 3: *Mann-Whitney-Wilcoxon* test: comparison between the results obtained for $(|P| = 100; E = 20,000)$ and $(|P| = 200; E = 20,000)$ with $\alpha = 0.05$.

From \mathcal{P}		From \mathcal{E}	
Metrics	P-value	Metrics	P-value
\overline{F}	0.528	\overline{F}	0.529
$\%A$	0.198	$\%A$	0.210
\overline{L}	0.684	\overline{L}	0.796
\overline{T}	0.631	\overline{T}	0.076
\overline{R}	0.037	\overline{R}	0.028

quently, we have studied different selection types: *Scaled Roulette Wheel* and *Tournament* with the different settings presented in Table 4.

Table 4: Parameters used to analyse the impact of the selective pressure on \mathcal{R} .

Parameters	Value(s)
GE	
$ P $	100
Effort (E)	20,000
% Selection (\mathcal{E})	20%
Selection type	{Scaled Roulette Wheel; Tournament}
% Selection (\mathcal{R})	{20%; 40%; 60%}
% Tournament	{10%; 25%; 50%}
Crossover type - P	Swap (Fig. 3a) - 75%
Mutation type - P	Int Flip (Fig. 3b) - 5%
Probabilistic SHACL	
Confidence level α	5%
Theoretical inconsistencies proportion p	50%

The results obtained with the *Scaled Roulette Wheel* selection are presented in Table 5. They highlight that the metrics are enhanced with a high selection rate ($S = 60\%$) even though no very good candidates have been found. It appears that a high selection rate enhances the exploration of the solution space, resulting in a relatively strong difference between the results observed for \mathcal{P} and the elitist subset \mathcal{E} .

The results obtained with the *Tournament* selection are presented in Table 6 and we identify the same trend. For $S = 20\%$, the global difference of results between \mathcal{P} and \mathcal{E} is quite low, enhancing the homogeneity of the population \mathcal{P} whereas a high selection rate ($S = 60\%$) reflects the heterogeneity of the population because the global difference of results between \mathcal{P} and \mathcal{E} is high. We can see from the difference between the average time for the whole population

Table 5: Results obtained using the *Scaled Roulette Wheel* selection and parameters presented in Table 4: best result for each metric is in bold and second best underlined.

	$S = 20\%$	$S = 40\%$	$S = 60\%$	
\mathcal{P}	\overline{F}	<u>0.86 ± 1.06</u>	1.78 ± 2.99	0.56 ± 0.16
	$\%A$	<u>8.7 ± 3.83</u>	7.3 ± 5.33	11.7 ± 3.06
	\overline{L}	7.85 ± 2.46	8.85 ± 1.89	<u>8.25 ± 2.6</u>
	\overline{T}	20.25 ± 6.26	<u>19.86 ± 8.61</u>	19.83 ± 5.78
\overline{R}	<u>0.04 ± 0.02</u>	<u>0.04 ± 0.03</u>	0.06 ± 0.02	
\mathcal{S}	\overline{F}	<u>4.27 ± 5.29</u>	8.87 ± 14.97	2.74 ± 0.8
	$\%A$	<u>43.5 ± 19.16</u>	36 ± 25.47	58 ± 14.94
	\overline{L}	<u>24.33 ± 4.36</u>	25.25 ± 5.5	23.94 ± 4.31
	\overline{T}	8.49 ± 1.29	<u>7.77 ± 1.46</u>	7.38 ± 1.19
\overline{R}	<u>0.04 ± 0.02</u>	0.04 ± 0.03	0.06 ± 0.02	

vs the average time for the elite that a high selection rate S tends to make the population vary considerably while maintaining a very good elite population which ensures shapes of good quality in the elite with a wide exploration in the global population, favouring the discovery of heterogeneous and potentially interesting shapes.

Table 6: Results obtained using the *Tournament* selection and parameters presented in Table 4: best result for each metric is in bold, second best underlined. The highlighted column corresponds to the *reference* results presented in Table 2.

	$S = 20\%$			$S = 40\%$			$S = 60\%$			
	$Tour = 10\%$	$Tour = 25\%$	$Tour = 50\%$	$Tour = 10\%$	$Tour = 25\%$	$Tour = 50\%$	$Tour = 10\%$	$Tour = 25\%$	$Tour = 50\%$	
\mathcal{P}	\overline{F}	1.78 ± 1.93	<u>1.82 ± 2.44</u>	1.08 ± 1.98	1.2 ± 1.54	1.52 ± 1.7	1.71 ± 2.42	0.99 ± 0.83	2.32 ± 2	1.78 ± 2.33
	$\%A$	11.1 ± 5.2	8.9 ± 5.45	8.1 ± 5.43	9.9 ± 4.51	10.1 ± 4.18	9.7 ± 4.27	8.3 ± 3.97	<u>10.7 ± 4.99</u>	8.4 ± 3.27
	\overline{L}	5.55 ± 2.6	6.84 ± 2.81	<u>6.98 ± 2.58</u>	6.11 ± 1.14	6.74 ± 2.07	6.03 ± 2.02	6.12 ± 1.32	6.36 ± 2.1	7.1 ± 1.32
	\overline{T}	22.92 ± 9.12	16.4 ± 3.4	<u>18.88 ± 6.41</u>	23.93 ± 5.93	21.94 ± 3.73	21.66 ± 7.37	28.44 ± 7.5	27.89 ± 6.46	27.9 ± 6.62
\overline{R}	0.06 ± 0.03	0.04 ± 0.03	0.04 ± 0.03	<u>0.05 ± 0.02</u>	<u>0.05 ± 0.02</u>	<u>0.05 ± 0.02</u>	0.04 ± 0.02	<u>0.05 ± 0.03</u>	0.04 ± 0.02	
\mathcal{S}	\overline{F}	8.83 ± 9.64	<u>9.02 ± 12.21</u>	5.34 ± 9.94	5.93 ± 7.69	7.41 ± 8.48	8.47 ± 12.13	4.9 ± 4.14	11.55 ± 10.02	8.86 ± 11.68
	$\%A$	55.5 ± 25.98	43.5 ± 26.46	39.5 ± 25.65	49 ± 21.96	49.5 ± 19.64	48.5 ± 21.35	41 ± 19.12	<u>52.5 ± 24.41</u>	42 ± 16.36
	\overline{L}	18.96 ± 6.49	19.1 ± 3.73	21.6 ± 4.95	21.22 ± 2.5	21.51 ± 4.88	19.77 ± 3.69	22.68 ± 4.24	<u>24.25 ± 5.67</u>	24.46 ± 3.09
	\overline{T}	7.54 ± 1	8.09 ± 1.07	7.76 ± 0.79	7.4 ± 1.52	7.68 ± 1.26	7.25 ± 1.51	6.93 ± 0.43	<u>7.06 ± 1.18</u>	8.02 ± 1.38
\overline{R}	0.06 ± 0.03	0.04 ± 0.03	0.04 ± 0.03	<u>0.05 ± 0.02</u>	<u>0.05 ± 0.02</u>	<u>0.05 ± 0.02</u>	0.04 ± 0.02	<u>0.05 ± 0.02</u>	0.04 ± 0.02	

Acceptable shapes Among all the conducted experiments, we have discovered a set of 1,766 distinct and acceptable shapes. An overview of these results is provided in Table 7. Some of these shapes have been accepted with a very high violation rate ($> 50\%$) but can be easily validated, e.g. the candidate implying the following rule (**gene expression profiling** \rightarrow **gene expression**) is easily understandable and acceptable even if it implies 52.6% of violations. Moreover, 46.38% of the whole has been accepted after performing hypothesis testing which

shows it has a valuable impact on the acceptance of shapes and so on the mining process.

However, some of these shapes require final validation from experts due to the complexity, e.g. the following rule (`chemokine` \rightarrow `cytokine`) has been automatically accepted and validated after some research: “*Chemokines [...] are a family of small cytokines*”⁵. Rule (`tlr9` \rightarrow `toll-like receptor`) has been accepted but must be validated by experts.

Table 7: Overview of the distinct and acceptable shapes discovered from all the performed experiments.

Metrics	\bar{F}	\bar{L}	\bar{T}	\bar{R}
Values	19.49	19.14	7.93	8.91

The discovery of very good candidates from some of these experiments impacts the standard deviation of many values \bar{F} and $\bar{\%A}$ (some of these are higher than the mean value). We also note some trivial shapes with identical classes for the `sh:targetClass` and the constraint `sh:hasValue` which implies a perfect acceptance of these candidates, *i.e.*, without any violations, and so a very good fitness value. These can be generated because of production rules selection (with a modulo operator) and a *quasi-infinite* range for codon definition. However, this is a fairly rare occurrence: we observe it among (only) 132 candidates from the 1,766 acceptable shapes, *i.e.*, 7.47%. We suggest accepting a low occurrence of these shapes being discovered (even if they are uninteresting) in order to avoid any negative impact on the exploration of the solution space.

The \bar{T} value presented in Table 7 and the correlation between the number of violations and the CPU time presented in Fig. 4 suggest that the CPU time required to invest in an evolutionary process is maximal at the beginning, then decreases as the average number of violations decreases (and therefore when the shapes become more and more acceptable). Considering this expected evolution and that the average time is low, this demonstrates the relevance of this evolutive approach to the discovery of SHACL shapes over an RDF data graph and appears to be suitable for scalability.

6 Conclusion

In this paper, we have proposed a framework using a grammatical evolution method to extract candidate SHACL shapes from an RDF dataset based on a manually defined BNF grammar. The proposed algorithm provides an effective response to the redundancy problem whereas the low locality appears to be *problem-dependent*. Additionally, the generative evolution allows to tackle the

⁵ <https://en.wikipedia.org/wiki/Chemokine>

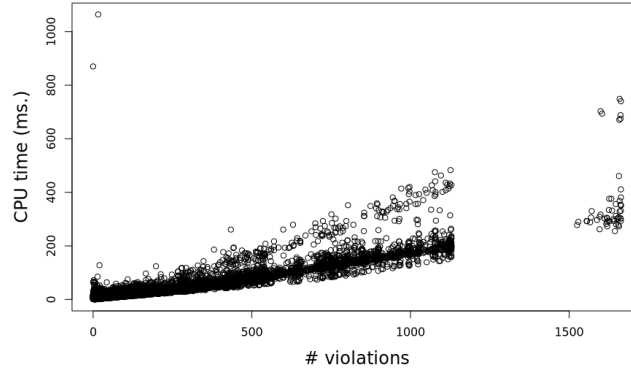


Fig. 4: CPU time spent for the SHACL validation of each discovered shape considering the number of violations

requirement of a broad exploration of the huge search space of possible SHACL shapes to discover acceptable ones. The framework uses a probabilistic SHACL validation process with an acceptability measure and a fitness function to evaluate candidate forms and retain the best ones. Experiments show that our approach captures interesting SHACL shapes describing domain constraints from a real-world RDF dataset. In addition, it provides an effective way to discover a large set of heterogeneous shapes in RDF data, weighting the errors that these may imply and adapting the mining of these shapes. Our future research will focus on studying the scalability of this approach to larger RDF datasets such as DBpedia ($> 50M$ triples). In addition, we plan to extend the approach to the exploration of complex shapes, e.g. shapes with multiple constraints. Finally, it appears important to study the application (and adaptation) of various algorithms proposed by the community, e.g. Structured Grammatical Evolution [15], to the SHACL shapes mining.

Supplemental Material Statement The source code, RDF datasets and obtained results are available in a public repository.⁶

Acknowledgements This work has been partially funded by the 3IA Côte d’Azur “Investments in the Future” project managed by the National Research Agency (ANR) with the reference number ANR-19-P3IA-0002.

References

1. shaclgen 0.2.5.2, <https://github.com/uwlib-cams/shaclgen>

⁶ https://github.com/RemiFELIN/RDFMining/tree/eurogp_2024

2. Ali, M.S., Kshirsagar, M., Naredo, E., Ryan, C.: Autoge: A tool for estimation of grammatical evolution models. In: *International Conference on Agents and Artificial Intelligence (2021)*, <https://api.semanticscholar.org/CorpusID:232106265>
3. Boneva, I., Dusart, J., Fernández Alvarez, D., Gayo, J.E.L.: Shape Designer for ShEx and SHACL Constraints. *ISWC 2019 - 18th International Semantic Web Conference (Oct 2019)*, <https://hal.science/hal-02268667>, poster
4. Cimmino, A., Fernández-Izquierdo, A., García-Castro, R.: Astrea: Automatic generation of SHACL shapes from ontologies. In: *ESWC. Lecture Notes in Computer Science*, vol. 12123, pp. 497–513. Springer (2020)
5. Cérés, R., Corby, O., Demairy, E.: Corese (Mar 2023), <https://github.com/Wimmics/corese>
6. Dick, G., Whigham, P.A.: Initialisation and grammar design in grammar-guided evolutionary computation (2022)
7. Felin, R., Faron, C., Tettamanzi, A.G.B.: A Framework to Include and Exploit Probabilistic Information in SHACL Validation Reports. In: *ESWC (2023)*
8. Fernandez-Álvarez, D., Labra-Gayo, J.E., Gayo-Avello, D.: Automatic extraction of shapes using shexer. *Knowledge-Based Systems* **238**, 107975 (2022). <https://doi.org/10.1016/j.knosys.2021.107975>
9. Harper, R.: Ge, explosive grammars and the lasting legacy of bad initialisation. In: *IEEE Congress on Evolutionary Computation*. pp. 1–8 (2010). <https://doi.org/10.1109/CEC.2010.5586336>
10. Kim, H.T., Ahn, C.W.: A new grammatical evolution based on probabilistic context-free grammar. In: Handa, H., Ishibuchi, H., Ong, Y.S., Tan, K.C. (eds.) *Proceedings of the 18th Asia Pacific Symposium on Intelligent and Evolutionary Systems - Volume 2*. pp. 1–12. Springer International Publishing, Cham (2015)
11. Kontokostas, D., Knublauch, H.: Shapes constraint language (SHACL). W3C recommendation, W3C (2017)
12. Lourenço, N., Assunção, F., Pereira, F., Costa, E., Machado, P.: Structured grammatical evolution: A dynamic approach, pp. 137–161 (01 2018). https://doi.org/10.1007/978-3-319-78717-6_6
13. Lourenço, N., Ferrer, J., Pereira, F.B., Costa, E.: A comparative study of different grammar-based genetic programming approaches. In: McDermott, J., Castelli, M., Sekanina, L., Haasdijk, E., García-Sánchez, P. (eds.) *Genetic Programming*. pp. 311–325. Springer International Publishing, Cham (2017)
14. Lourenço, N., Pereira, F., Costa, E.: Unveiling the properties of structured grammatical evolution. *Genetic Programming and Evolvable Machines* **17** (09 2016). <https://doi.org/10.1007/s10710-015-9262-4>
15. Lourenço, N., Pereira, F.B., Costa, E.: Sge: A structured representation for grammatical evolution. In: Bonnevey, S., Legrand, P., Monmarché, N., Lutton, E., Schoenauer, M. (eds.) *Artificial Evolution*. pp. 136–148. Springer International Publishing, Cham (2016)
16. Medvet, E.: A comparative analysis of dynamic locality and redundancy in grammatical evolution. In: McDermott, J., Castelli, M., Sekanina, L., Haasdijk, E., García-Sánchez, P. (eds.) *Genetic Programming*. pp. 326–342. Springer International Publishing, Cham (2017)
17. Mégane, J., Lourenço, N., Machado, P.: Probabilistic grammatical evolution (2021)
18. Mégane, J., Lourenço, N., Machado, P.: Co-evolutionary probabilistic structured grammatical evolution. In: *Proceedings of the Genetic and Evolutionary Computation Conference. ACM (jul 2022)*. <https://doi.org/10.1145/3512290.3528833>

19. Mégane, J., Lourenço, N., Machado, P.: Probabilistic structured grammatical evolution. In: 2022 IEEE Congress on Evolutionary Computation (CEC). IEEE (jul 2022). <https://doi.org/10.1109/cec55065.2022.9870397>
20. Michel, F., Gandon, F., Ah-Kane, V., Bobasheva, A., Cabrio, E., Corby, O., Gazzotti, R., Giboin, A., Marro, S., Mayer, T., Simon, M., Villata, S., Winckler, M.: Covid-on-the-web: Knowledge graph and services to advance COVID-19 research. In: ISWC (2). Lecture Notes in Computer Science, vol. 12507, pp. 294–310. Springer (2020)
21. Nguyen, T.H., Tettamanzi, A.G.B.: An Evolutionary Approach to Class Disjointness Axiom Discovery. In: Barnaghi, P.M., Gottlob, G., Manolopoulos, Y., Tzouramanis, T., Vakali, A. (eds.) WI 2019 - IEEE/WIC/ACM International Conference on Web Intelligence. pp. 68–75. ACM, Thessaloniki, Greece (Oct 2019). <https://doi.org/10.1145/3350546.3352502>
22. Nguyen, T.H., Tettamanzi, A.G.B.: Grammatical Evolution to Mine OWL Disjointness Axioms Involving Complex Concept Expressions. In: CEC 2020 - IEEE Congress on Evolutionary Computation. pp. 1–8. IEEE, Glasgow, United Kingdom (Jul 2020). <https://doi.org/10.1109/CEC48606.2020.9185681>
23. Nicolau, M., Agapitos, A.: Understanding grammatical evolution: Grammar design, pp. 23–53 (01 2018). https://doi.org/10.1007/978-3-319-78717-6_2
24. Nicolau, M., O’Neill, M., Brabazon, A.: Termination in grammatical evolution: grammar design, wrapping, and tails. pp. 1–8 (06 2012). <https://doi.org/10.1109/CEC.2012.6256563>
25. Omran, P., Taylor, K., Rodríguez Méndez, S., Haller, A.: Learning shacl shapes from knowledge graphs. *Semantic Web* **14**, 1–21 (09 2022). <https://doi.org/10.3233/SW-223063>
26. O’Neill, M., Ryan, C.: Grammatical evolution. *IEEE Trans. Evol. Comput.* **5**(4), 349–358 (2001)
27. Pandit, H., O’Sullivan, D., Lewis, D.: Using ontology design patterns to define shacl shapes. In: WOP@ISWC. pp. 67–71. Monterey California, USA (2018)
28. Rabbani, K., Lissandrini, M., Hose, K.: SHACL and shex in the wild: A community survey on validating shapes generation and adoption. In: WWW (Companion Volume). pp. 260–263. ACM (2022)
29. Rothlauf, F., Oetzel, M.: On the locality of grammatical evolution. In: Collet, P., Tomassini, M., Ebner, M., Gustafson, S., Ekárt, A. (eds.) *Genetic Programming*. pp. 320–330. Springer Berlin Heidelberg, Berlin, Heidelberg (2006)
30. Ryan, C., O’Neill, M., Collins, J.: Introduction to 20 Years of Grammatical Evolution, pp. 1–21. Springer International Publishing, Cham (2018). https://doi.org/10.1007/978-3-319-78717-6_1