

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

Herminio García-González^{Corresp., 1}, **Iovka Boneva**², **Slawek Staworko**², **José Emilio Labra-Gayo**¹, **Juan Manuel Cueva Lovelle**¹

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

² University of Lille, INRIA, Lille, Nord-Pas-de-Calais, France

Corresponding Author: Herminio García-González
Email address: garciaherminio@uniovi.es

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

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

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

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

²University of Lille, INRIA, Lille, France

Corresponding author:

Herminio García-González¹

Email address: garciaherminio@uniovi.es

ABSTRACT

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

1 INTRODUCTION

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

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

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

The first approach to perform such data transformation was to use *ad-hoc* solutions which were designed to take one data source and transform it to an RDF output. This supposed the creation of a dedicated script for 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

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

More recent proposals allow to directly transform multiple data sources into a single representation. Some of them are based on dedicated DSL in which a single script defines the multi-source transformation, others propose graphical interfaces. This is an improvement compared to previous techniques as in principle it allows for faster development and improved maintainability (Meester et al., 2019). However, the adoption of such languages depends also on their *usability*.

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

- RQ1: Is ShExML ~~language design improving usability~~ for first-time users over other languages?
- RQ2: If true, can a relation be established between expressiveness and usability for first-time users?
- RQ3: Which parts of ShExML—and of other languages—can be improved to increase usability?

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

The rest of the paper is structured as follows: Section 2 studies the related work, in Section 3 the three languages are compared alongside a expressiveness comparison between them, Section 4 describes ShExML further features, in Section 5 we describe the methodology followed in the study, in Section 6 the results are exposed with an statistical analysis of them, in Section 7 we discuss and interpret the results and in Section 8 we draw some conclusions and future lines of this work.

2 BACKGROUND

The background of this work can be divided into two categories: one to one transformations and many to one transformations. In addition, we examine some studies made in the field of heterogeneous data mapping tools.

2.1 One to one transformations

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

2.1.1 From XML to RDF

In XML ecosystem many conversions and tools were proposed from which the following are some of them:

Miletic et al. (2007) describe their experience with the transformation of RDF to XML (and vice versa) and from XML Schema to RDF Schema. Deursen et al. (2008) propose a transformation from XML to RDF which is based on an ontology and a mapping document. An approach to convert XML to RDF using XML Schema is reported by Battle (2004, 2006). Thuy et al. (2008) describe how they perform a translation from XML to RDF using a matching between XML Schema and RDF Schema. The same procedure was firstly proved with a matching between DTD and RDF Schema by the same authors in (Thuy et al., 2007). Breitling (2009) reports a technique for the transformation between XML and RDF by means of the XSLT technology which is applied to the astronomy use case. Another approach that uses XSLT attached to schemata definitions is described by Sperberg-McQueen and Miller (2004). Bischof et al. (2012) present XSPARQL, a framework that enables the transformation between XML and RDF by using XQuery and SPARQL to solve the drawbacks of using XSLT for these transformations.

2.1.2 From JSON to RDF

Although in JSON ecosystem there are less proposed conversions and tools, there are ~~some experiences~~ that should be mentioned.

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

2.1.3 From tabular form to RDF

The importance of CSV (along with its spreadsheet counterpart) has influenced ~~some experiences~~ in this ecosystem:

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

2.1.4 From Databases to RDF

Along with the XML ecosystem, relational database transformation to RDF is one of the ~~largest worked~~ fields:

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

~~More~~ tools and comparisons of tools for the purpose of lifting from relational databases to RDF ~~can be consulted on~~ (Michel et al., 2014; Hert et al., 2011; Sahoo et al., 2009).

2.2 Many to one transformations

Many to one transformations is a recent topic which has evolved ~~from~~ the problem that one to one transformations need a different solution for each format and ~~the consequent maintainability problem~~.

2.2.1 Source-centric approaches

Source-centric approaches are those that, even giving the possibility of transforming multiple data sources to multiple serialisation formats, they base their transformation mechanism in one to one transformations. This can deliver optimal results—if exported to RDF—due to RDF compositional property. Some of the tools available are: OpenRefine⁴ which allows to perform data cleanup and transformation to other formats, DataTank⁵ which offers transformation of data by means of a RESTful architecture, Virtuoso Sparger⁶ ~~which~~ is a middleware component of Virtuoso able to transform from a data input format to another serialisation format, RDFizers⁷ ~~which by means of the~~ Open Semantic Framework offer hundreds of different format converters to RDF. Datalift (Scharffe et al., 2012) framework ~~which~~ offers ~~also~~ the possibility of transforming raw data to semantic interlinked data sources.

2.2.2 Text-based approaches

The use of a mapping language as the way to define all the mappings for various data sources was first introduced by RML (Dimou et al., 2014) which extends R2RML syntax (Turtle based) to cover

³<http://tarql.github.io/>

⁴<http://openrefine.org/>

⁵<http://thedataatank.com/>

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

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

heterogeneous data sources. With RML it is possible to gather data from: XML, JSON, CSV, Databases and so on; and put them together in the same RDF output. A similar approach was also followed in KR2RML (Slepicka et al., 2015) which proposed a source-agnostic processor facilitating data cleaning and transformation. To deal with non-relation databases, Michel et al. (2015) presented xR2RML which extends R2RML and RML specifications to achieve this goal. Then, SPARQL-Generate (Lefrançois et al., 2016) was proposed which extends SPARQL syntax to serve as a mapping language for heterogeneous data. This solution has the advantage of using a very well-known syntax in the semantic web community and to be more efficient than RML (Lefrançois et al., 2017). To offer a simpler solution for users of text-based approaches, YARRRML (Heyvaert et al., 2018) was introduced which offers a YAML based syntax and performs a translation to RML rules.

2.2.3 Graphical-based approaches

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

2.3 Related studies

Some studies for this kind of tools and languages were proposed. Lefrançois et al. (2017) tested the performance of SPARQL-Generate in comparison with RML. The results showed that SPARQL-Generate has a better performance when transforming more than 1500 CSV rows in comparison with RML. They also claimed that SPARQL-Generate is easier to learn and use but they only performed a limited cognitive complexity query/mapping study. RMLEditor, a graphical tool to generate RML rules was proposed by Heyvaert et al. (2016). They performed a usability evaluation over this tool with semantic web experts and non-experts. In the case of semantic web experts they also evaluate the differences between the textual approach (RML) and this new visual one. However, RMLEditor was neither compared with other similar tools nor with other languages. Heyvaert et al. (2018) proposed YARRRML as a human-readable text-based representation which offers an easier layer on top of RML and R2RML. However, the authors did not run any evaluation on this language. Meester et al. (2019) made a comparative characteristic analysis of different mapping languages. However, a qualitative analysis is not performed and usability is only mentioned in NF1 "Easy to use by Semantic Web experts" which only YARRRML and SPARQL-Generate achieve. The other mapping languages offer an RDF-based syntax.

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

3 PRESENTATION OF THE LANGUAGES UNDER STUDY

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

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

Listing 1. JSON films file

```
{
```

```

189     "films": [
190         {
191             "id": 3,
192             "title": "Inception",
193             "date": "2010",
194             "countryOfOrigin": "USA",
195             "director": "Christopher Nolan",
196             "screenwriter": "Christopher Nolan"
197         },
198         {
199             "id": 4,
200             "title": "The Prestige",
201             "date": "2006",
202             "countryOfOrigin": "USA",
203             "director": "Christopher Nolan",
204             "screenwriter": ["Christopher Nolan",
205                             "Jonathan Nolan"]
206         }
207     ]
208 }

```

Listing 2. XML films file

```

210 <films>
211   <film id="1">
212     <name>Dunkirk</name>
213     <year>2017</year>
214     <country>USA</country>
215     <director>Christopher Nolan</director>
216     <screenwriters>
217       <screenwriter>Christopher Nolan</screenwriter>
218     </screenwriters>
219   </film>
220   <film id="2">
221     <name>Interstellar</name>
222     <year>2014</year>
223     <country>USA</country>
224     <director>Christopher Nolan</director>
225     <screenwriters>
226       <screenwriter>Christopher Nolan</screenwriter>
227       <screenwriter>Jonathan Nolan</screenwriter>
228     </screenwriters>
229   </film>
230 </films>
231
232

```

Listing 3. RDF output

```

233 @prefix :      <http://example.com/> .
234
235
236 :4      :country      "USA" ;
237         :screenwriter  "Jonathan Nolan" ,
238                       "Christopher Nolan" ;
239         :director      "Christopher Nolan" ;
240         :name          "The Prestige" ;
241         :year          :2006 .
242
243 :3      :country      "USA" ;
244         :screenwriter  "Christopher Nolan" ;
245         :director      "Christopher Nolan" ;
246         :name          "Inception" ;
247         :year          :2010 .
248
249 :2      :country      "USA" ;
250         :screenwriter  "Jonathan Nolan" ,
251                       "Christopher Nolan" ;
252         :director      "Christopher Nolan" ;
253         :name          "Interstellar" ;
254         :year          :2014 .
255

```

```

256 :1      :country      "USA" ;
257      :screenwriter  "Christopher Nolan" ;
258      :director      "Christopher Nolan" ;
259      :name           "Dunkirk" ;
260      :year           :2017 .

```

3.1 YARRRML

Listing 4. YARRRML transformation script for the films example

```

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

```

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

3.2 SPARQL Generate

Listing 5. SPARQL Generate transformation script for the films example

```

297 BASE <http://example.com/>
298
299 PREFIX iter: <http://w3id.org/sparql-generate/iter/>
300 PREFIX fun: <http://w3id.org/sparql-generate/fn/>
301 PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
302 PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
303 PREFIX : <http://example.com/>
304 PREFIX dbr: <http://dbpedia.org/resource/>
305 PREFIX schema: <http://schema.org/>
306 PREFIX sc: <http://purl.org/science/owl/sciencecommons/>
307
308 GENERATE {
309   ?id_json :name ?name_json ;
310           :year ?year_json ;
311           :director ?director_json ;
312           :country ?country_json .
313
314   GENERATE {
315     ?id_json :screenwriter ?screenwriter_json .
316   }
317   ITERATOR iter:Split(?screenwriters_json, ",")

```



```

318     AS ?screenwriters_json_iterator
319 WHERE {
320     BIND(REPLACE(?screenwriters_json_iterator,
321         "\\[[\\]]|\\\"", ""))
322     AS ?screenwriter_json)
323 } .
324
325 ?id_xml :name ?name_xml ;
326       :year ?year_xml ;
327       :director ?director_xml ;
328       :country ?country_xml .
329
330 GENERATE {
331     ?id_xml :screenwriter ?screenwriter_xml .
332 }
333 ITERATOR iter:XPath(?film_xml,
334     "/film/screenwriters[*]/screenwriter")
335     AS ?screenwriters_xml_iterator
336 WHERE {
337     BIND(fun:XPath(?screenwriters_xml_iterator,
338         "/screenwriter/text()") AS ?screenwriter_xml)
339 } .
340
341 }
342 ITERATOR iter:JSONPath(
343     <https://raw.githubusercontent.com/herminiogg/ShExML/
344     master/src/test/resources/filmsPaper.json>,
345     "$.films[*]") AS ?film_json
346 ITERATOR iter:XPath(
347     <https://raw.githubusercontent.com/herminiogg/ShExML/
348     master/src/test/resources/filmsPaper.xml>,
349     "//film") AS ?film_xml
350 WHERE {
351     BIND(IRI(CONCAT("http://example.com/",
352         STR(fun:JSONPath(?film_json,"$.id")))) AS ?id_json)
353     BIND(fun:JSONPath(?film_json, "$.title") AS ?name_json)
354     BIND(fun:JSONPath(?film_json, "$.director")
355         AS ?director_json)
356     BIND(IRI(CONCAT("http://example.com/",
357         fun:JSONPath(?film_json, "$.date")) AS ?year_json)
358     BIND(fun:JSONPath(?film_json, "$.countryOfOrigin")
359         AS ?country_json)
360     BIND(fun:JSONPath(?film_json, "$.director")
361         AS ?directors_json)
362     BIND(fun:JSONPath(?film_json, "$.screenwriter")
363         AS ?screenwriters_json)
364     BIND(IRI(CONCAT("http://example.com/",
365         fun:XPath(?film_xml, "/film/@id")) AS ?id_xml)
366     BIND(fun:XPath(?film_xml, "/film/name/text()")
367         AS ?name_xml)
368     BIND(fun:XPath(?film_xml, "/film/director/text()")
369         AS ?director_xml)
370     BIND(IRI(CONCAT("http://example.com/",
371         fun:XPath(?film_xml, "/film/year/text()"))
372         AS ?year_xml)
373     BIND(fun:XPath(?film_xml, "/film/country/text()")
374         AS ?country_xml)
375 }
376

```

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

3.3 ShExML

Listing 6. ShExML transformation script for the films example

```

385 PREFIX : <http://example.com/>
386 SOURCE films_xml_file <
387     https://raw.githubusercontent.com/herminiogg/
388     ShExML/master/src/test/resources/filmsPaper.xml>
389 SOURCE films_json_file <
390     https://raw.githubusercontent.com/herminiogg/
391     ShExML/master/src/test/resources/filmsPaper.json>
392 ITERATOR film_xml <xpath: //film> {
393     FIELD id <@id>
394     FIELD name <name>
395     FIELD year <year>
396     FIELD country <country>
397     FIELD director <director>
398     FIELD screenwriters <screenwriters/screenwriter>
399 }
400 ITERATOR film_json <jsonpath: $.films[*]> {
401     FIELD id <id>
402     FIELD name <title>
403     FIELD year <date>
404     FIELD country <countryOfOrigin>
405     FIELD director <director>
406     FIELD screenwriters <screenwriter>
407 }
408 EXPRESSION films <films_xml_file.film_xml
409     UNION films_json_file.film_json>
410
411 :Films :[films.id] {
412     :name [films.name] ;
413     :year [films.year] ;
414     :country [films.country] ;
415     :director [films.director] ;
416     :screenwriter [films.screenwriters] ;
417 }
418

```

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

3.4 Comparing expressiveness of the languages

Regarding the expressiveness—defined as “the breadth of ideas that can be described in a language” (Leitão and Proença, 2014)—comparison we are going to see which statements can be expressed in one of the languages and how they can be expressed—or not—in the other languages.

Iterators, sources, fields, unions and so on are common to the three languages as they have the same objective. They have different syntaxes, as it can be seen in the three examples, but from an expressiveness point of view these features have no differences.

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

Multiple results: The handling of multiple results, like it occurs on the screenwriters case, is different between SPARQL Generate and the two other languages. In YARRRML and ShExML if a query returns

multiple results they are treated like a list of them. However, in SPARQL Generate this functionality must be explicitly declared like it can be seen in Listing 5. It leads to complex iterator definitions like the one used in JSON screenwriters one.

Transformations: The possibility of transforming the output to another values by means of a function is something very useful for different purposes when building a knowledge graph. Therefore, in YARRRML this is supported through the FnO mechanism (Meester et al., 2017) which offers a way to define functions inside mapping languages in a declarative fashion. SPARQL Generate offers some functions for strings embedded inside the SPARQL binding functions mechanism; however, it is possible to extend the language through the SPARQL 1.1 extension mechanism. In the case of ShExML, only Matchers and String operations (see Section 4) are offered for transformation purposes.

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

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

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

Further expressiveness of SPARQL Generate: Apart from what has been presented in the previous point, SPARQL Generate, as being based on SPARQL 1.1, offers more expressiveness than the other two languages. One possibility that emerges from that is the use of the defined variables along all the script. For example, it is possible to define an iterator of numbers and then use that numbers to request different parts of an API. This versatility enables the creation of very complex and rich scripts that can cover a lot of use cases. It is natural to expect that learning to use the full capabilities of SPARQL Generate is complex, as the language is expressive. In our experiments, however, only some basic features of the language were required, and it appears that SPARQL Generate design does not allow to easily solve the easy use cases.

4 FURTHER SHEXML CONSTRUCTIONS

In this section we are going to present some ShExML features that were not showed in the previous section but can be important to understand the full language and its possibilities⁸.

4.1 Nested iterators

In ShExML it is possible to nest iterators to cover more complex structures. In this way, if a file contains nested multiple elements that must be iterated, they can be handled and the results showed. See Listing 7 for an example.

Listing 7. Example of nested iterators

```

ITERATOR film_xml <xpath: //film> {
  FIELD id <@id>
  FIELD name <name>
  FIELD year <year>
  FIELD country <country>
  ITERATOR actors <cast/actor> {

```

⁸Full examples and working demos can be consulted at:
<http://herminiogg.github.io/shexml-paper-2019-data/paper-examples/>

```

496         FIELD name <name>
497         FIELD role <role>
498     }
499 }
500

```

4.2 Shape linking

As in ShEx, in ShExML it is possible to ~~link a shape to another shape~~. In the case of ShExML it provides a syntactic sugar mechanism to avoid repeating twice—or more—the same expression. The resulting values for the objects where the shape link is placed will be extracted from the subject expression of the linked shape. An example of this syntax can be seen in Listing 8.

Listing 8. Example of shape linking

```

506 :Films :[films.id] {
507     :name [films.name] ;
508     :year :[films.year] ;
509     :country [films.country] ;
510     :actor @:Actor ;
511 }
512
513
514 :Actor :[films.actors.name] {
515     :name [films.actors.name] ;
516     :role [films.actors.role] ;
517 }
518

```

4.3 Expressions over fields

The manner of using expressions in the previous example where the iterators are treated as a whole is, as we called it, expressions over iterators. However, it is also possible to use expressions over fields. For example, the names of the films can be extracted with the statement showed in Listing 9.

Listing 9. Example of union over fields

```

523
524 EXPRESSION films_names <films_xml_file.film_xml.name
525 UNION films_json_file.film_json.name>
526

```

4.4 Basic expressions

A basic expression is the one that extracts the content of an iterator or a field but does not perform any operation with the values. Therefore, it can be used in the case of a single translation where no merging is needed. See Listing 10 for an example.

Listing 10. Example of basic expression

```

531
532 EXPRESSION films <films_xml_file.film_xml>
533

```

4.5 Join expressions

A join expression allows to extract identifiers from another source when a common attribute is present. This mechanism is inspired from RML in which it is called primary interlinking⁹. For example, if in a file A there is an id and a name but in a file B there is only a name it is possible to substitute the name of file B for the corresponding id of file A in presence of equal names (because we assume that they are the same entity). Example in Listing 11 shows how this expression can be used.

Listing 11. Example of join expression

```

540
541 EXPRESSION venues_union <
542     performances_json.performances_iterator.venues_ids
543     UNION events_xml.events_iterator.venues_names
544     JOIN performances_json.performances_iterator
545     .venues_names_json>
546

```

⁹http://rml.io/RML_examples.html

4.6 String operation expressions

A string operation expression allows to combine the values of different fields, for instance we can define a film id by concatenating the title and the year, as shown in Listing 12.

Listing 12. Example of string operation expression

```
EXPRESSION films_ids <films_xml_file.film_xml.name
+ "-" + films_xml_file.film_xml.year>
```

4.7 Matchers

Matchers offer the possibility to change some result to another one. They are intended for the matching of existing IRIs in the LOD Cloud and that is the reason for their name. So, in case of having the name of a country using acronyms we can match it with the local part existing in some online repository (e.g., Wikidata, DBpedia, Geonames, etc.). Listing 13 shows an example of this functionality for the countries case.

Listing 13. Example of matcher

```
MATCHER country_names <USA, US AS United_States &
FR, FRA AS France>

:Films :[films.id] {
  :country dbr:[films.country MATCHING country_names] ;
}
```

5 METHODOLOGY

In order to test the hypothesis mentioned at the beginning of this paper, an experiment was carried out. The University of Oviedo granted ethical approval to carry out the described study. Verbal consent was requested before starting the experiment.

5.1 Experiment protocol

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

A small manual was developed to guide the students along the experiment and to inform them about the input files and which are the expected outputs¹⁴. This manual contained two tasks to perform during the experiment which were designed to be performed sequentially, i.e., the student should finish the first task before starting with the second one. The first task was the mapping and integration of two files (JSON and XML) with information about books. The final output should be equal to the one given in the guide. The second task was to modify the script done in the previous task so that the prices are separated and can be compared between markets.

The study was designed as a mixed method approach, including a quantitative analysis and a qualitative analysis. For the quantitative analysis measures, Mousotron¹⁵ was used which allows to register the number of keystrokes, the distance travelled by the mouse and so on. In addition, the elapsed time was calculated from the Office 365 forms. For the qualitative analysis two Office 365 forms were used with questions based on a Likert scale (see questions in Table 1).

¹⁰<http://rml.io/yarrmrl/>

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.blacksunsoftware.com/mousotron.html>

5.2 Conduction

The sample consisted on 20 students (4 women and 13 men) of the MSc in Web Engineering first course at the University of Oviedo¹⁶. Most of them have a bachelor degree (240 ECTS credits) in computer science field or similar fields. They were receiving a semantic web course of two weeks—a total of 30 hours (3 hours per day)—where they were introduced to semantic technologies like: RDF, SPARQL, ShEx, etc. Before this course they had not previous knowledge on semantic web technologies. The experiment was hosted the final day of the mentioned course.

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

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

1. Open the assigned tool on the dedicated webpage and clear the given example.
2. Open Mousotron and reset it.
3. Proceed with task 1 (start time registered for elapsed time calculation).
4. Once task 1 is finished, capture Mousotron results (screenshot) and fill the first Office 365 questionnaire.
5. Reset Mousotron and proceed with task 2.
6. Once task 2 is finished, capture Mousotron results (screenshot) and fill the second Office 365 questionnaire.

5.2.1 Analysis

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

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

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

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

¹⁶<http://miw.uniovi.es/>

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

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

636 For the analysis ~~the IBM SPSS software, in its~~ version 24, was used. We planned a One Way ANOVA
637 test within the three groups in the quantitative analysis where a normal distribution was found and the
638 Kruskal-Wallis test where not. The qualitative analysis comparison between three groups was established
639 using the Kruskal-Wallis test. The report and analysis of the results was made using Field (2013) guidance
640 and using the suggested APA style as a standard manner to report statistical results.

641 **5.3 Threat to validity**

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

643 **5.3.1 Internal validity**

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

- 645 • More expertise in some specific tool: In semantic web area—as in other areas—people tend to be
646 more expert in some specific technologies and languages. The derived risk is that this expertise can
647 have an influence on final results. To alleviate this we have selected MSc students that are studying
648 the same introductory semantic web course and we have assigned the tools in a random manner.
- 649 • Not homogeneous group: It is possible that the selected group is not homogeneous on skills and
650 previous knowledge. To mitigate this we have applied the same measures as for the previous threat:
651 Students of a semantic web course and a randomised tool assignment.
- 652 • Unfamiliar environment: In usability studies, unfamiliar environments can play a role on final
653 conclusions. Therefore, we opted to run the experiment in a well-known environment for the
654 students, that is, their whole-year classroom.
- 655 • More guide and information about one tool: As we have designed one of the languages, it could
656 lead to a bias in information delivery. To try to mitigate this threat we developed three identical
657 manuals for each tool. Questions and doubts were answered equally for all the students and tools.

658 **5.3.2 External validity**

659 Following the measures taken in the internal validity threats we identified the corresponding external
660 validity ones:

- 661 • Very focused sample: As we have restricted the profile of the sample to students of a MSc course
662 which are more or less within the same knowledge level, there is the risk that these findings cannot
663 be extrapolated for other samples or populations. It is possible that for Semantic Web practitioners—
664 with different interests and expertises—these findings are not applicable. However, the intention of
665 this study was to evaluate usability on first-time users as a first step for future and bigger analysis.

666 **6 RESULTS**

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

671 **Task 1:** As previously stated, the number of students that finished—correctly or not—the proposed
672 task was 17. Descriptive statistics can be seen in Table 2. Comparison of three groups was made
673 by means of a One Way ANOVA which results showed significant differences on elapsed seconds
674 $F(2, 14) = 6.00, p = .013, \omega = .60$. As completeness percentage and precision are not following a

¹⁷Original datasets available on:

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

Table 1. 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

normal distribution on SPARQL Generate group ($W(4) = .63, p = .001$ and $W(4) = .63, p = .001$), the comparison was established by means of the Kruskal-Wallis test which showed significant differences in both variables ($H(2) = 9.73, p = .008$ and $H(2) = 9.68, p = .008$). Post hoc test for elapsed seconds using the Gabriel's criterion showed significant differences between ShExML group and YARRRML group ($p = .016$). Post hoc test for completeness percentage and precision using the Bonferroni's criterion showed significant differences between ShExML and SPARQL Generate ($p = .012, r = .87$ and $p = .012, r = .87$). Likert scale questionnaire results ($\alpha = 0,73$) (see Fig. 1) were analysed using Kruskal-Wallis test which resulted in significant differences between groups for variables general satisfaction level ($H(2) = 6.28, p = .043$), easiness of use ($H(2) = 9.82, p = .007$), mapping definition easiness ($H(2) = 10.25, p = .006$) and learnability ($H(2) = 8.63, p = .013$). Bonferroni's criterion was used as post hoc test for the variables with significant differences. For general satisfaction level significant differences were found between ShExML and YARRRML ($p = .039, r = .69$). For easiness of use significant differences were found between ShExML and YARRRML ($p = .011, r = .81$). For mapping definition easiness significant differences were found between ShExML and SPARQL Generate ($p = .013, r = .90$) and between ShExML and YARRRML ($p = .037, r = .69$). For learnability significant differences were found between ShExML and SPARQL Generate ($p = .042, r = .78$) and between ShExML and YARRRML ($p = .040, r = .69$).

Task 2: In this task only 7 students reached this step: 6 for ShExML and 1 for YARRRML. Descriptive statistics of this task can be seen in Table 3. No significant differences were found in any of the variables. In subjective variable analysis (see Fig. 2) no significant differences were found.

7 DISCUSSION

7.1 Statistical results discussion

Results of task 1 show that variables like keystrokes, left button clicks, right button clicks, mouse wheel scroll and meters travelled by the mouse, do not have a significant variability depending on the used tool. This suggests that web interfaces used as online development environments are more or less homogeneous and do not have an impact on the development of the scripts. However, keystrokes variable results should be taken cautiously because on SPARQL Generate the mean of completeness percentage was very low; therefore, the final solution may be longer. On its side, elapsed seconds, completeness percentage and precision have significant differences between groups which demonstrate that the selected language has an influence on these variables. Moreover, we can see that elapsed seconds has a medium size effect ($\omega = .60$). Post hoc results show that there are significant differences between ShExML and YARRRML which suggests that YARRRML users tend to use more time than ShExML users. In the case of comparisons with SPARQL Generate there are not significant differences which can be due to the small sample size and the low completeness percentage. Differences between ShExML and SPARQL Generate on completeness percentage and precision remark that SPARQL Generate users were not able to reach the same solutions as ShExML users which have the highest mean on both variables. However, between ShExML and YARRRML groups there were not such significant differences which is in line with the great variability of those two variables.

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

Table 2. 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

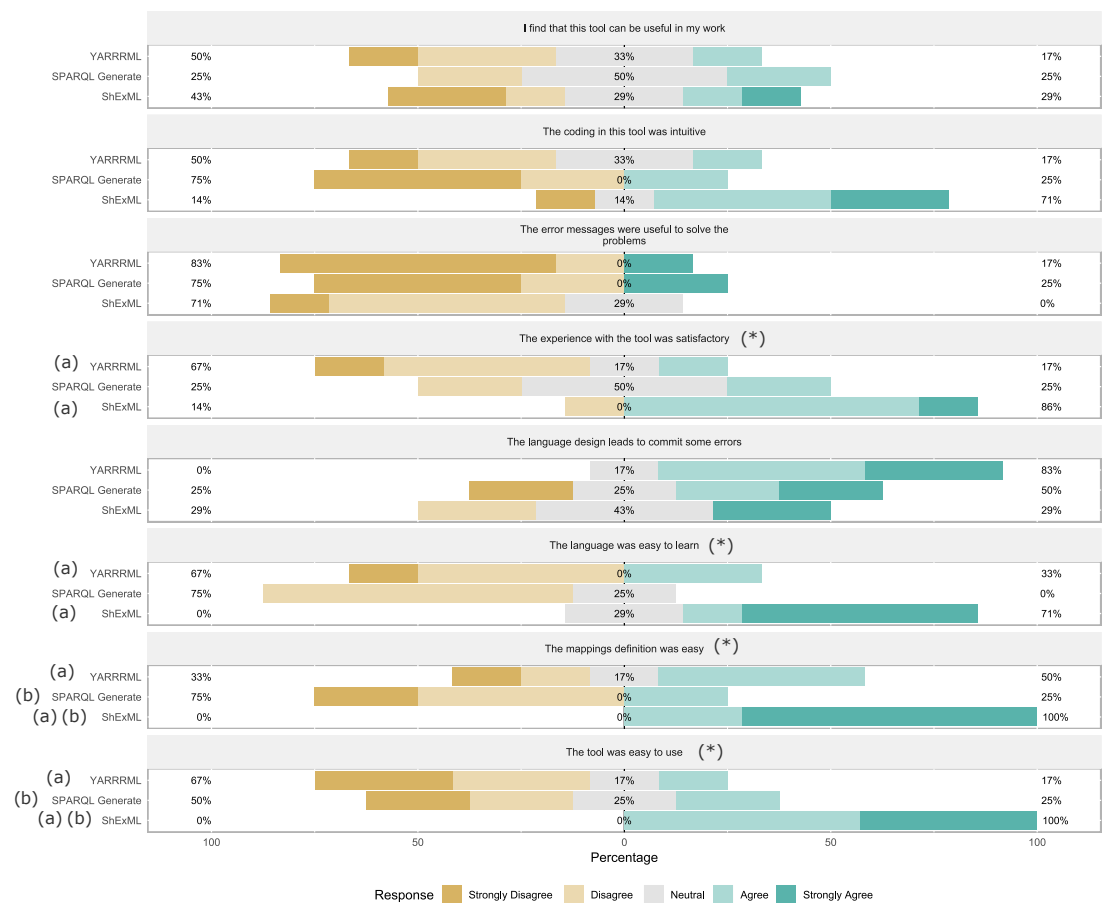


Figure 1. Task 1 results for Likert scale questionnaire where results are divided into questions and groups. (*) means significant differences between groups and (a) and (b) means significant differences in the post hoc test between the marked groups at the level of significance ($\alpha = .05$)

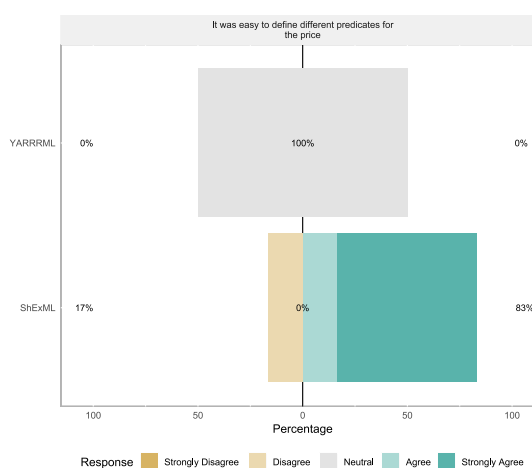


Figure 2. Task 2 results for Likert scale questionnaire where results are divided into the two groups.

Table 3. 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

only one individual made this step. However, completeness percentage and precision show us that some students ~~reach the final~~ solution with ShExML, whereas in YARRRML group and in SPARQL group they did not. This leads to a conclusion that only the ShExML group manage to find a solution to both proposed tasks. Nevertheless, these conclusions must be validated with bigger experiments to have valid statistical significances.

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

Regarding the subjective analysis, significant differences were found between groups in general satisfaction level, difficulty of use, easiness of use and **learnability**.

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

In the case of easiness of use significant differences were found between ShExML and YARRRML which suggests that ShExML users found this language easier to use than YARRRML users did with their language counterpart. In this case, like in the previous variable, differences could not be established between SPARQL Generate and the two other groups ~~which also caught attention~~ due to their results ~~in the commanded task~~. In mapping definition easiness differences were established between ShExML group and YARRRML group and between ShExML group and SPARQL Generate group which indicates that ShExML users found mappings easier to define in ShExML than in the other two languages. ~~Also to mention~~ that users did not find differences on mapping definition easiness between YARRRML and SPARQL Generate ~~which may be motivated by the reason that~~ SPARQL Generate users did not use the whole language.

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

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

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

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

Looking to modifiability variable of task 2, ShExML users tend to rate this feature with high marks and YARRRML user gave a response of 3 in a 5 point Likert scale which is in line with his/her completeness percentage mark. As with the objective results of task 2, subjective results should be ~~expanded in future~~ bigger experiments to corroborate these early findings.

7.2 Alignment with expressiveness comparison

In the light of the statistical analysis outcome, SPARQL Generate design has showed to have a bad impact on first-time users ~~usability~~. This has led to the abandon of three users ~~plus the impossibility to finish~~ of the rest of the group. Although **expressiveness is a good feature on a language, these results caught attention on how this expressiveness should be carefully designed in the language in order to not affect variables such as the one measured in this study and, therefore, affect other variables like the adoption of the tool**. In the case of YARRRML language, although it has been designed with human-friendliness in mind, in our experiment it has not reached the expected results in comparison with ShExML. However, it has better results than SPARQL Generate ~~with a decrease on expressiveness but a greater expressiveness~~ than ShExML. Nevertheless, it does not seem that expressiveness could explain the differences between

YARRRML and ShExML as the features used on the experiment do not differ on expressiveness. Instead other syntax details may be affecting the differences between these two groups such as: the use of keywords that made the language more self explanatory and the modularity used on iterators which reminds of object-oriented programming languages.

These results highlight the importance on how expressive features are designed and included in a language. Therefore, SPARQL Generate with a high expressiveness and being a highly flexible language tends to have a bad influence on users' usability due to its verbose syntax. Comparing ShExML and YARRRML we see that these differences are smaller than with SPARQL Generate and that expressiveness does not seem to be the variable affecting YARRRML usability. Thus, we can conclude—and answer the RQ2—that it is not the expressiveness of the language which affects usability of first-time users but its design.

8 CONCLUSIONS AND FUTURE WORK

In this work we have established a comparison between three heterogeneous data mapping languages on the usability level. Statistical results have shown that better results, and speed on finding this solution, are related to ShExML users whereas SPARQL Generate users were not able to find any solution. In the case of YARRRML users, on the analysis of the mentioned variable, they performed better than SPARQL Generate users but worse than ShExML users finding partial solutions to the given problem.

This study gives us a first approach to the topic of usability on first-time users with programming and Linked Data background in this kind of languages and how the expressiveness can have an influence on these results. It also reflects the importance that usability has on the accuracy of the encountered solutions. Therefore, as this study suffers from a low sample size boundary which could limit the extensibility of the found answers; as future work, there is the need to expand this experiment with higher sample sizes and different use cases that could verify these early findings.

Alongside with a bigger experiment, it is also interesting to invest on the improvement of these aspects that were worst rated in the three languages; namely: error proneness and the error reporting system.

This work highlights the importance of usability on this kind of languages and how it could affect the adoption of these tools.

ACKNOWLEDGEMENTS

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

FUNDING

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

REFERENCES

- Auer, S., Dietzold, S., Lehmann, J., Hellmann, S., and Aumüller, D. (2009). Triplify: light-weight linked data publication from relational databases. In *Proceedings of the 18th International Conference on World Wide Web, WWW 2009, Madrid, Spain, April 20-24, 2009*, pages 621–630.
- Battle, S. (2004). Round-tripping between XML and RDF. In *International Semantic Web Conference (ISWC), Hiroshima, Japan*.
- Battle, S. (2006). Gloze: XML to RDF and back again. In *Proceedings of the First Jena User Conference*.
- Berners-Lee, T., Hendler, J., and Lassila, O. (2001). The semantic web. *Scientific american*, 284(5):28–37.
- Bischof, S., Decker, S., Krennwallner, T., Lopes, N., and Polleres, A. (2012). Mapping between RDF and XML with XSPARQL. *Journal on Data Semantics*, 1(3):147–185.
- Bizer, C. and Seaborne, A. (2004). D2rq-treating non-rdf databases as virtual rdf graphs. In *Proceedings of the 3rd international semantic web conference (ISWC2004)*, volume 2004. Proceedings of ISWC2004.

- 817 Boneva, I., Labra Gayo, J. E., and Prud'hommeaux, E. G. (2017). Semantics and validation of shapes
818 schemas for rdf. In d'Amato, C., Fernandez, M., Tamma, V., Lecue, F., Cudré-Mauroux, P., Sequeda,
819 J., Lange, C., and Heflin, J., editors, *The Semantic Web – ISWC 2017*, pages 104–120, Cham. Springer
820 International Publishing.
- 821 Breitling, F. (2009). A standard transformation from XML to RDF via XSLT. *Astronomische Nachrichten*,
822 330(7):755–760.
- 823 Das, S., Sundara, S., and Cyganiak, R. (2012). R2RML: RDB to RDF Mapping Language. [https :
824 //www.w3.org/TR/r2rml/](https://www.w3.org/TR/r2rml/).
- 825 Deursen, D. V., Poppe, C., Martens, G., Mannens, E., and de Walle, R. V. (2008). XML to RDF
826 Conversion: A Generic Approach. In *Automated solutions for Cross Media Content and Multi-channel
827 Distribution, 2008. AXMEDIS '08. International Conference on*, pages 138–144, Washington.
- 828 Dimou, A., Sande, M. V., Colpaert, P., Verborgh, R., Mannens, E., and de Walle, R. V. (2014). RML: A
829 generic language for integrated RDF mappings of heterogeneous data. In *Proceedings of the Workshop
830 on Linked Data on the Web co-located with the 23rd International World Wide Web Conference (WWW
831 2014)*, Seoul, Korea, April 8, 2014.
- 832 Ermilov, I., Auer, S., and Stadler, C. (2013). CSV2RDF: User-driven CSV to RDF mass conversion
833 framework. In *Proceedings of the ISEM*, volume 13, pages 04–06, Graz, Austria.
- 834 Fagin, R., Kolaitis, P. G., Miller, R. J., and Popa, L. (2005). Data exchange: semantics and query
835 answering. *Theor. Comput. Sci.*, 336(1):89–124.
- 836 Field, A. (2013). *Discovering statistics using IBM SPSS statistics*. Sage.
- 837 Fiorelli, M., Lorenzetti, T., Paziienza, M. T., Stellato, A., and Turbati, A. (2015). Sheet2rdf: a flexible and
838 dynamic spreadsheet import&lifting framework for RDF. In *Current Approaches in Applied Artificial
839 Intelligence - 28th International Conference on Industrial, Engineering and Other Applications of
840 Applied Intelligent Systems, IEA/AIE 2015, Seoul, South Korea, June 10-12, 2015, Proceedings*, pages
841 131–140.
- 842 Freire, F., Freire, C., and Souza, D. (2017). Enhancing JSON to RDF data conversion with entity type
843 recognition. In *Proceedings of the 13th International Conference on Web Information Systems and
844 Technologies, WEBIST 2017, Porto, Portugal, April 25-27, 2017*, pages 97–106.
- 845 García-González, H., Fernández-Álvarez, D., and Gayo, J. E. L. (2018). ShExML: An Heterogeneous
846 Data Mapping Language based on ShEx. In *Proceedings of the EKAU 2018 Posters and Demonstrations
847 Session co-located with 21st International Conference on Knowledge Engineering and Knowledge
848 Management (EKAU 2018)*, Nancy, France, November 12-16, 2018., pages 9–12.
- 849 Halevy, A. Y. (2001). Answering queries using views: A survey. *The VLDB Journal*, 10(4):270–294.
- 850 Han, L., Finin, T., Parr, C. S., Sachs, J., and Joshi, A. (2008). RDF123: from spreadsheets to RDF. In
851 *The Semantic Web - ISWC 2008, 7th International Semantic Web Conference, ISWC 2008, Karlsruhe,
852 Germany, October 26-30, 2008. Proceedings*, pages 451–466.
- 853 Hert, M., Reif, G., and Gall, H. C. (2011). A comparison of rdb-to-rdf mapping languages. In *Proceedings
854 of the 7th International Conference on Semantic Systems*, pages 25–32. ACM.
- 855 Heyvaert, P., Dimou, A., Herregodts, A., Verborgh, R., Schuurman, D., Mannens, E., and de Walle, R. V.
856 (2016). Rmleditor: A graph-based mapping editor for linked data mappings. In *The Semantic Web.
857 Latest Advances and New Domains - 13th International Conference, ESWC 2016, Heraklion, Crete,
858 Greece, May 29 - June 2, 2016, Proceedings*, pages 709–723.
- 859 Heyvaert, P., Meester, B. D., Dimou, A., and Verborgh, R. (2018). Declarative rules for linked data
860 generation at your fingertips! In *The Semantic Web: ESWC 2018 Satellite Events - ESWC 2018 Satellite
861 Events, Heraklion, Crete, Greece, June 3-7, 2018, Revised Selected Papers*, pages 213–217.
- 862 Lefrançois, M., Zimmermann, A., and Bakerally, N. (2016). Flexible RDF generation from RDF
863 and heterogeneous data sources with sparql-generate. In *Knowledge Engineering and Knowledge
864 Management - EKAU 2016 Satellite Events, EKM and Drift-an-LOD, Bologna, Italy, November 19-23,
865 2016, Revised Selected Papers*, pages 131–135.
- 866 Lefrançois, M., Zimmermann, A., and Bakerally, N. (2017). A SPARQL extension for generating RDF
867 from heterogeneous formats. In *The Semantic Web - 14th International Conference, ESWC 2017,
868 Portorož, Slovenia, May 28 - June 1, 2017, Proceedings, Part I*, pages 35–50.
- 869 Leitão, A. and Proença, S. (2014). On the expressive power of programming languages for generative
870 design-the case of higher-order functions.
- 871 Meester, B. D., Heyvaert, P., Verborgh, R., and Dimou, A. (2019). Mapping languages: Analysis of

- comparative characteristics. In Chaves-Fraga, D., Heyvaert, P., Priyatna, F., Sequeda, J. F., Dimou, A., Jabeen, H., Graux, D., Sejdiu, G., Saleem, M., and Lehmann, J., editors, *Joint Proceedings of the 1st International Workshop on Knowledge Graph Building and 1st International Workshop on Large Scale RDF Analytics co-located with 16th Extended Semantic Web Conference (ESWC 2019)*, Portorož, Slovenia, June 3, 2019, volume 2489 of *CEUR Workshop Proceedings*, pages 37–45. CEUR-WS.org.
- Meester, B. D., Maroy, W., Dimou, A., Verborgh, R., and Mannens, E. (2017). RML and fno: Shaping dbpedia declaratively. In *The Semantic Web: ESWC 2017 Satellite Events - ESWC 2017 Satellite Events*, Portorož, Slovenia, May 28 - June 1, 2017, *Revised Selected Papers*, pages 172–177.
- Michel, F., Djimenou, L., Faron-Zucker, C., and Montagnat, J. (2015). Translation of relational and non-relational databases into RDF with xr2rml. In Monfort, V., Krempels, K., Majchrzak, T. A., and Turk, Z., editors, *WEBIST 2015 - Proceedings of the 11th International Conference on Web Information Systems and Technologies*, Lisbon, Portugal, 20-22 May, 2015, pages 443–454. SciTePress.
- Michel, F., Montagnat, J., and Zucker, C. F. (2014). *A survey of RDB to RDF translation approaches and tools*. PhD thesis, I3S.
- Miletic, I., Vujasinovic, M., Ivezic, N., and Marjanovic, Z. (2007). Enabling Semantic Mediation for Business Applications: XML-RDF, RDF-XML and XSD-RDFS transformations. In Gonçalves, R. J., Müller, J. P., Mertins, K., and Zelm, M., editors, *Enterprise Interoperability II: New Challenges and Approaches*, pages 483–494. Springer London, London.
- Müller, H., Cabral, L., Morshed, A., and Shu, Y. (2013). From restful to SPARQL: A case study on generating semantic sensor data. In *Proceedings of the 6th International Workshop on Semantic Sensor Networks co-located with the 12th International Semantic Web Conference (ISWC 2013)*, Sydney, Australia, October 22nd, 2013., pages 51–66.
- Prud'hommeaux, E., Labra Gayo, J. E., and Solbrig, H. (2014). Shape Expressions: An RDF Validation and Transformation Language. In *Proceedings of the 10th International Conference on Semantic Systems, SEM '14*, pages 32–40, New York, NY, USA. ACM.
- Reinsel, D., Gantz, J., and Rydning, J. (2018). The Digitization of the World. From Edge to Core. Technical report, Seagate, IDC. Last accessed: October 28, 2019.
- Sahoo, S. S., Halb, W., Hellmann, S., Idehen, K., Thibodeau Jr, T., Auer, S., Sequeda, J., and Ezzat, A. (2009). A survey of current approaches for mapping of relational databases to rdf. *W3C RDB2RDF Incubator Group Report*, 1:113–130.
- Scharffe, F., Bihanic, L., Képéklian, G., Atemezing, G., Troncy, R., Cotton, F., Gandon, F., Villata, S., Euzenat, J., Fan, Z., Bucher, B., Hamdi, F., Vandenbussche, P., and Vatan, B. (2012). Enabling linked data publication with the datalift platform. In *Semantic Cities, Papers from the 2012 AAAI Workshop*, Toronto, Ontario, Canada, July 22-23, 2012.
- Sequeda, J. F., Arenas, M., and Miranker, D. P. (2012). On directly mapping relational databases to RDF and OWL. In *Proceedings of the 21st World Wide Web Conference 2012, WWW 2012, Lyon, France, April 16-20, 2012*, pages 649–658.
- Slepicka, J., Yin, C., Szekely, P. A., and Knoblock, C. A. (2015). KR2RML: an alternative interpretation of R2RML for heterogenous sources. In Hartig, O., Sequeda, J. F., and Hogan, A., editors, *Proceedings of the 6th International Workshop on Consuming Linked Data co-located with 14th International Semantic Web Conference (ISWC 2015)*, Bethlehem, Pennsylvania, USA, October 12th, 2015, volume 1426 of *CEUR Workshop Proceedings*. CEUR-WS.org.
- Sperberg-McQueen, C. M. and Miller, E. (2004). On mapping from colloquial XML to RDF using XSLT. In *Extreme Markup Languages®*.
- Stadler, C., Unbehauen, J., Westphal, P., Sherif, M. A., and Lehmann, J. (2015). Simplified RDB2RDF mapping. In Bizer, C., Auer, S., Berners-Lee, T., and Heath, T., editors, *Proceedings of the Workshop on Linked Data on the Web, LDOW 2015, co-located with the 24th International World Wide Web Conference (WWW 2015)*, Florence, Italy, May 19th, 2015, volume 1409 of *CEUR Workshop Proceedings*. CEUR-WS.org.
- Tandy, J., Herman, I., and Kellogg, G. (2015). Generating rdf from tabular data on the web, w3c recommendation 17 december 2015. w3c recommendation. <https://www.w3.org/TR/2015/REC-csv2rdf-20151217>.
- Theocharis, S. and Tsihrintzis, G. A. (2016). Rdf serialization from JSON data: The case of JSON data in diavgeia.gov.gr. In *7th International Conference on Information, Intelligence, Systems & Applications, IISA 2016, Chalkidiki, Greece, July 13-15, 2016*, pages 1–6.

- 927 Thuy, P. T. T., Lee, Y.-K., Lee, S., and Jeong, B.-S. (2007). Transforming valid XML documents into RDF
928 via RDF schema. In *Next Generation Web Services Practices, 2007. NWeSP 2007. Third International*
929 *Conference on*, pages 35–40. IEEE.
- 930 Thuy, P. T. T., Lee, Y.-K., Lee, S., and Jeong, B.-S. (2008). Exploiting XML schema for interpreting
931 XML documents as RDF. In *Services Computing, 2008. SCC'08. IEEE International Conference on*,
932 volume 2, pages 555–558. IEEE.