

1 ShExML: Improving the usability of 2 heterogeneous data mapping languages for 3 first-time users

4 Herminio García-González¹, Iovka Boneva², Sławek Staworko², José
5 Emilio Labra-Gayo¹, and Juan Manuel Cueva Lovelle¹

6 ¹Department of Computer Science, University of Oviedo, Oviedo, Asturias, Spain

7 ²University of Lille, INRIA, Lille, France

8 Corresponding author:

9 Herminio García-González¹

10 Email address: garciaherminio@uniovi.es

11 ABSTRACT

12 Integration of heterogeneous data sources in a single representation is an active field with many different
13 tools and techniques. In the case of text-based approaches—those that base the definition of the
14 mappings and the integration on a DSL—there is a lack of usability studies. In this work we have
15 conducted a usability experiment ($n = 17$) on three different languages: ShExML (our own language),
16 YARRRML and SPARQL-Generate. Results show that ShExML users tend to perform better than those
17 of YARRRML and SPARQL-Generate. This study sheds light on usability aspects of these languages
18 design and remarks some aspects of improvement.

19 1 INTRODUCTION

20 Data integration is the problem of mapping data from different sources so that they can be used through a
21 single interface (Halevy, 2001). In particular, data exchange is the process of transforming source data to
22 a target data model, so that it can be integrated in existing applications (Fagin et al., 2005). Modern data
23 exchange solutions require from the user to define a *mapping* from the source data model to the target
24 data model, which is then used by the system to perform the actual data transformation. This process is
25 crucial to many applications nowadays as the number of heterogeneous data sources is growing (Reinsel
26 et al., 2018).

27 Although many technologies have appeared through the years, the emergence of the semantic web
28 (Berners-Lee et al., 2001) offered new perspectives for data integration. The semantic web principle
29 recommends to represent entities through a unique Internationalized Resource Identifier (IRI) which
30 allows creation of implicit links between distinct datasets simply by reusing existing IRIs. Moreover, the
31 Resource Description Framework (RDF), which is the advocated data format for the semantic web, is
32 compositional, meaning that one can simply fuse data sources without the use of a specific merger. These
33 characteristics make RDF a privileged format for data integration, thus a target for data exchange and
34 transformation.

35 The most notable example of an RDF based data integration system is Wikidata¹ where multiple
36 contributors—humans or robots—transform data from different sources and integrate it to the Wikidata
37 data store. Another example is the `data.bnf.fr`² project that exposes in RDF format the catalog of
38 the French National Library (BNF) by interlinking it with other datasets around the world.

39 Initially, the only way to perform these data transformations was to use *ad-hoc* scripts designed to take
40 one data source and transform it to an RDF output. This supposed the creation of a dedicated script for
41 every new input data source that needed to be converted. Such solutions are slow and costly to develop.

¹<https://www.wikidata.org/>

²<https://data.bnf.fr/en/about> for more information on the project

42 Later on, Domain Specific Language (DSL) approaches emerged which are able to define a translation
43 in a declarative fashion instead of an imperative one. This technique lowers the development time, but a
44 script for every different data source is still needed, which can be a maintenance issue.

45 More recent systems allow direct transformation of multiple data sources into a single representation.
46 Some of them provide dedicated DSLs in which a single script defines the multi-source transformation,
47 others provide graphical interfaces. This is an improvement compared to previous techniques as in
48 principle it allows for faster development and improved maintainability (Meester et al., 2019). However,
49 the adoption of such systems depends also on their *usability* (Hanenberg, 2010).

50 With usability in mind we have designed the ShExML (García-González et al., 2018) language that
51 allows transformation and integration of data from XML and JSON sources in a single RDF output.
52 ShExML uses Shape Expressions (ShEx) (Prud'hommeaux et al., 2014) for defining the desired structure
53 of the output. ShExML has text based syntax (in contrast to graphical tools) and is intended for users that
54 prefer this kind of representation. Our hypothesis is that for first-time users with some programming and
55 Linked Data background, data integration is performed more easily using ShExML than using one of the
56 existing alternatives. The consequent research questions that we study in the current paper are:

- 57 • RQ1: Is ShExML more usable for first-time users over other languages?
- 58 • RQ2: If true, can a relation be established between features support and usability for first-time
59 users?
- 60 • RQ3: Which parts of ShExML—and of other languages—can be improved to increase usability?

61 In the case of this work we are going to focus on usability of tools based on a DSL and see how the
62 design of the language can have an effect on usability and associated measures such as: development time,
63 learning curve, etc.

64 The rest of the paper is structured as follows: Section 2 studies the related work, in Section 3 the
65 three languages are compared alongside a features comparison between them, in Section 4 we describe
66 the methodology followed in the study, in Section 5 the results are presented along with their statistical
67 analysis. In Section 6 we discuss and interpret the results and in Section 7 we draw some conclusions and
68 propose some future lines from this work.

69 2 BACKGROUND

70 We first review available tools and systems for generating RDF from different systems for data representa-
71 tion. These can be divided into one-to-one and many-to-one transformations. We also survey existing
72 studies on the effectiveness of heterogeneous data mapping tools.

73 2.1 One to one transformations

74 Much research work has been done in this topic where conversions and technologies were proposed to
75 transform from a structured format (e.g., XML, JSON, CSV, Databases, etc.) to RDF.

76 2.1.1 From XML to RDF

77 In XML ecosystem many conversions and tools have been proposed:

78 Miletic et al. (2007) describe their experience with the transformation of RDF to XML (and vice
79 versa) and from XML Schema to RDF Schema. Deursen et al. (2008) propose a transformation from
80 XML to RDF which is based on an ontology and a mapping document. An approach to convert XML
81 to RDF using XML Schema is reported by Battle (2004, 2006). Thuy et al. (2008) describe how they
82 perform a translation from XML to RDF using a matching between XML Schema and RDF Schema. The
83 same procedure was firstly proved with a matching between DTD and RDF Schema by the same authors
84 in (Thuy et al., 2007). Breitling (2009) reports a technique for the transformation between XML and RDF
85 by means of the XSLT technology which is applied to astronomy data. Another approach that uses XSLT
86 attached to schemata definitions is described by Sperberg-McQueen and Miller (2004). However, use of
87 XSLT for lifting purposes tends to end up in complex and non flexible stylesheets. Thus, Bischof et al.
88 (2012) present XSPARQL, a framework that enables the transformation between XML and RDF by using
89 XQuery and SPARQL to overcome the drawbacks of using XSLT for these transformations.

90 **2.1.2 From JSON to RDF**

91 Although in the JSON ecosystem there are less proposed conversions and tools, there are some works that
92 should be mentioned.

93 Müller et al. (2013) present a transformation of a RESTful API serving interlinked JSON documents
94 to RDF for sensor data. An RDF production methodology from JSON data tested on the Greek open data
95 repository is presented by Theocharis and Tsihrintzis (2016). Freire et al. (2017) report a tool able to
96 identify JSON metadata, align them with vocabulary and convert it to RDF; in addition, they identify the
97 most appropriate entity type for the JSON objects.

98 **2.1.3 From tabular form to RDF**

99 The importance of CSV (along with its spreadsheet counterparts) has influenced work in this ecosystem:

100 Ermilov et al. (2013) present a mapping language whose processor is able to convert from tabular
101 data to RDF. A tool for translating spreadsheets to RDF without the assumption of identical vocabulary
102 per row is described by Han et al. (2008). Fiorelli et al. (2015) report a platform to import and lift from
103 spreadsheet to RDF with a human-computer interface. Using SPARQL 1.1 syntax TARQL³ offers an
104 engine to transform from CSV to RDF. CSVW proposed a W3C Recommendation to define CSV to RDF
105 transformations using a dedicated DSL (Tandy et al., 2015).

106 **2.1.4 From Databases to RDF**

107 Along with the XML ecosystem, relational database transformation to RDF is another field:

108 Bizer and Seaborne (2004) present a platform to access relational databases as a virtual RDF store. A
109 mechanism to directly map relational databases to RDF and OWL is described by Sequeda et al. (2012);
110 this direct mapping produces a OWL ontology which is used as the basis for the mapping to RDF. Triplify
111 (Auer et al., 2009) allows to publish relational data as Linked Data converting HTTP-URI requests to
112 relational database queries. One of the most relevant proposals is R2RML (Das et al., 2012) that became a
113 W3C Recommendation in 2012. R2RML offers a standard language to define conversions from relational
114 databases to RDF. In order to offer a more intuitive way to declare mapping from databases to RDF,
115 Stadler et al. (2015) presented SML which bases its mappings into SQL views and SPARQL construct
116 queries.

117 More comprehensive reviews of tools and comparisons of tools for the purpose of lifting from
118 relational databases to RDF are presented by (Michel et al., 2014; Hert et al., 2011; Sahoo et al., 2009).

119 **2.2 Many to one transformations**

120 Many to one transformations is a recent topic which has evolved to overcome the problem that one to one
121 transformations need a different solution for each format and that subsequently must be maintained.

122 **2.2.1 Source-centric approaches**

123 Source-centric approaches are those that, even giving the possibility of transforming multiple data sources
124 to multiple serialisation formats, they base their transformation mechanism in one to one transformations.
125 This can deliver optimal results—if exported to RDF—due to RDF compositional property. Some of
126 the tools available are: OpenRefine⁴ which allows to perform data cleanup and transformation to other
127 formats, DataTank⁵ which offers transformation of data by means of a RESTful architecture, Virtuoso
128 Sponger⁶ is a middleware component of Virtuoso able to transform from a data input format to another
129 serialisation format, RDFizers⁷ employs the Open Semantic Framework to offer hundreds of different
130 format converters to RDF. The Datalift (Scharffe et al., 2012) framework also offers the possibility of
131 transforming raw data to semantic interlinked data sources.

132 **2.2.2 Text-based approaches**

133 The use of a mapping language as the way to define all the mappings for various data sources was
134 first introduced by RML (Dimou et al., 2014) which extends R2RML syntax (Turtle based) to cover
135 heterogeneous data sources. With RML implementations it is possible to gather data from: XML, JSON,
136 CSV, Databases and so on; and put them together in the same RDF output. A similar approach was also

³<http://tarql.github.io/>

⁴<http://openrefine.org/>

⁵<http://thedataatank.com/>

⁶<http://vos.openlinksw.com/owiki/wiki/VOS/VirtSponger>

⁷<http://wiki.opensemanticframework.org/index.php/RDFizers>

137 followed in KR2RML (Slepicka et al., 2015) which proposed an alternative interpretation of R2RML
138 rules paired with a source-agnostic processor facilitating data cleaning and transformation. To deal with
139 non-relational databases, Michel et al. (2015) presented xR2RML language which extends R2RML and
140 RML specifications. Then, SPARQL-Generate (Lefrançois et al., 2016) was proposed which extends
141 SPARQL syntax to serve as a mapping language for heterogeneous data. This solution has the advantage
142 of using a very well-known syntax in the semantic web community and that its implementation is more
143 efficient than RML main one (i.e., RMLMapper⁸) (Lefrançois et al., 2017). To offer a simpler solution for
144 users of text-based approaches, YARRRML (Heyvaert et al., 2018) was introduced which offers a YAML
145 based syntax and its processor⁹ performs a translation to RML rules.

146 **2.2.3 Graphical-based approaches**

147 Graphical tools offer an easier way to interact with the mapping engine and are more accessible to
148 non-expert users. Some of the tools mentioned in the previous source-centric approaches section have a
149 graphical interfaces, like OpenRefine and DataTank. RMLEditor (Heyvaert et al., 2016) offers a graphical
150 interface for the creation of RML rules.

151 **2.3 Related studies**

152 Some studies have been made to evaluate available tools and languages. Lefrançois et al. (2017) compared
153 SPARQL-Generate implementation to RMLMapper. Their results showed that SPARQL-Generate has
154 a better computational performance when transforming more than 1500 CSV rows in comparison with
155 RMLMapper. They also concluded that SPARQL-Generate language is easier to learn and use for semantic
156 web practitioners (who are likely already familiar with SPARQL), but this was based on a limited analysis
157 of the cognitive complexity of query/mappings in the two languages. RMLEditor, a graphical tool to
158 generate RML rules was proposed by Heyvaert et al. (2016). They performed a usability evaluation
159 for their tool with semantic web experts and non-experts. In the case of semantic web experts they
160 also evaluate the differences between the textual approach (RML) and this new visual one. However,
161 RMLEditor was neither compared with other similar tools nor RML with other languages. Heyvaert
162 et al. (2018) proposed YARRRML as a human-readable text-based representation which offers an easier
163 layer on top of RML and R2RML. However, the authors did not present any evaluation of this language.
164 Meester et al. (2019) made a comparative characteristic analysis of different mapping languages. However,
165 a qualitative analysis is not performed and usability is only mentioned in NF1 "Easy to use by Semantic
166 Web experts" which only YARRRML and SPARQL-Generate achieve.

167 Thus, to the best of our knowledge no usability study was performed in these languages which share
168 the easiness of use as one of their goals. Therefore, we introduce this study as a first step into the usability
169 evaluation of heterogeneous data mapping languages.

170 **3 PRESENTATION OF THE LANGUAGES UNDER STUDY**

171 In this section we compare YARRRML, SPARQL-Generate and ShExML syntax by means of a simple
172 example. These three tools each offer a DSL able to define mappings for heterogeneous data sources like
173 we have seen in the previous section and their designers share the goal to be user friendly (Meester et al.,
174 2019; García-González et al., 2018). RML and similar alternatives are not included in the comparison
175 because they have a verbose syntax very close to the RDF data model. While it might be an interesting
176 solution for users without any programming knowledge but familiar with RDF, we consider it more like a
177 lower level middle language to compile to rather than a language to be used by programmers and data
178 engineers. Indeed, YARRRML and ShExML engines are able to compile their mappings to RML.

179 For the sake of the example two small files on JSON and XML are presented in Listing 1 and Listing
180 2 respectively. Each one of these files define two films with 6 attributes—that could differ on name and
181 structure—that will be translated to the RDF output showed in Listing 3. In this example, and with the
182 aim to keep it simple, different ids are used in each entity; however, it is possible to use objects with same
183 ids that could be merged into a single entity or divided into different new entities depending on users'
184 intention.

Listing 1. JSON films file

⁸<https://github.com/RMLio/RML-Mapper>

⁹<https://github.com/RMLio/yarrml-parser>

```

185 {
186   {
187     "films": [
188       {
189         "id": 3,
190         "title": "Inception",
191         "date": "2010",
192         "countryOfOrigin": "USA",
193         "director": "Christopher Nolan",
194         "screenwriter": "Christopher Nolan"
195       },
196       {
197         "id": 4,
198         "title": "The Prestige",
199         "date": "2006",
200         "countryOfOrigin": "USA",
201         "director": "Christopher Nolan",
202         "screenwriter": ["Christopher Nolan",
203                           "Jonathan Nolan"]
204       }
205     ]
206   }
207 }

```

Listing 2. XML films file

```

208 <films>
209   <film id="1">
210     <name>Dunkirk</name>
211     <year>2017</year>
212     <country>USA</country>
213     <director>Christopher Nolan</director>
214     <screenwriters>
215       <screenwriter>Christopher Nolan</screenwriter>
216     </screenwriters>
217   </film>
218   <film id="2">
219     <name>Interstellar</name>
220     <year>2014</year>
221     <country>USA</country>
222     <director>Christopher Nolan</director>
223     <screenwriters>
224       <screenwriter>Christopher Nolan</screenwriter>
225       <screenwriter>Jonathan Nolan</screenwriter>
226     </screenwriters>
227   </film>
228 </films>
229
230

```

Listing 3. RDF output

```

231 @prefix : <http://example.com/> .
232
233
234 :4      :country      "USA" ;
235         :screenwriter  "Jonathan Nolan" ,
236                   "Christopher Nolan" ;
237         :director     "Christopher Nolan" ;
238         :name         "The Prestige" ;
239         :year         :2006 .
240
241 :3      :country      "USA" ;
242         :screenwriter  "Christopher Nolan" ;
243         :director     "Christopher Nolan" ;
244         :name         "Inception" ;
245         :year         :2010 .
246
247 :2      :country      "USA" ;
248         :screenwriter  "Jonathan Nolan" ,
249                   "Christopher Nolan" ;
250         :director     "Christopher Nolan" ;
251         :name         "Interstellar" ;

```

```

252         :year                               :2014 .
253
254 :1      :country                            "USA" ;
255       :screenwriter                        "Christopher Nolan" ;
256       :director                            "Christopher Nolan" ;
257       :name                                "Dunkirk" ;
258       :year                               :2017 .

```

260 3.1 YARRRML

Listing 4. YARRRML transformation script for the films example

```

261
262 prefixes:
263   ex: "http://example.com/"
264
265 mappings:
266   films_json:
267     sources:
268       - ['films.json~jsonpath', '$.films[*]']
269     s: ex:${id}
270     po:
271       - [ex:name, $(title)]
272       - [ex:year, ex:$(date)~iri]
273       - [ex:director, $(director)]
274       - [ex:screenwriter, $(screenwriter)]
275       - [ex:country, $(countryOfOrigin)]
276   films_xml:
277     sources:
278       - ['films.xml~xpath', '//film']
279     s: ex:${@id}
280     po:
281       - [ex:name, $(name)]
282       - [ex:year, ex:$(year)~iri]
283       - [ex:director, $(director)]
284       - [ex:screenwriter, $(screenwriters/screenwriter)]
285       - [ex:country, $(country)]

```

287 YARRRML is designed with human-readability in mind which is achieved through a YAML based
 288 syntax. Listing 4 shows the mappings `films_json` and `films_xml` for our films example. Each
 289 mapping starts with a source definition that contains the query to be used as iterator, e.g., `//film`. It is
 290 followed by the definition of the output given by a subject definition (`s:`) and a number of associated
 291 predicate-object definitions (`po:`). Subject and predicate-object definitions can use “partial” queries
 292 relative to the iterator to populate the subject and object values. This way of defining mappings is very
 293 close to RML; YARRRML actually does not provide an execution engine but is translated to RML.

294 3.2 SPARQL-Generate

Listing 5. SPARQL-Generate transformation script for the films example

```

295
296 BASE <http://example.com/>
297 PREFIX iter: <http://w3id.org/sparql-generate/iter/>
298 PREFIX fun: <http://w3id.org/sparql-generate/fn/>
299 PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
300 PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
301 PREFIX : <http://example.com/>
302 PREFIX dbr: <http://dbpedia.org/resource/>
303 PREFIX schema: <http://schema.org/>
304 PREFIX sc: <http://purl.org/science/owl/sciencecommons/>
305
306 GENERATE {
307   ?id_json :name ?name_json ;
308           :year ?year_json ;
309           :director ?director_json ;
310           :country ?country_json .
311
312 GENERATE {
313   ?id_json :screenwriter ?screenwriter_json .

```

```

314 }
315 ITERATOR iter:Split(?screenwriters_json, ",")
316 AS ?screenwriters_json_iterator
317 WHERE {
318     BIND(REPLACE(?screenwriters_json_iterator,
319         "\\[[\\]]|\\\"", ""))
320     AS ?screenwriter_json)
321 } .
322
323 ?id_xml :name ?name_xml ;
324         :year ?year_xml ;
325         :director ?director_xml ;
326         :country ?country_xml .
327
328 GENERATE {
329     ?id_xml :screenwriter ?screenwriter_xml .
330 }
331 ITERATOR iter:XPath(?film_xml,
332     "/film/screenwriters[*]/screenwriter")
333 AS ?screenwriters_xml_iterator
334 WHERE {
335     BIND(fun:XPath(?screenwriters_xml_iterator,
336         "/screenwriter/text()") AS ?screenwriter_xml)
337 } .
338
339 }
340 ITERATOR iter:JSONPath(
341     <https://raw.githubusercontent.com/hermiolog/ShExML/
342     master/src/test/resources/filmsPaper.json>,
343     "$.films[*]") AS ?film_json
344 ITERATOR iter:XPath(
345     <https://raw.githubusercontent.com/hermiolog/ShExML/
346     master/src/test/resources/filmsPaper.xml>,
347     "//film") AS ?film_xml
348 WHERE {
349     BIND(IRI(CONCAT("http://example.com/",
350         STR(fun:JSONPath(?film_json, "$.id")))) AS ?id_json)
351     BIND(fun:JSONPath(?film_json, "$.title") AS ?name_json)
352     BIND(fun:JSONPath(?film_json, "$.director")
353         AS ?director_json)
354     BIND(IRI(CONCAT("http://example.com/",
355         fun:JSONPath(?film_json, "$.date"))) AS ?year_json)
356     BIND(fun:JSONPath(?film_json, "$.countryOfOrigin")
357         AS ?country_json)
358     BIND(fun:JSONPath(?film_json, "$.director")
359         AS ?directors_json)
360     BIND(fun:JSONPath(?film_json, "$.screenwriter")
361         AS ?screenwriters_json)
362     BIND(IRI(CONCAT("http://example.com/",
363         fun:XPath(?film_xml, "/film/@id"))) AS ?id_xml)
364     BIND(fun:XPath(?film_xml, "/film/name/text()")
365         AS ?name_xml)
366     BIND(fun:XPath(?film_xml, "/film/director/text()")
367         AS ?director_xml)
368     BIND(IRI(CONCAT("http://example.com/",
369         fun:XPath(?film_xml, "/film/year/text()")))
370         AS ?year_xml)
371     BIND(fun:XPath(?film_xml, "/film/country/text()")
372         AS ?country_xml)
373 }

```

375 SPARQL-Generate is an extension of SPARQL 1.1 for querying heterogeneous data sources and creating
376 RDF and text. It offers a set of SPARQL binding functions and SPARQL iterator functions to achieve this
377 goal. The mapping for our films example is shown in Listing 5. The output of the mapping is given within
378 the GENERATE clauses and can use variables and IRIs, while queries, IRI and variable declarations are
379 declared in the WHERE clause. SPARQL-Generate is an expressive language that can be further extended
380 using the SPARQL 1.1 extension system. On the other side, SPARQL-Generate scripts tend to be verbose
381 compared to the other two languages studied in this paper.

Listing 6. ShExML transformation script for the films example

```

383 PREFIX : <http://example.com/>
384 SOURCE films_xml_file <
385     https://raw.githubusercontent.com/herminiogg/
386     ShExML/master/src/test/resources/filmsPaper.xml>
387 SOURCE films_json_file <
388     https://raw.githubusercontent.com/herminiogg/
389     ShExML/master/src/test/resources/filmsPaper.json>
390 ITERATOR film_xml <xpath: //film> {
391     FIELD id <@id>
392     FIELD name <name>
393     FIELD year <year>
394     FIELD country <country>
395     FIELD director <director>
396     FIELD screenwriters <screenwriters/screenwriter>
397 }
398 ITERATOR film_json <jsonpath: $.films[*]> {
399     FIELD id <id>
400     FIELD name <title>
401     FIELD year <date>
402     FIELD country <countryOfOrigin>
403     FIELD director <director>
404     FIELD screenwriters <screenwriter>
405 }
406 EXPRESSION films <films_xml_file.film_xml
407     UNION films_json_file.film_json>
408
409 :Films :[films.id] {
410     :name [films.name] ;
411     :year :[films.year] ;
412     :country [films.country] ;
413     :director [films.director] ;
414     :screenwriter [films.screenwriters] ;
415 }
416

```

418 ShExML, our proposed language, can be used to map XML and JSON documents to RDF. The ShExML
419 mapping for the films example is presented in Listing 6. It consists of source definitions followed by
420 iterator definitions. The latter define structured objects which fields are populated with the results of
421 source queries. The output of the mapping is described using a Shape Expression (ShEx) (Prud’hommeaux
422 et al., 2014; Boneva et al., 2017) which can refer to the previously defined fields. The originality of
423 ShExML, compared to the other two languages studied here, is that the output is defined only once even
424 when several sources are used. This is a design choice that allows the user to separate concerns: how to
425 structure the output on the one hand, and how to extract the data on the other hand.

426 3.4 Comparing languages features

427 In this subsection we compare languages features and what operations are supported or not in each
428 language (see Table 1).

429 Iterators, sources, fields, unions and so on are common to the three languages as they have the same
430 objective. They have different syntaxes, as it can be seen in the three examples, but from a functionality
431 point of view there are no differences.

432 **Source and output definition and their artefacts:** As we saw, the mechanism to define the form
433 of the RDF output has different flavour in the three languages: subject and predicate-object definitions
434 for every source in YARRRML; GENERATE clauses for every source in SPARQL-Generate; a single
435 Shape Expression in ShExML. Additionally, the three languages offer slightly different operators for
436 constructing the output values. All of them typically obtain IRIs by concatenating a source value to
437 some prefix, and reuse literal values as is. YARRRML supports the generation of multiple named graphs
438 whereas SPARQL-Generate can only generate one named graph at a time and ShExML only generates
439 RDF datasets.

440 **Multiple results:** The handling of multiple results, like it occurs on the screenwriters case, is different
441 between SPARQL-Generate and the two other languages. In YARRRML and ShExML if a query returns
442 multiple results they are treated like a list of them. However, in SPARQL-Generate this functionality must

Features	ShExML	YARRRML	SPARQL-Generate
Defining output	Shape expression	Subject and predicate-object definitions	Generate clause
Source and output definition	Prefix and value generation expression (concatenation)	Prefix and value generation expression (array)	Variable (previous use of concat function) or string interpolation
Datatypes & Language tags	Yes	Yes	Yes
Multiple results from a query	Treated like an array	Treated like an array	Need to iterate over the results
Transformations	Limited (Matchers and String operators).	FnO hub	Functions for strings and extension mechanism
Output formats	RDF	RDF	RDF and any text-based format
Translation	RML	RML	No translation offered
Link between mappings	Shape Linking and JOIN keyword (do not fully cover YARRRML feature)	Yes (conditions allowed)	Nested generate clauses, filter clauses and extension mechanism
Conditional mapping generation	No	Yes (Function and conditional clause)	Yes (Filter clause and extension mechanism)

Table 1. Features comparison between the three languages

443 be explicitly declared like it can be seen in Listing 5. It leads to complex iterator definitions like the one
444 used in JSON screenwriters one.

445 **Transformations:** The possibility of transforming the output to another values by means of a
446 function is something very useful for different purposes when building a knowledge graph. Therefore,
447 in YARRRML this is supported through the FnO mechanism (Meester et al., 2017) which offers a way
448 to define functions inside mapping languages in a declarative fashion. SPARQL-Generate offers some
449 functions for strings embedded inside the SPARQL binding functions mechanism; however, it is possible
450 to extend the language through the SPARQL 1.1 extension mechanism. In the case of ShExML, only
451 Matchers and String operations are offered for transformation purposes.

452 **Other formats output:** Output format on YARRRML and ShExML is limited to RDF; whereas, in
453 SPARQL-Generate it is possible to also generate plain text, enabling the potential transformation to a lot
454 of different formats. In this aspect, SPARQL-Generate presents a much more flexible output. Conversely,
455 YARRRML and ShExML engines offer a translation of their mappings to RML rules which improves
456 interoperability with other solutions.

457 **Link to other mappings:** In YARRRML there is the possibility to link mappings between them. This
458 functionality is provided by giving the name of the mapping to be linked and the condition that must
459 be satisfied (e.g., ID of mapping A equal to ID of mapping B). This can be useful when the subject is
460 generated with a certain attribute but this attribute does not appear on the other file so the linking should
461 be done using another attribute. In ShExML this can be partially achieved by Shape linking—which
462 is a syntactic sugar to avoid repeating an expression twice—and by the Join clause which gives an
463 implementation for primary interlinking covering a subset of what is covered with YARRRML mapping
464 linking. In SPARQL-Generate this can be achieved using nested Generate clauses and Filter clauses.

465 **Conditional mapping generation:** Sometimes there is the need to generate triples only in the case
466 that some condition is fulfilled. In YARRRML this is achieved using the conditional clause and a
467 function. In SPARQL-Generate this can be obtained with the SPARQL 1.1 Filter clauses and also with
468 the extensibility mechanism offered by the language. In ShExML this is not possible currently.

469 **Further features of SPARQL-Generate:** Apart from what has been presented in the previous point,
470 SPARQL-Generate, as being based on SPARQL 1.1, offers more expressiveness than the other two
471 languages. One possibility that emerges from that is the use of defined variables. For example, it is
472 possible to define an iterator of numbers and then use that numbers to request different parts of an API.
473 This versatility enables the creation of very complex and rich scripts that can cover a lot of use cases.
474 It is natural to expect that learning to use the full capabilities of SPARQL-Generate is complex, as the
475 language offers a lot of features. In our experiments, however, only some basic features of the language
476 were required and, as is shown in Section 5, it appears that SPARQL-Generate design did not help test
477 subjects to solve the proposed tasks easily.

478 4 METHODOLOGY

479 In order to test our hypothesis that ShExML is easier for first-time users only experienced in programming
480 and the basics of linked data, an experiment was carried out. The University of Oviedo granted ethical
481 approval to carry out the described study. Verbal consent was requested before starting the experiment.

482 **4.1 Experiment design**

483 The selected tools were YARRRML¹⁰, SPARQL-Generate¹¹ and ShExML¹². We decided not to include
484 RML¹³ and similar alternatives for the same reason mentioned on Section 3. Three manuals were designed
485 for the students based on the example about films that described how the integration can be done with
486 each tool¹⁴. The experiment was designed to be performed in each tool dedicated online environment,
487 which are available through the Internet as a webpage.

488 In addition, a small manual was developed to guide the students along the experiment and to inform
489 them about the input files and which are the expected outputs¹⁴. This manual contained two tasks to
490 perform during the experiment which were designed to be performed sequentially, i.e., the student should
491 finish the first task before starting with the second one. The first task was the mapping and integration
492 of two files (JSON and XML) with information about books which should be mapped in a unique RDF
493 graph. The final output should be equal to the one given in the guide. The second task was to modify the
494 script done in the previous task so that the prices are separated and can be compared between markets.
495 In other words, that multiple prices are tagged individually referring to the market where the specific
496 price was found, like they were in the input files. This second task gives us an intuition on how easy is to
497 modify an existing set of data mapping rules in each language.

498 The study was designed as a mixed method approach, including a quantitative analysis and a qualitative
499 analysis. For the quantitative analysis measures, Mousotron¹⁵ was used which allows to register the
500 number of keystrokes, the distance travelled by the mouse and so on. For the qualitative analysis two
501 Office 365 forms were used with questions based on a Likert scale (see questions in Table 2). In addition,
502 the elapsed time was calculated from timestamps in the Office 365 forms.

503 **4.2 Conduction**

504 The sample consisted on 20 students (4 women and 13 men) of the MSc in Web Engineering first-year
505 course (out of two years) at the University of Oviedo¹⁶. Most of them have a bachelor degree (240
506 ECTS credits) in computer science or similar fields. They were receiving a semantic web course of two
507 weeks—a total of 30 hours (3 hours per day)—where they were introduced to semantic technologies
508 like: RDF, SPARQL, ShEx, etc. Before this course they had not previous knowledge on semantic web
509 technologies. Regarding prior knowledge of YAMML by subjects, even though it is normally known and
510 used by developers, we could not assure prior knowledge on it. The experiment was hosted the final day
511 of the course.

512 The experiment was conducted in their usual classroom and with their whole-year-assigned computers.
513 So that they were in a comfortable environment and with a computer they are familiar with. The three
514 tools were assigned to the students in a random manner. Each student received the printed manual for
515 its assigned tool and they were given a time of 20 minutes to read it, test the language in the online
516 environment, and ask doubts and questions. Once these 20 minutes were elapsed the printed experiment
517 guide was given to the students and they were explained about the experiment proceeding with indications
518 about Mousotron operation.

519 In particular the procedure followed to perform the whole experiment was:

- 520 1. Open the assigned tool on the dedicated webpage and clear the given example.
- 521 2. Open Mousotron and reset it.
- 522 3. Proceed with task 1 (start time registered for elapsed time calculation).
- 523 4. Once task 1 is finished, capture Mousotron results (screenshot) and fill the first Office 365 question-
524 naire.
- 525 5. Reset Mousotron and proceed with task 2.

¹⁰<http://rml.io/yarrrml/>

¹¹<https://ci.mines-stetienne.fr/sparql-generate/>

¹²<http://shexml.herminiogarcia.com/>

¹³<http://rml.io/>

¹⁴Material can be consulted on:

<https://github.com/herminiogg/shexml-paper-2019-data/tree/master/experiment-material>

¹⁵<http://www.blacksoftware.com/mousotron.html>

¹⁶<http://miw.uniovi.es/>

526 6. Once task 2 is finished, capture Mousotron results (screenshot) and fill the second Office 365
527 questionnaire.

528 **4.2.1 Analysis**

529 The quantitative results were dump into an Excel sheet and anonymised. Although many results can be
530 used as given by the students, some of them need to be calculated. This is the case of elapsed time (on both
531 tasks), completeness percentage and precision. Elapsed time in the first task (t_{t1}) was calculated as the
532 subtraction of questionnaire 1 beginning time (st_{q1}) and experiment start time (st_e), i.e., ($t_{t1} = st_{q1} - st_e$).
533 Elapsed time in the second task (t_{t2}) was calculated as the subtraction of questionnaire 1 ending time
534 (et_{q1}) and questionnaire 2 beginning time (st_{q2}), i.e., ($t_{t2} = st_{q2} - et_{q1}$).

535 Completeness percentage was calculated from three measures: the proportion of correctly generated
536 triples contributed 50%, the proportion of data correctly translated contributed 25% and the proportion of
537 correctly generated prefixes and datatypes as a 25%. This design gives more importance to the structure,
538 which is the main goal when using these tools. Other aspects, like correct data (i.e., the object part of a
539 triple), prefixes (i.e., using the correct predicate for the subject, the predicate and the object in case of an
540 IRI) and the datatype (i.e., putting the correct xsd type in case of a literal object) are a little less valued
541 as these errors could come more easily from a distraction or an oversight. Let CP be the completeness
542 percentage, t the number of triples, d the number of data gaps and $p\&dt$ the number of prefixes and
543 datatypes, so the calculation of the completeness percentage can be expressed as:

$$CP = 0.5 * \frac{t_{total} - t_{generated}}{t_{total}} + 0.25 * \frac{d_{total} - d_{generated}}{d_{total}} + 0.25 * \frac{p\&dt_{total} - p\&dt_{generated}}{p\&dt_{total}}$$

544 Finally, precision was calculated as the division of current student elapsed time by minimum elapsed
545 time of all students, multiplied by the completeness percentage. This precision formulation gives us
546 an intuition on how fast was some student in comparison with the fastest student and with a correction
547 depending on how well his/her solution was. Let t_{sn} be the elapsed time of student n and CP_{sn} the
548 completeness percentage of student n calculated with the previous formula.

$$Precision_{sn} = \frac{t_{sn}}{\min(\{t_{s1}, \dots, t_{sn}\})} * CP_{sn}$$

549 The results of the qualitative analysis were only anonymised as they can be directly used from the
550 Office 365 output.

551 For the analysis the IBM SPSS version 24 was used. We planned a One Way ANOVA test within the
552 three groups in the quantitative analysis where a normal distribution was found and the Kruskal-Wallis
553 test where not. The qualitative analysis comparison between three groups was established using the
554 Kruskal-Wallis test. The report and analysis of the results was made using Field (2013) as guidance and
555 using the suggested APA style as a standard manner to report statistical results.

556 **4.3 Threat to validity**

557 In this experiment we have identified the following threats to its validity.

558 **4.3.1 Internal validity**

559 We have identified the following internal validity threats in the experiment design:

- 560 • More expertise in some specific tool: In semantic web area—as in other areas—people tend to be
561 more expert in some specific technologies and languages. The derived risk is that this expertise can
562 have an influence on final results. To alleviate this we have selected MSc students that are studying
563 the same introductory semantic web course and we have assigned the tools in a random manner.
- 564 • Not homogeneous group: It is possible that the selected group is not homogeneous on skills and
565 previous knowledge. To mitigate this we have applied the same measures as for the previous threat:
566 Students of a semantic web course and a randomised tool assignment.
- 567 • Unfamiliar environment: In usability studies, unfamiliar environments can play a role on final
568 conclusions. Therefore, we opted to run the experiment in a well-known environment for the
569 students, that is, their whole-year classroom.

Table 2. Statements to evaluate by the students based on a 5 point Likert scale

Questionnaire	Statement	Obtained Variable
1	The experience with the tool was satisfactory	General satisfaction level
1	The tool was easy to use	Easiness of use
1	The mapping definitions was easy	Mapping definition easiness
1	The language was easy to learn	Learnability
1	I find that these tool can be useful in my work	Applicability
1	The coding in this tool was intuitive	Intuitiveness
1	The language design leads to commit some errors	Error proneness
1	The error messages were useful to solve the problems	Error reporting usefulness
2	It was easy to define different predicates for the price	Modifiability

- 570 • More guide and information about one tool: As we have designed one of the languages, it could
571 lead to a bias in information delivery. To try to mitigate this threat we developed three identical
572 manuals for each tool. Questions and doubts were answered equally for all the students and tools.

573 4.3.2 External validity

574 Following the measures taken in the internal validity threats we identified the corresponding external
575 validity ones:

- 576 • Very focused sample: As we have restricted the profile of the sample to students of a MSc course
577 which are more or less within the same knowledge level, there is the risk that these findings cannot
578 be extrapolated for other samples or populations. It is possible that for semantic web practitioners—
579 with different interests and expertises—these findings are not applicable. However, the intention of
580 this study was to evaluate usability with first-time users as a first step to guide future studies.

581 5 RESULTS

582 From the 20 students of the sample¹⁷, in the first task, 3 of them left the experiment without making any
583 questionnaire, 2 for SPARQL-Generate and 1 for YARRRML. In the second task, only 7 out of the 20
584 students made the questionnaire, 6 for ShExML and 1 for YARRRML. The statistical analysis was made
585 using the IBM SPSS software, version 24.

586 **Task 1:** As previously stated, the number of students that finished—correctly or not—the proposed
587 task was 17. Descriptive statistics can be seen in Table 3. Comparison of three groups was made
588 by means of a One Way ANOVA which results showed significant differences on elapsed seconds
589 $F(2, 14) = 6.00, p = .013, \omega = .60$. As completeness percentage and precision are not following a
590 normal distribution on SPARQL-Generate group ($W(4) = .63, p = .001$ and $W(4) = .63, p = .001$), the
591 comparison was established by means of the Kruskal-Wallis test which showed significant differences
592 in both variables ($H(2) = 9.73, p = .008$ and $H(2) = 9.68, p = .008$). Post hoc test for elapsed seconds
593 using the Gabriel’s criterion showed significant differences between ShExML group and YARRRML
594 group ($p = .016$). Post hoc test for completeness percentage and precision using the Bonferroni’s
595 criterion showed significant differences between ShExML and SPARQL-Generate ($p = .012, r = .87$ and
596 $p = .012, r = .87$). Likert scale questionnaire results ($\alpha = 0,73$) (see Fig. 1) were analysed using Kruskal-
597 Wallis test which resulted in significant differences between groups for variables general satisfaction
598 level ($H(2) = 6.28, p = .043$), easiness of use ($H(2) = 9.82, p = .007$), mapping definition easiness
599 ($H(2) = 10.25, p = .006$) and learnability ($H(2) = 8.63, p = .013$). Bonferroni’s criterion was used as post
600 hoc test for the variables with significant differences. For general satisfaction level significant differences
601 were found between ShExML and YARRRML ($p = .039, r = .69$). For easiness of use significant
602 differences were found between ShExML and YARRRML ($p = .011, r = .81$). For mapping definition
603 easiness significant differences were found between ShExML and SPARQL-Generate ($p = .013, r = .90$)
604 and between ShExML and YARRRML ($p = .037, r = .69$). For learnability significant differences
605 were found between ShExML and SPARQL-Generate ($p = .042, r = .78$) and between ShExML and
606 YARRRML ($p = .040, r = .69$).

¹⁷Original datasets available on:
<https://github.com/herminiogg/shexml-paper-2019-data/tree/master/datasets>

Table 3. Descriptive statistics for task 1 objective results where n is the sample size, \bar{x} is the mean, s is the standard deviation, max is the maximum value of the sample and min is the minimum value of the sample. (*) means significant differences between groups and (a) means significant differences in the post hoc test between the marked groups at the level of significance ($\alpha = .05$). Differences in totals are due to malfunctions while operating capture software.

Measure	Group	n	\bar{x}	s	max	min
Elapsed seconds (*)	ShExML (a)	7	1560.1429	541.57376	2192	782
	YARRRML (a)	6	2443.8333	375.44502	2896	1891
	SPARQL-Generate	4	2292.7500	533.49063	2769	1634
	Total	17	2044.4118	620.68370	2896	782
Keystrokes	ShExML	6	1138.50	610.588	2287	674
	YARRRML	4	1187	449.649	1795	810
	SPARQL-Generate	3	1125.67	121.476	1265	1042
	Total	13	1150.46	457.183	2287	674
Left button clicks	ShExML	6	176.50	112.169	327	58
	YARRRML	4	318.75	177.989	551	170
	SPARQL-Generate	3	166	78.791	254	102
	Total	13	217.85	138.267	551	58
Right button clicks	ShExML	6	2.17	2.137	6	0
	YARRRML	4	2.25	1.708	4	0
	SPARQL-Generate	2	4.50	2.121	6	3
	Total	12	2.58	2.021	6	0
Mouse wheel scroll	ShExML	6	148	183.737	486	13
	YARRRML	4	679.25	606.711	1404	101
	SPARQL-Generate	3	199	131.160	348	101
	Total	13	323.23	412.819	1404	13
Meters travelled by the mouse	ShExML	7	30.400	24.318	70.079	0
	YARRRML	6	43.454	43.144	101.767	0
	SPARQL-Generate	4	21.220	16.526	37.680	0
	Total	17	32.847	30.550	101.767	0
Completeness percentage (*)	ShExML (a)	7	0.771	0.296	1	0.19
	YARRRML	6	0.323	0.366	0.82	0
	SPARQL-Generate (a)	4	0.02	0.04	0.08	0
	Total	17	0.436	0.415	1	0
Precision (*)	ShExML (a)	7	0.495	0.286	1	0.07
	YARRRML	6	0.131	0.160	0.38	0
	SPARQL-Generate (a)	4	0.005	0.01	0.02	0
	Total	17	0.251	0.292	1	0

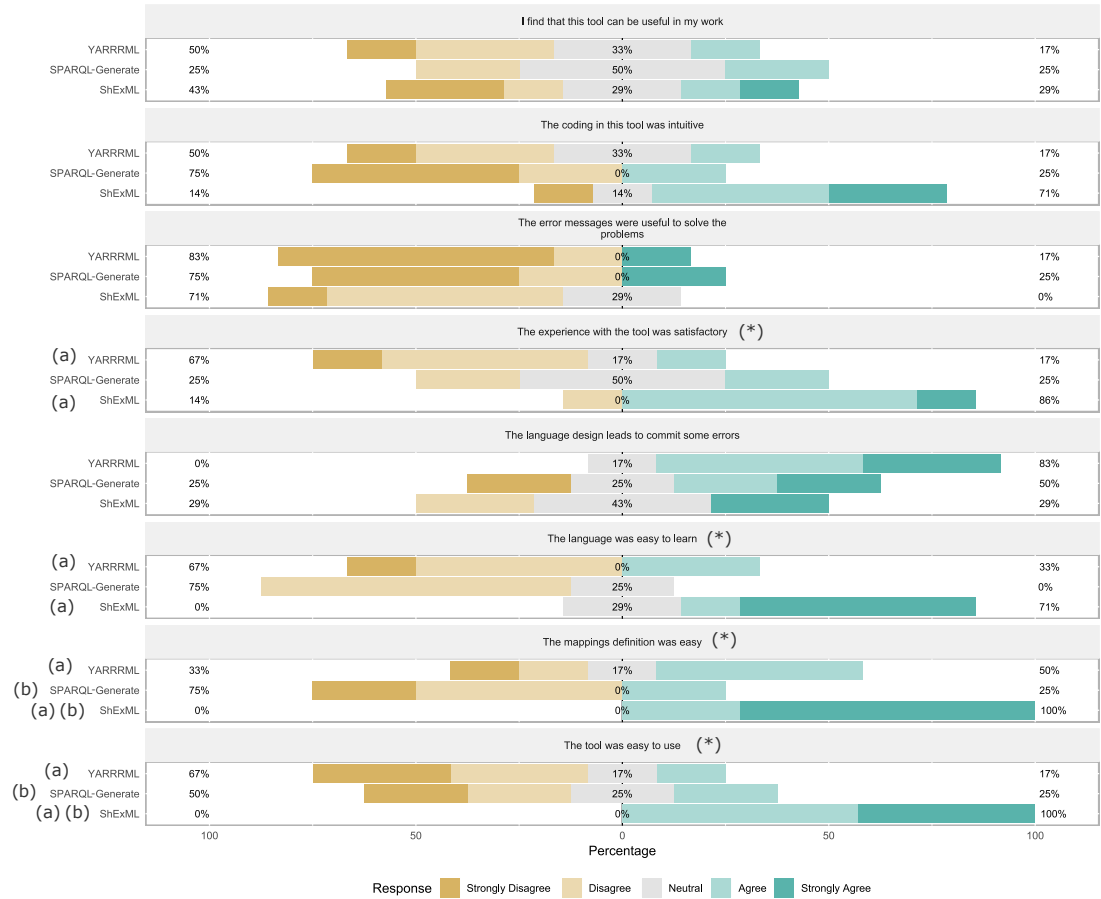


Figure 1. Task 1 results for Likert scale questionnaire where results are divided into questions and groups. (*) means significant differences between groups and (a) and (b) means significant differences in the post hoc test between the marked groups at the level of significance ($\alpha = .05$)

Table 4. Descriptive statistics for task 2 objective results where n is the sample size, \bar{x} is the mean, s is the standard deviation, max is the maximum value of the sample and min is the minimum value of the sample. Differences in totals are due to malfunctions while operating capture software.

Measure	Group	n	\bar{x}	s	max	min
Elapsed seconds	ShExML	6	325.5	328.9248	879	3
	YARRRML	1	47	0	47	47
	Total	7	285.7143	318.1822	879	3
Keystrokes	ShExML	5	206.40	175.832	438	43
	YARRRML	1	91	0	91	91
	Total	6	187.17	164.174	438	43
Left button clicks	ShExML	5	61.80	81.417	207	16
	YARRRML	1	43	0	43	43
	Total	6	58.67	73.225	207	16
Right button clicks	ShExML	5	0.40	0.548	1	0
	YARRRML	1	0	0	0	0
	Total	6	0.33	0.516	1	0
Mouse wheel scroll	ShExML	5	123.80	129.494	288	0
	YARRRML	1	41	0	41	41
	Total	6	110	120.655	288	0
Meters travelled by the mouse	ShExML	6	9.7629	13.8829	37.7565	0
	YARRRML	1	11.7563	0	11.7563	11.7563
	Total	7	10.0477	12.6957	37.7565	0
Completeness percentage	ShExML	6	0.73	0.3904	1	0
	YARRRML	1	0	0	0	0
	Total	7	0.6257	0.4507	1	0
Precision	ShExML	6	0.4683	0.37467	1	0
	YARRRML	1	0	0	0	0
	Total	7	0.4014	0.38512	1	0

607 **Task 2:** In this task only 7 students reached this step: 6 for ShExML and 1 for YARRRML.
608 Descriptive statistics of this task can be seen in Table 4. No significant differences were found in any of
609 the variables. In subjective variable analysis (see Fig. 2) no significant differences were found.

610 6 DISCUSSION

611 6.1 Statistical results discussion

612 Results of task 1 show that variables like keystrokes, left button clicks, right button clicks, mouse wheel
613 scroll and meters travelled by the mouse, do not have a significant variability depending on the used tool.
614 This suggests that web interfaces used as online development environments are more or less homogeneous
615 and do not have an impact on the development of the scripts. However, keystrokes variable results
616 should be considered with caution because for SPARQL-Generate the mean of completeness percentages
617 was very low; therefore, achieving a final solution may involve more keystrokes. On the other hand,
618 elapsed seconds, completeness percentage and precision show significant differences between groups
619 which suggest that the selected language has an influence on these variables. Moreover, we can see that
620 elapsed seconds has a medium size effect ($\omega = .60$). Post hoc results show that there are significant
621 differences between ShExML and YARRRML which suggests that YARRRML users tend to need more
622 time than ShExML users for these tests. In the case of comparisons with SPARQL-Generate there are not
623 significant differences which can be due to the small sample size and the low completeness percentage.
624 Differences between ShExML and SPARQL-Generate for completeness percentage and precision suggest
625 that SPARQL-Generate users were not as achieve working solutions as ShExML users, which have the
626 highest mean on both variables. However, between ShExML and YARRRML groups there were no
627 significant differences which is in line with the great variability of those two variables.

628 Results of task 2 do not show any significant difference between the ShExML group and the
629 YARRRML group. This can be explained by the low sample size in the YARRRML group where

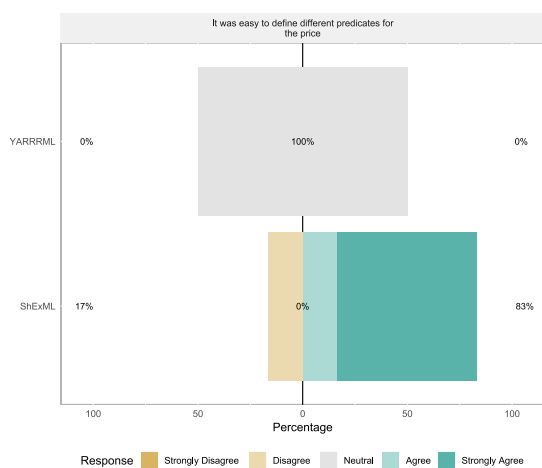


Figure 2. Task 2 results for Likert scale questionnaire where results are divided into the two groups.

630 only one individual made this step. However, completeness percentage and precision show us that some
 631 students did achieve a correct solution with ShExML, whereas in YARRRML group and in SPARQL group
 632 they did not. This leads to the conclusion that only the ShExML group managed to find a working solution
 633 to both proposed tasks. Nevertheless, these conclusions must be validated with bigger experiments to
 634 have statistical confidence.

635 The differences in completeness percentage and precision between ShExML and SPARQL-Generate
 636 and also between ShExML and YARRRML in elapsed seconds can lead us to the conclusion that
 637 usability on first-time users is improved by using ShExML over the other two languages, which answers
 638 RQ1. Moreover, this conclusion is reinforced by the situation that in task 2 neither YARRRML nor
 639 SPARQL-Generate users were able to find a solution to this task.

640 Regarding the subjective analysis, significant differences were found between groups in general
 641 satisfaction level, difficulty of use, easiness of use and learnability (as perceived by the students).

642 On general satisfaction level significant differences were found between ShExML and YARRRML
 643 which indicates that ShExML users were more satisfied with the overall use of the tool respect to the
 644 YARRRML users. Differences between SPARQL-Generate users and the two other groups could not be
 645 established due to their low completeness percentage and precision rates.

646 In the case of easiness of use significant differences were found between ShExML and YARRRML
 647 which suggests that ShExML users found this language easier to use than YARRRML users did with
 648 their language counterpart. In this case, like in the previous variable, significant differences could not be
 649 established between SPARQL-Generate and the two other groups due to low completeness percentage. In
 650 mapping definition easiness differences were established between ShExML group and YARRRML group
 651 and between ShExML group and SPARQL-Generate group which indicates that ShExML users found
 652 mappings easier to define in ShExML than in the other two languages. We also note that users did not
 653 find differences on mapping definition easiness between YARRRML and SPARQL-Generate, this may be
 654 because SPARQL-Generate users did not use the whole language.

655 On learnability significant differences were found between ShExML and SPARQL-Generate and
 656 between ShExML and YARRRML which suggests that the users found easier to learn ShExML than
 657 the other two languages. However, no significant differences were found between YARRRML and
 658 SPARQL-Generate which seems strange due to the difference of verbosity between the two languages.

659 Differences on subjective analysis between ShExML and YARRRML on general satisfaction level,
 660 mapping definition easiness, easiness of use and learnability, and between ShExML and SPARQL-
 661 Generate on mapping definition easiness and learnability comes to corroborate what we have elucidated
 662 with the objective analysis answering RQ1.

663 Review of the other variables shows that the users do not see much applicability on the three languages,
 664 that the design of the languages leads users to commit some errors during the development of the script
 665 and that the error reporting system in the three of them is not very useful to solve the incoming problems.

666 The feedback received from the users in the error proneness and error reporting usefulness variables
667 determines that these two aspects are the ones that should be improved in the three languages to improve
668 their usability. This comes to answer the RQ3.

669 For the modifiability variable assessed in task 2, ShExML users tend to rate this feature with high
670 marks whereas the single YARRRML user gave a response of 3 in a 5 point Likert scale which is in line
671 with his/her completeness percentage mark. As with the objective results of task 2, these subjective results
672 should be further validated in future bigger experiments to corroborate these early findings.

673 **6.2 Alignment with features comparison**

674 In the light of the statistical analysis outcome, SPARQL-Generate design has been shown to have a
675 negative impact on first-time users. This led to three users abandoning the task and low completeness
676 scores for the rest of the group. Although having more features in a language is something good and
677 desirable, these results caught attention on how these features should be carefully designed and included
678 in the language in order to improve ease of use, and thus overall adoption of the tool. In the case of
679 YARRRML language, although it has been designed with human-friendliness in mind, in our experiment
680 it has not reached the expected results in comparison with ShExML. However, it has better results than
681 SPARQL-Generate, suggesting it is less complex to use than that language, but still more complex than
682 ShExML. Nevertheless, it does not seem that supported features could explain the differences between
683 YARRRML and ShExML as the features used on the experiment are more or less equal. Instead other
684 syntax details may be affecting the differences between these two groups such as: the use of keywords
685 that made the language more self explanatory and the modularity used on iterators which reminds of
686 object-oriented programming languages. However, this would require a broader study taking into account
687 programming style background of participants and their own style preferences using techniques like a
688 cognitive complexity architecture (Hansen et al., 2012) to identify how each feature and its design is
689 affecting the usability of each specific language.

690 These results highlight the importance on how features are designed and included in a language.
691 Therefore, SPARQL-Generate with more features and being a highly flexible language tends to have a
692 bad influence on users' usability. Comparing ShExML and YARRRML we see that these differences are
693 smaller than with SPARQL-Generate and that features support does not seem to be the variable affecting
694 YARRRML usability. Thus, we can conclude—and answer the RQ2—that it is not the features supported
695 by a language which affects usability of first-time users but their design.

696 **7 CONCLUSIONS AND FUTURE WORK**

697 In this work we have compared the usability of three heterogeneous data mapping languages. The findings
698 of our user study were that better results, and speed on finding this solution, are related to ShExML users
699 whereas SPARQL-Generate users were not able to find any solution under study conditions. In the case of
700 YARRRML users, they performed better than SPARQL-Generate users but worse than ShExML users
701 finding partial solutions to the given problem.

702 This study is (to our knowledge) the first to explore the topic of usability for first-time users with
703 programming and Linked Data background in these kind of languages. It also reflects the importance that
704 usability has on the accuracy of the encountered solutions and how features should be carefully designed
705 in a language to not impact negatively on its usability.

706 As future work, bigger experiments should be carried out with an emphasis on programming style
707 background and styles (using cognitive complexity frameworks) to corroborate and expand these early
708 findings. In addition, improving these aspects that were worst rated in the three languages (i.e., error
709 proneness and the error reporting system) would enhance perceived user friendliness.

710 This work highlights the importance of usability on these kind of languages and how it could affect
711 their adoption.

712 **ACKNOWLEDGEMENTS**

713 We want to thank the students of the Master's Degree in Web Engineering for their willingness to
714 participate in the experiment described in this work.

715 FUNDING

716 This work has been funded by the Principality of Asturias through the Severo Ochoa call (grant BP17-
717 29), by the Ministry of Economy, Industry and Competitiveness under the call of "Programa Estatal
718 de I+D+i Orientada a los Retos de la Sociedad" (project TIN2017-88877-R), the CPER Nord-Pas de
719 Calais/FEDER DATA Advanced data science and technologies 2015-2020, and the ANR project DataCert
720 ANR-15-CE39-0009. There was no additional external funding received for this study.

721 REFERENCES

- 722 Auer, S., Dietzold, S., Lehmann, J., Hellmann, S., and Aumüller, D. (2009). Triplify: light-weight linked
723 data publication from relational databases. In *Proceedings of the 18th International Conference on*
724 *World Wide Web, WWW 2009, Madrid, Spain, April 20-24, 2009*, pages 621–630.
- 725 Battle, S. (2004). Round-tripping between XML and RDF. In *International Semantic Web Conference*
726 *(ISWC), Hiroshima, Japan*.
- 727 Battle, S. (2006). Gloze: XML to RDF and back again. In *Proceedings of the First Jena User Conference*.
- 728 Berners-Lee, T., Hendler, J., and Lassila, O. (2001). The semantic web. *Scientific american*, 284(5):28–37.
- 729 Bischof, S., Decker, S., Krennwallner, T., Lopes, N., and Polleres, A. (2012). Mapping between RDF and
730 XML with XSPARQL. *Journal on Data Semantics*, 1(3):147–185.
- 731 Bizer, C. and Seaborne, A. (2004). D2rq-treating non-rdf databases as virtual rdf graphs. In *Proceedings of*
732 *the 3rd international semantic web conference (ISWC2004)*, volume 2004. Proceedings of ISWC2004.
- 733 Boneva, I., Labra Gayo, J. E., and Prud'hommeaux, E. G. (2017). Semantics and validation of shapes
734 schemas for rdf. In d'Amato, C., Fernandez, M., Tamma, V., Lecue, F., Cudré-Mauroux, P., Sequeda,
735 J., Lange, C., and Heflin, J., editors, *The Semantic Web – ISWC 2017*, pages 104–120, Cham. Springer
736 International Publishing.
- 737 Breitling, F. (2009). A standard transformation from XML to RDF via XSLT. *Astronomische Nachrichten*,
738 330(7):755–760.
- 739 Das, S., Sundara, S., and Cyganiak, R. (2012). R2RML: RDB to RDF Mapping Language. <https://www.w3.org/TR/r2rml/>.
- 741 Deursen, D. V., Poppe, C., Martens, G., Mannens, E., and de Walle, R. V. (2008). XML to RDF
742 Conversion: A Generic Approach. In *Automated solutions for Cross Media Content and Multi-channel*
743 *Distribution, 2008. AXMEDIS '08. International Conference on*, pages 138–144, Washington.
- 744 Dimou, A., Sande, M. V., Colpaert, P., Verborgh, R., Mannens, E., and de Walle, R. V. (2014). RML: A
745 generic language for integrated RDF mappings of heterogeneous data. In *Proceedings of the Workshop*
746 *on Linked Data on the Web co-located with the 23rd International World Wide Web Conference (WWW*
747 *2014), Seoul, Korea, April 8, 2014*.
- 748 Ermilov, I., Auer, S., and Stadler, C. (2013). CSV2RDF: User-driven CSV to RDF mass conversion
749 framework. In *Proceedings of the ISEM*, volume 13, pages 04–06, Graz, Austria.
- 750 Fagin, R., Kolaitis, P. G., Miller, R. J., and Popa, L. (2005). Data exchange: semantics and query
751 answering. *Theor. Comput. Sci.*, 336(1):89–124.
- 752 Field, A. (2013). *Discovering statistics using IBM SPSS statistics*. Sage.
- 753 Fiorelli, M., Lorenzetti, T., Paziienza, M. T., Stellato, A., and Turbati, A. (2015). Sheet2rdf: a flexible and
754 dynamic spreadsheet import&lifting framework for RDF. In *Current Approaches in Applied Artificial*
755 *Intelligence - 28th International Conference on Industrial, Engineering and Other Applications of*
756 *Applied Intelligent Systems, IEA/AIE 2015, Seoul, South Korea, June 10-12, 2015, Proceedings*, pages
757 131–140.
- 758 Freire, F., Freire, C., and Souza, D. (2017). Enhancing JSON to RDF data conversion with entity type
759 recognition. In *Proceedings of the 13th International Conference on Web Information Systems and*
760 *Technologies, WEBIST 2017, Porto, Portugal, April 25-27, 2017*, pages 97–106.
- 761 García-González, H., Fernández-Álvarez, D., and Gayo, J. E. L. (2018). ShExML: An Heterogeneous
762 Data Mapping Language based on ShEx. In *Proceedings of the EKAW 2018 Posters and Demonstrations*
763 *Session co-located with 21st International Conference on Knowledge Engineering and Knowledge*
764 *Management (EKAW 2018), Nancy, France, November 12-16, 2018.*, pages 9–12.
- 765 Halevy, A. Y. (2001). Answering queries using views: A survey. *The VLDB Journal*, 10(4):270–294.
- 766 Han, L., Finin, T., Parr, C. S., Sachs, J., and Joshi, A. (2008). RDF123: from spreadsheets to RDF. In

- 767 *The Semantic Web - ISWC 2008, 7th International Semantic Web Conference, ISWC 2008, Karlsruhe,*
768 *Germany, October 26-30, 2008. Proceedings*, pages 451–466.
- 769 Hanenberg, S. (2010). Faith, hope, and love: an essay on software science’s neglect of human factors. In
770 Cook, W. R., Clarke, S., and Rinard, M. C., editors, *Proceedings of the 25th Annual ACM SIGPLAN*
771 *Conference on Object-Oriented Programming, Systems, Languages, and Applications, OOPSLA 2010,*
772 *October 17-21, 2010, Reno/Tahoe, Nevada, USA*, pages 933–946. ACM.
- 773 Hansen, M. E., Lumsdaine, A., and Goldstone, R. L. (2012). Cognitive architectures: a way forward for
774 the psychology of programming. In Leavens, G. T. and Edwards, J., editors, *ACM Symposium on New*
775 *Ideas in Programming and Reflections on Software, Onward! 2012, part of SPLASH '12, Tucson, AZ,*
776 *USA, October 21-26, 2012*, pages 27–38. ACM.
- 777 Hert, M., Reif, G., and Gall, H. C. (2011). A comparison of rdb-to-rdf mapping languages. In *Proceedings*
778 *of the 7th International Conference on Semantic Systems*, pages 25–32. ACM.
- 779 Heyvaert, P., Dimou, A., Herregodts, A., Verborgh, R., Schuurman, D., Mannens, E., and de Walle, R. V.
780 (2016). Rmleditor: A graph-based mapping editor for linked data mappings. In *The Semantic Web.*
781 *Latest Advances and New Domains - 13th International Conference, ESWC 2016, Heraklion, Crete,*
782 *Greece, May 29 - June 2, 2016, Proceedings*, pages 709–723.
- 783 Heyvaert, P., Meester, B. D., Dimou, A., and Verborgh, R. (2018). Declarative rules for linked data
784 generation at your fingertips! In *The Semantic Web: ESWC 2018 Satellite Events - ESWC 2018 Satellite*
785 *Events, Heraklion, Crete, Greece, June 3-7, 2018, Revised Selected Papers*, pages 213–217.
- 786 Lefrançois, M., Zimmermann, A., and Bakerally, N. (2016). Flexible RDF generation from RDF
787 and heterogeneous data sources with sparql-generate. In *Knowledge Engineering and Knowledge*
788 *Management - EKAW 2016 Satellite Events, EKM and Drift-an-LOD, Bologna, Italy, November 19-23,*
789 *2016, Revised Selected Papers*, pages 131–135.
- 790 Lefrançois, M., Zimmermann, A., and Bakerally, N. (2017). A SPARQL extension for generating RDF
791 from heterogeneous formats. In *The Semantic Web - 14th International Conference, ESWC 2017,*
792 *Portorož, Slovenia, May 28 - June 1, 2017, Proceedings, Part I*, pages 35–50.
- 793 Meester, B. D., Heyvaert, P., Verborgh, R., and Dimou, A. (2019). Mapping languages: Analysis of
794 comparative characteristics. In Chaves-Fraga, D., Heyvaert, P., Priyatna, F., Sequeda, J. F., Dimou,
795 A., Jabeen, H., Graux, D., Sejdiu, G., Saleem, M., and Lehmann, J., editors, *Joint Proceedings of the*
796 *1st International Workshop on Knowledge Graph Building and 1st International Workshop on Large*
797 *Scale RDF Analytics co-located with 16th Extended Semantic Web Conference (ESWC 2019), Portorož,*
798 *Slovenia, June 3, 2019*, volume 2489 of *CEUR Workshop Proceedings*, pages 37–45. CEUR-WS.org.
- 799 Meester, B. D., Maroy, W., Dimou, A., Verborgh, R., and Mannens, E. (2017). RML and fno: Shaping
800 dbpedia declaratively. In *The Semantic Web: ESWC 2017 Satellite Events - ESWC 2017 Satellite Events,*
801 *Portorož, Slovenia, May 28 - June 1, 2017, Revised Selected Papers*, pages 172–177.
- 802 Michel, F., Djimenou, L., Faron-Zucker, C., and Montagnat, J. (2015). Translation of relational and
803 non-relational databases into RDF with xr2rml. In Monfort, V., Krempels, K., Majchrzak, T. A., and
804 Turk, Z., editors, *WEBIST 2015 - Proceedings of the 11th International Conference on Web Information*
805 *Systems and Technologies, Lisbon, Portugal, 20-22 May, 2015*, pages 443–454. SciTePress.
- 806 Michel, F., Montagnat, J., and Zucker, C. F. (2014). *A survey of RDB to RDF translation approaches and*
807 *tools*. PhD thesis, I3S.
- 808 Miletic, I., Vujasinovic, M., Ivezic, N., and Marjanovic, Z. (2007). Enabling Semantic Mediation for
809 Business Applications: XML-RDF, RDF-XML and XSD-RDFS transformations. In Gonçalves, R. J.,
810 Müller, J. P., Mertins, K., and Zelm, M., editors, *Enterprise Interoperability II: New Challenges and*
811 *Approaches*, pages 483–494. Springer London, London.
- 812 Müller, H., Cabral, L., Morshed, A., and Shu, Y. (2013). From restful to SPARQL: A case study on
813 generating semantic sensor data. In *Proceedings of the 6th International Workshop on Semantic Sensor*
814 *Networks co-located with the 12th International Semantic Web Conference (ISWC 2013), Sydney,*
815 *Australia, October 22nd, 2013.*, pages 51–66.
- 816 Prud’hommeaux, E., Labra Gayo, J. E., and Solbrig, H. (2014). Shape Expressions: An RDF Validation
817 and Transformation Language. In *Proceedings of the 10th International Conference on Semantic*
818 *Systems, SEM ’14*, pages 32–40, New York, NY, USA. ACM.
- 819 Reinsel, D., Gantz, J., and Rydning, J. (2018). The Digitization of the World. From Edge to Core.
820 Technical report, Seagate, IDC. Last accessed: October 28, 2019.
- 821 Sahoo, S. S., Halb, W., Hellmann, S., Idehen, K., Thibodeau Jr, T., Auer, S., Sequeda, J., and Ezzat, A.

- 822 (2009). A survey of current approaches for mapping of relational databases to rdf. *W3C RDB2RDF*
823 *Incubator Group Report*, 1:113–130.
- 824 Scharffe, F., Bihanic, L., Képéklian, G., Ateazing, G., Troncy, R., Cotton, F., Gandon, F., Villata, S.,
825 Euzenat, J., Fan, Z., Bucher, B., Hamdi, F., Vandenbussche, P., and Vatan, B. (2012). Enabling linked
826 data publication with the datalift platform. In *Semantic Cities, Papers from the 2012 AAAI Workshop,*
827 *Toronto, Ontario, Canada, July 22-23, 2012.*
- 828 Sequeda, J. F., Arenas, M., and Miranker, D. P. (2012). On directly mapping relational databases to RDF
829 and OWL. In *Proceedings of the 21st World Wide Web Conference 2012, WWW 2012, Lyon, France,*
830 *April 16-20, 2012,* pages 649–658.
- 831 Slepicka, J., Yin, C., Szekely, P. A., and Knoblock, C. A. (2015). KR2RML: an alternative interpretation
832 of R2RML for heterogenous sources. In Hartig, O., Sequeda, J. F., and Hogan, A., editors, *Proceedings*
833 *of the 6th International Workshop on Consuming Linked Data co-located with 14th International*
834 *Semantic Web Conference (ISWC 2105), Bethlehem, Pennsylvania, USA, October 12th, 2015,* volume
835 1426 of *CEUR Workshop Proceedings.* CEUR-WS.org.
- 836 Sperberg-McQueen, C. M. and Miller, E. (2004). On mapping from colloquial XML to RDF using XSLT.
837 In *Extreme Markup Languages®.*
- 838 Stadler, C., Unbehauen, J., Westphal, P., Sherif, M. A., and Lehmann, J. (2015). Simplified RDB2RDF
839 mapping. In Bizer, C., Auer, S., Berners-Lee, T., and Heath, T., editors, *Proceedings of the Workshop*
840 *on Linked Data on the Web, LDOW 2015, co-located with the 24th International World Wide Web Con-*
841 *ference (WWW 2015), Florence, Italy, May 19th, 2015,* volume 1409 of *CEUR Workshop Proceedings.*
842 CEUR-WS.org.
- 843 Tandy, J., Herman, I., and Kellogg, G. (2015). Generating rdf from tabular data on the web, w3c
844 recommendation 17 december 2015. w3c recommendation. [https://www.w3.org/TR/2015/REC-](https://www.w3.org/TR/2015/REC-csv2rdf-20151217)
845 [csv2rdf-20151217.](https://www.w3.org/TR/2015/REC-csv2rdf-20151217)
- 846 Theocharis, S. and Tsihrintzis, G. A. (2016). Rdf serialization from JSON data: The case of JSON data in
847 diavgeia.gov.gr. In *7th International Conference on Information, Intelligence, Systems & Applications,*
848 *IISA 2016, Chalkidiki, Greece, July 13-15, 2016,* pages 1–6.
- 849 Thuy, P. T. T., Lee, Y.-K., Lee, S., and Jeong, B.-S. (2007). Transforming valid XML documents into RDF
850 via RDF schema. In *Next Generation Web Services Practices, 2007. NWeSP 2007. Third International*
851 *Conference on,* pages 35–40. IEEE.
- 852 Thuy, P. T. T., Lee, Y.-K., Lee, S., and Jeong, B.-S. (2008). Exploiting XML schema for interpreting
853 XML documents as RDF. In *Services Computing, 2008. SCC'08. IEEE International Conference on,*
854 volume 2, pages 555–558. IEEE.