

1 DataGen: JSON/XML Dataset Generator

2 **Filipa Alves dos Santos**

3 University of Minho, Portugal

4 a83631@alunos.uminho.pt

5 **Hugo André Coelho Cardoso**

6 University of Minho, Portugal

7 a85006@alunos.uminho.pt

8 **João da Cunha e Costa**

9 University of Minho, Portugal

10 a84775@alunos.uminho.pt

11 **Válter Ferreira Picas Carvalho**

12 University of Minho, Portugal

13 a84464@alunos.uminho.pt

14 **José Carlos Ramalho** 

15 Department of Informatics, University of Minho, Portugal

16 jcr@di.uminho.pt

17 — Abstract —

18 In this document, it is described the steps towards DataGen’s implementation.

19 DataGen is a versatile and powerful tool that allows for quick prototyping and testing of software
20 applications, since currently too few solutions offer both the complexity and scalability necessary
21 to generate adequate datasets in order to feed a data API or a more complex APP, enabling their
22 testing with appropriate data volume and complexity.

23 DataGen’s core is a Domain Specific Language (DSL) that was created to specify datasets.
24 This language suffered several updates: repeating fields (with no limit), fuzzy fields (statistically
25 generated), lists, high order functions over lists, custom made transformation functions. The final
26 result is a diversified algebra that allows the generation of very complex datasets coping with very
27 convoluted requirements. Throughout the paper, several examples of the possibilities will be given.

28 After generating a dataset, DataGen gives the user the possibility to generate a RESTful data
29 API with it, creating a running prototype.

30 This solution has already been used in real life cases, described with more detail throughout
31 the paper, in which it was able to create the intended datasets successfully. These allowed the
32 application’s performance to be tested and for the right adjustments to be made.

33 The tool is currently being deployed for general use.

34 **2012 ACM Subject Classification** Software and its engineering → Domain specific languages; Theory
35 of computation → Grammars and context-free languages; Information systems → Open source
36 software

37 **Keywords and phrases** JSON, XML, Data Generation, Open Source, REST API, Strapi, JavaScript,
38 Node.js, Vue.js, Scalability, Fault Tolerance, Dataset, DSL, PEG.js, MongoDB

39 **Digital Object Identifier** 10.4230/OASICS.SLATE.2021.5

40 **1** Introduction

41 Every application and software developed should be thoroughly tested before release, in order
42 to determine the system’s ability to withstand realistic amounts of data and traffic, and
43 that implies the usage of representative datasets that fit its use cases. The creation of said
44 datasets is a laborious and drawn out process, as it implies firstly generating the test data
45 in some way, in a file format compatible with the system. As it stands, there are currently



© F. Santos, H. Cardoso, J. Costa, V. Carvalho, J.C. Ramalho;
licensed under Creative Commons License CC-BY

10th Symposium on Languages, Applications and Technologies (SLATE 2021).

Editors: Ricardo Queirós, Mário Pinto, Alberto Simões, Filipe Portela, and Maria João Pereira; Article No. 5;
pp. 5:1–5:15



OpenAccess Series in Informatics

OASICS Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

46 no efficient, intuitive and scalable tools to do this and so, developers often end up having
47 to create test records either by hand, which is incredibly inefficient and time-consuming, or
48 by using existing tools with some clever tactics to manage their shortcomings. As a matter
49 of fact, many projects are not able to progress to the development stage due to the lack of
50 adequate and sufficient data [3].

51 Even with a reliable generation tool, the user might want to have control over the data in
52 the resulting records, being able to observe and manipulate it freely, through CRUD requests,
53 and adapting it however they want. Though there are products capable of doing this - for
54 example the package json-server, which creates a full fake REST API -, its usage entails the
55 user manually starting the application every time they want to edit the dataset's contents
56 and ensuring the data's format is compliant with the software they're using, which ends up
57 requiring a lot of extra effort and setup on the user's side.

58 As such, the team came up with the idea of coupling the generation process and the
59 REST API, in a way that allows the user to automatically generate data compatible with an
60 integrated RESTful API - for which the software Strapi was chosen, as will be explained in
61 more detail in section 3.4 -, allowing the user to load the records into an API server and
62 perform CRUD operations over it, either through the user interface or via HTTP requests.

63 This paper will cover the development of the resulting application, DataGen - a more
64 versatile and powerful tool for data generation, according to the user's specifications, and
65 subsequent handling -, seeking to explain the decisions that were taken for its architecture,
66 as well as all the implemented functionalities.

67 **2 Related Work**

68 User privacy concerns [13] have been a major topic of discussion over the last decades, which
69 lead to the creation of strict regulations regarding the way sensitive data should be handled,
70 such as the EU's General Directive on Data Protection GDPR [2]. These regulations keep
71 entities from disclosing private and sensitive information, which in turn hurts new growing
72 ideas and projects that would require access to similar data. As it stands, not only is the
73 lack of available data a big problem in this context, as are the data sharing agreements
74 themselves, as their ratification tends to take, on average, a year and a half [10], which proves
75 to be fatal for many small and upcoming companies.

76 To circumvent these concerns, organisations have been increasingly adopting synthetic
77 data generation [16], an approach that was originally suggested by Rubin in 1993 [1] in
78 order to create realistic data from existing models without compromising user privacy. The
79 prospect of being able to create viable information that does not relate to actual individuals is
80 a very enticing solution for the privacy conundrum. If perfected, it could potentially provide
81 the capacity of generating sizeable amounts of data according to any use cases, bypassing the
82 need to access real users' information. As such, this approach has been increasingly adopted
83 and put to the test, in order to measure its efficiency with real-life cases [5, 14, 18].

84 Dataset generation has become a requisite in many development projects. However, the
85 majority of developed tools produce datasets for specific contexts like intrusion detection
86 systems in IoT [4], crops and weed detection [6], vehicular social networks based on floating
87 car data [12], 5G channel and context metrics [17], GitHub projects [9], to cite a few. Many
88 others exist in different contexts like medicine, bioinformatics, weather forecasts, color
89 constancy, market analysis, etc.

90 Most of the researched tools are domain specific. The team's goal is to create a generic
91 tool but powerful enough to generate datasets for several and different contexts and with

92 many levels of complexity. There are some tools available, many online, but they cover the
93 generation of simple datasets, many times flat datasets.

94 The main source of inspiration for this project was an already existing web application
95 called "JSON Generator", developed by Vazha Omanashvili [11], as it is the most interesting
96 dataset generation tool that was found. It features a DSL (Domain Specific Language) that's
97 parsed and converted into a JSON dataset, allowing the user to generate an array of objects
98 that follow some predefined structure, specified in the DSL.

99 In terms of utility, it is a very powerful tool that can generate very complex data structures
100 for any application with relatively little effort. However, it had some shortcomings which
101 initially inspired the team to develop an alternative solution that attempts to address them.

102 Specifically, these shortcomings are:

- 103 1. Limited size for the 'repeat' statement (100 total). Arguably, the size of the dataset itself
104 is one of the most important features to take into account. For applications on a larger
105 scale, having a small amount of array elements does not allow for more realistic tests,
106 as very few, if any, that are developed in a production environment use as little as 100
107 entries created by the aforementioned directive.
- 108 2. It only generates JSON. Despite being one of the most popular file formats, there are
109 many others that could be useful to the user as they might want the data to be in a
110 different format for their application without converting the JSON themselves, such as
111 XML (another open text format [19]). This allows for further flexibility and expands the
112 possible use cases that it provides.
- 113 3. Does not generate a RESTful API for the dataset. Many users may optionally want their
114 newly generated dataset hosted and exposed by a RESTful API for direct usage in their
115 web applications, or to download a custom one created specifically for their dataset for
116 later deployment on a platform of their choosing.
- 117 4. Does not take into account fuzzy generation. Some elements of a dataset may not be
118 deterministic and are represented by probabilities. For example, there may a field that
119 exists only if another property has a specific value and the application should be able to
120 take that into account.
- 121 5. It does not have much data available. For instance, the user might want to use names
122 from a list of famous people for their dataset, as it allows for a more realistic generation
123 and consistency, which this tool currently does not provide.
- 124 6. It is not multilingual. The data that this application uses is only available in English, it
125 would be more user-friendly to give them the choice to use their preferred language for
126 the dataset instead of forcing it to just one.
- 127 7. Does not take into account integration on applications. The generation and download of
128 a dataset requires the consumer to use the website's interface - this is not ideal as many
129 applications may want to use HTTP requests to automate this process for internal usage.
- 130 8. Does not support functional programming features. Functions such as 'map', 'reduce' and
131 'filter' that are staple in the functional paradigm due to their simplicity and effectiveness
132 are not present in this application. This would allow the user to chain statements and
133 transform fields to the result they want, granting the application the ability to deal with
134 more complex use cases.

135 With clear goals and an initial application to take inspiration from, the team proceeded
136 to the development stage by deciding on its architecture (i.e. programming languages,
137 frameworks, external tools, etc), which will be explained in the following sections.

138 3 DataGen Development

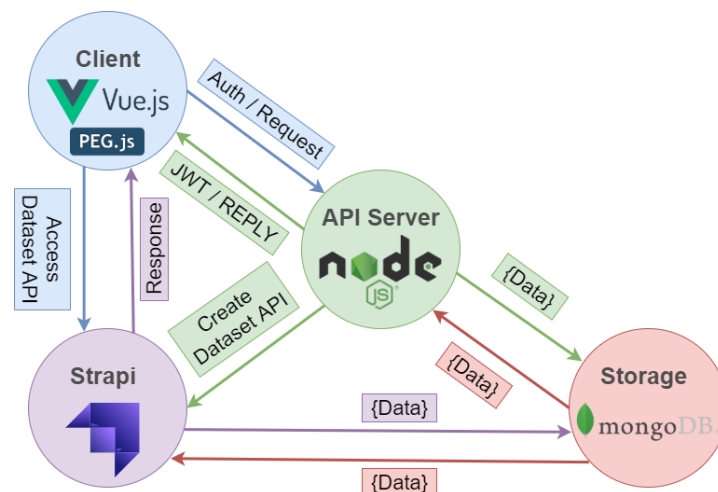
139 Building the application from the ground up requires a divide and conquer approach, since
 140 having a monolithic single server architecture will lead to a less stable user experience, due
 141 to the lack of scalability.

142 Having the application compartmentalized in multiple servers, each with their specific
 143 function, allows for a much more sensible and balanced architecture since it leaves room for
 144 the possibility of individually upgrading each of them, leading to a much higher scalability
 145 and fault tolerance, as the failure of one component does not compromise the functionality
 146 of the entire application, allowing for quicker and easier routine maintenance.

147 The following subsections will explain how the current architecture was achieved and the
 148 technological decisions behind each of the components.

149 3.1 Architecture

150 The picture below shows the general scope of the architecture, albeit simplified. This
 151 architecture allows for major upgrades, such as load balancers on both the back and front-end
 152 since they are both **stateless** (the JWT approach allows the servers to not maintain sessions
 153 with each user's state) and using MongoDB as a distributed database - sharded cluster.



■ Figure 1 Current architecture of DataGen

154 3.2 Front-End

155 The first component needed for the application is the front-end server, which is responsible
 156 for building and showing the website's interface to the user, making it the entry point to the
 157 application and its scripted behaviour.

158 3.2.1 Grammar

159 The application uses a PEG.js grammar-based parser [8] [15] to process the user's input
 160 and generate the intended dataset. The aforementioned grammar defines a domain-specific
 161 language (DSL), with JSON-like syntax, providing many features that allow for the generation
 162 of complex and diverse datasets. These features include relational and logic capabilities,

163 providing means for the datasets to satisfy several forms of constraints - which push towards
 164 the use of some declarative framework for this specification -, as well as functional capabilities,
 165 allowing for easier management and processing of certain properties of the datasets.

166 The first and most fundamental of said features is the JSON-similar syntax - the user
 167 can specify key-value properties, where the value may take any of the basic JSON types and
 168 data structures, from integers to objects and arrays. The user may also nest these to create
 169 a structure with any depth they may want.

```

name: {
  first: ["John", "William"],
  last: "Doe"
},
age: 21

```

170 To specify the size of the dataset, or a nested array, there is the **repeat** statement, where
 171 the user indicates the structure that they want replicated (which may range from a primitive
 172 JSON type to a complex object), as well as the number of copies, or range of numbers.

```

names: [ 'repeat(150,200)': {
  first: '{{firstName()}}',
  last: '{{surname()}}'
} ]

```

173 They may also specify as many collections as they want in a single model (collections are
 174 the key-value properties on the uppermost level of the model) and the application will return
 175 the resulting dataset in json-server syntax - an object with a property for each collection.
 176 During the parsing of the input, the application recursively builds both the final dataset
 177 and the Strapi model for the specified structure, concurrently, in order to allow for posterior
 178 integration in a RESTful API.

```

{
  names: [ 'repeat(10)': '{{fullName()}}' ],
  animals: [ 'repeat(20)': '{{animal()}}' ]
}

```

179 To define the value of a property, the user may also use interpolation. To access an
 180 interpolation function, it must be encased in double curly braces. There are two types of
 181 interpolation functions:

- 182 ■ functions that generate spontaneous values during execution, according to the user's
 183 instructions - for example, there is a random integer generation function where the user
 184 must provide a range of values for the intended result:

```

id: '{{objectId()}}',
int: '{{integer(50,100)}}',
random: '{{random(23, "hello", [1,2,3], true)}}'

```

5:6 DataGen

185 ■ functions that return random values from a group of datasets incorporated in the applica-
186 tion behind an API, where each dataset has information of a given category, for example
187 names and political parties.

```
name: '{{fullName()}}',  
party: '{{political_party()}}'
```

188 These interpolation functions may also be interwoven together and with normal strings to
189 generate more structured strings, such as addresses. Some of them take arguments, in which
190 case the user can either manually introduce the value or reference another property defined
191 above, through a variable **this**, allowing them to establish relations between properties.

```
parish: '{{pt_parish()}}',  
district: '{{pt_district("parish", this.parish)}}',  
address: 'St. {{fullName()}}, {{pt_city("district", this.district)}}'
```

192 In regard to the API of datasets, the team did an extensive search for useful datasets,
193 used the well-structured ones it found and salvaged whatever information it could from
194 others that were maybe less organized, processing this information to remove errors and
195 normalize its content, before joining it with other data of the same topic to create bigger,
196 more complete datasets for the user to use.

197 The team also created some original datasets by hand, for topics deemed appropriate,
198 and manually introduced bilingual support, for both portuguese and english, in all datasets
199 made available in the application, in order to let the user choose whichever language suits
200 best their goal. To indicate their language of choice, the user's model must start with the
201 following syntax:

```
<!LANGUAGE pt> or <!LANGUAGE en>
```

202 Currently, DataGen has support datasets of all the following categories: actors, animals,
203 brands, buzzwords, capitals, cities, car brands, continents, countries, cultural landmarks,
204 governmental entities, hackers, job titles, months, musicians, names, nationalities, political
205 parties, portuguese businessmen, portuguese districts, cities, counties and parishes, por-
206 tuguese public entities, portuguese politicians, portuguese public figures, portuguese top
207 100 celebrities, religions, soccer clubs, soccer players, sports, top 100 celebrities, weekdays,
208 writers.

209 The grammar also makes available a feature named **unique()**, to which the user may
210 provide an interpolation function, or a string interpolated with such a function, as an
211 argument. **unique** guarantees that the interpolation functions on which it is applied always
212 return unique values. This is especially relevant when it comes to interpolation functions
213 that fetch random data from the support datasets inside a **repeat** statement, as there is no
214 guarantee that there won't be duplicates among the fetched records and the user might not
215 want that.

216 As such, **unique** only has any effect when applied on dataset interpolation functions or
217 with **random** (which can be observed in one of the examples from last page). As long as
218 it's one of those (possibly interpolated with normal strings) and there are sufficient distinct
219 entries for the entire **repeat** statement, this tool guarantees that all objects in the resulting
220 dataset will have a different value in the property in question. If the user uses a string

221 with more than one interpolation function, there is also no effect - there may be repeated
222 combinations in the end.

223 Below are two examples: the first one depicts the correct usage of the **unique** feature;
224 the second shows instances of a wrong approach (not enough distinct values for the repeat
225 statement; not a dataset interpolation function or the random function; more than one
226 interpolation function) that will either not work or not assure any mutually exclusive
227 guarantee for the resulting values:

```
[ 'repeat(6)': {
  continent: unique('{{continent()}}'),
  country: unique('Country: {{country()}}'),
  random: unique('{{random(1,2,3,4,5,6)}}')
} ]

[ 'repeat(10)': {
  continent: unique('{{continent()}}'),
  int: unique('{{integer(5,20)}}'),
  name: unique('{{firstName()}} {{surname()}}')
} ]
```

228 Back to the properties of the model, the user may also use JavaScript functions to define
229 their value. There are two types of functions: signed functions, where the name of the
230 method corresponds to the key of the property, while the result of the body of the function
231 translates to its value, and anonymous arrow functions, which are used to indicate solely the
232 value of a property (the key needs to be precised separately beforehand).

```
name: "Oliver",
email(gen) {
  var i = gen.integer(1,30);
  return '{{this.name}}.{{gen.surname()}}${i}@gmail.com'.toLowerCase();
},
probability: gen => { return Math.random() * 100; }
```

233 Inside these functions, the user is free to write JavaScript code that will be later executed to
234 determine the value of the property. This way, more complex algorithms may be incorporated
235 into the construction logic of the dataset, allowing for a more nuanced and versatile generation
236 tool. Given that the user has access to the whole Javascript syntax, they may make use
237 of relational and logical operators to elaborate conditions on the intended data, as well as
238 functional methods (for example "map" or "filter", which Javascript implements).

239 Inside these blocks of code, the user has full access to any property declared above in
240 the DSL model, through the variable **this**, as well as any interpolation function available in
241 the parser, through a **gen** variable - whenever using a function, the user must declare this
242 argument in its signature, which they may then use inside to run said interpolation functions.
243 All of this can be observed in the example above.

244 The grammar also allows fuzzy generation of properties, i.e. constraining the existence of
245 certain properties based on logical conditions or probabilities. As of now, the grammar has
246 four different tools for this purpose:

5:8 DataGen

- 247 ■ **missing/having** statements - as an argument, they receive the probability of the prop-
248 erties inside (not) existing in the final dataset; this probability is calculated for each
249 element, originating a dataset where some elements have said properties and others don't:

```
missing(50) { prop1: 1, prop2: 2 },  
having(80) { prop3: 3 }
```

- 250 ■ **if... else if... else** statements - these work just as in any programming language: the
251 user can use relational operators and other conditional statements to create conditions
252 and string multiple of them together with the help of logical operators. The final object
253 will have the properties specified in the first condition that evaluates to true (or eventually
254 none of them, if all conditions prove to be false). In these conditions, similar to the
255 functions, the user has unrestricted access to all properties declared above in the DSL
256 model, as well as the interpolation functions, which gives them the ability to causally
257 relate different properties:

```
type: '{random("A","B","C")}',  
if (this.type == "A") { A: "type is A" }  
else if (this.type == "B") { B: "type is B" }  
else { C: "type is C" }
```

- 258 ■ the **or** statement - the grammar makes available this logical operator for quick prototyping
259 of mutually exclusive properties, where only one will be selected for each object (note
260 that it doesn't make sense to create an **and** statement, since that translates to simply
261 listing the wanted properties normally in the DSL model):

```
or() {  
  prop1: 1,  
  prop2: 2,  
  prop3: 3  
}
```

- 262 ■ the **at_least** statement - inside this, the user writes a set of properties and gives the
263 statement an integer argument specifying the minimum number of those properties that
264 must be present in the final object. The parser selects that number of properties randomly
265 from the given set:

```
at_least(2) {  
  prop1: 1,  
  prop2: 2,  
  prop3: 3  
}
```

- 266 Finally, the grammar also provides an implementation of the fundamental functional
267 programming features - **map**, **filter** and **reduce**. The user may chain together one or several
268 of these functions with an array value (from any of the various array-creating features made

269 available). Shorthand syntax is not allowed, so curly braces must always be opened for
270 the code block. Aside from that, these features work exactly like the native Javascript
271 implementations: the user may either declare the function inside or use anonymous syntax
272 for the variables; they may declare only the current value or any of the additional, albeit
273 less commonly used, variables. Examples of several possible usages of these features can be
274 observed below:

```
map: range(5).map(value => { return value+1 } ),
filter: [0,1,2].filter(function(value, index) {return [0,1,2][index]>0}),
reduce: range(5).reduce((accum, value, index, array) => {
    return accum + array[index] } ),
combined: range(5).map((value) => { return value+3 } )
    .filter(x => { return x >= 5})
    .map(x => { return x*2 } ).reduce((a,c) => {return a+c})
```

275 3.2.2 Interoperability

276 After processing the model, the parser generates an intermediary data structure with the
277 final dataset, which can then be translated to either JSON or XML, according to the user's
278 preference. The parser also generates another data structure with the corresponding Strapi
279 model, for possible later integration in its RESTful API.

280 Note that the application's purpose is to create datasets according to the user's instructions,
281 in either JSON or XML. Although the model specification may involve Javascript code, under
282 the form of functions or conditions, as explained in the previous subsection, this does not
283 correlate to the target application whatsoever. DataGen merely generates test datasets - it
284 can be used for any kind of application that accepts data in JSON/XML format, whether it
285 be an application written in Javascript, Python, C++ or some other language.

286 3.2.3 Client-Side Generation

287 This project was developed with the intent of being a web application, more specifically a
288 website with user-friendly features. A server-sided approach would imply parsing the DSL
289 models on the back-end server, which wouldn't be sensible as the created PEG.js parser
290 doesn't require access to any private services (i.e. databases) hidden by the back-end itself.

291 Therefore, a client-sided approach makes the most sense for this application in particular,
292 freeing resources (mainly the CPU and memory modules) from the back-end server and
293 shifting the computation of the generated dataset to the client in their browser, using the
294 back-end as an API server.

295 There are many frameworks aimed at client-sided browser experiences, however, it was
296 decided that Vue.js would be the most adequate for this application. It allows for reactive two-
297 way data binding - connection between model data updates and the view (UI) which creates
298 a more user-friendly experience, since it shows changes on the DOM as they occur, instead
299 of forcing a reload on the page (as in the traditional served-sided approach). Other reasons
300 such as flexibility - on the components' styling and scripted behaviour - and performance
301 - it's more efficient than React and Angular - were also a deciding factor on picking this
302 specific framework.

303 After deciding which framework to use, the team started developing the interface itself,
304 which currently has the following features:

- 305 ■ Authentication - it uses JWT (JSON Web Tokens) that are sent in the 'Authorization' header on every HTTP request that needs to be validated by the back-end (i.e accesses restricted user data);
- 306
- 307
- 308 ■ Generation and download of the dataset and/or its API - as previously mentioned, it uses a PEG.js parser for converting the DSL into JSON or XML, as well as Strapi for the API (which will be explained in section 3.4);
- 309
- 310
- 311 ■ Saving DSL models - authenticated users may save their created models and optionally make them available for others to use, being able to do CRUD operations on those they own;
- 312
- 313
- 314 ■ Documentation - since the DSL has a defined structure, it is essential that the user has access to these guidelines at any time;
- 315
- 316 ■ Team description - the user may want to contact the team directly so there is a dedicated page for them to find all of this information and a brief description of the team.
- 317

318 The front-end needs to access persistent information such as user data and saved models
319 which is accessible through the back-end's RESTful API, viewed in more detail in section 3.3.

320 3.3 Back-End

321 The application needs a back-end server for multiple purposes, namely storing data, authenticating users and generating the API for the datasets.

322

323 None of the above require intensive computational power for each request sent by the user, which was one of the main reasons why the team chose a Node.js server - it is built to deal with any incoming request that is not CPU intensive due to its single-threaded, event-driven, non-blocking IO model - and because it is scalable, has good performance in terms of speed and has a wide selection of packages (available on the **npm** registry).

324

325

326

327 For storing all the data that needs to be persistent, the back-end server accesses a MongoDB server instance, which was chosen due to its scalability (the data is not coupled relationally, which means that each document may be in a different node instance without any conflicts since they are self-contained) and direct compatibility with Node.js since they both accept JSON documents.

328

329

330

331

332 Currently the application uses three collections on the MongoDB database:

- 333 ■ **users** - stores all user specific data, which is their name, email, password (hashed) and the dates of register and most recent access;
- 334
- 335
- 336 ■ **models** - stores all DSL models and whom (user) they belong to, their visibility (public or private), title, description and register date;
- 337
- 338 ■ **blacklist** - stores users' JWT tokens after they log-out and their expiry date so that they are automatically removed from the database after they expire.
- 339

340

341 Authenticating the user allows them to access their saved DSL models and perform CRUD operations on them. Due to its Node.js integration, a JWT (JSON Web Token) approach was the chosen strategy to authenticate a user - after they submit their credentials, the system compares them to the ones saved on the database and if they match, the token is returned in the HTTP response. This token is needed for any further request that accesses critical information for that same user and expires after a short time for precaution and safety. After they log-out, it is added to a blacklist to provide extra security, since it does not allow for a user that got access to another user's JWT (if they log-out before the short expiration date) to submit requests signed with it.

342

343

344

345

346

347

348

349 Generating the API is a more complex process and has a dedicated subsection (3.4) which
350 explains the steps followed in order to obtain a fully functional REST API for any generated
351 dataset.

352 **3.4 Strapi API**

353 DataGen also provides another important functionality, which is generating a data API from
354 the dataset previously created. It's a useful feature since a lot of users may want to perform
355 CRUD operations on the data they requested or even utilize the API themselves for further
356 work.

357 The tool chosen to create this API was Strapi [20], one of the best content management
358 systems currently. Strapi automatically generates a REST API and allows multiple APIs
359 to run simultaneously. It's also very simple to configure and supports different database
360 systems like PostgreSQL, MySQL, SQLite and MongoDB, being that the latter was the one
361 used in this project. JSON-server was also considered as a tool but lacked scalability, as it
362 only allows a single API to run at a time, which wouldn't be ideal at all. However, Strapi
363 presented its own set of challenges, like the difficult way in which it stores structured data
364 (an array, for example) and how data is imported, all of which were surpassed successfully.

365 The process of building the API begins within the grammar section of the project, since
366 the Strapi model is written in a recursive way, at the same time the dataset itself is being
367 built. This strategy was a big time save in terms of the program's execution. For example,
368 any time an array is encountered, because Strapi doesn't have its own type to represent it, a
369 component is created with the array's elements and a reference to that component is written
370 in the model itself. This recursive process keeps on going with this same logic until it reaches
371 the root, which corresponds to the collection.

372 After the model is created, this data is redirected to an auxiliary application that processes
373 and rearranges it to be in the usual Strapi format. The data consists in the finished model
374 and also an array filled with all the components created. At this point, the user can download
375 a zipped version of the API model, if they so intend, and easily run it on their personal
376 device.

377 Furthermore, DataGen populates the newly created API with the generated dataset
378 through a number of POST operations. Because of Strapi's lack of methods for importing
379 whole files as data, this cycle of POST requests was the solution found to provide a temporary
380 populated API REST, with all the standard HTTP methods functional.

381 **4 Results**

382 One of the priorities during development was to test DataGen with real cases from early on,
383 in order to not only validate the concept and its operability, but also to observe what kind
384 of restrictions and requests were more frequent in the creation of test datasets, as a means
385 to gather reliable external feedback on useful capabilities that should be implemented.

386 The team made contact with other parties and received requests to create test datasets
387 for real systems, using DataGen, which involved the usage of complicated generation logic
388 with many restrictions. These opportunities helped further DataGen's growth, as new ideas
389 arised from the analysis of the requirements and were brought to life in the application, as
390 well as proved the application's credibility, given that the results obtained were adequate
391 and very positive.

392 In the interest of keeping this paper concise, it will be shown only the most important
393 parts of one of the most complex of these application cases, that of elimination records [7].

5:12 DataGen

394 Elimination records are a structure that must be created and carefully filled out in order
395 to safely eliminate documentation that reaches the end of its administrative conservation
396 deadline. This is an increasingly important tool nowadays, since most public information
397 has shifted to being stored in digital format and the correct method for storing such data is
398 often not followed, which increases the risk of long-term information loss. In order to correct
399 this, the deletion of outdated documentation is just as important as the storage of new one.

400 The generation of these documents implies a complex logic, with many properties that
401 directly relate between themselves according to their values and others whose value must
402 belong to a very rigid group of possibilities. Each record has a legitimation source, whose
403 type can take one of five different string values. According to the record's source type, its
404 funds (public entities) vary from a single entity, in some cases, to an array of several:

```
legitimationSource: {
  type: '{{random("PGD/LC", "TS/LC", "PGD", "RADA", "RADA/CLAV")}}',
  funds(gen) {
    if (["PGD/LC", "TS/LC", "PGD"].includes(this.legitimationSource.type))
      return [gen.pt_entity()]
    else {
      var arr = [], i
      for (i=0; i < gen.integer(1,5); i++) arr.push(gen.pt_entity())
      return arr
    }
  }
}
```

405 Moving on, each record has an array of classes. In case the legitimation source's type is
406 "PGD/LC" or "TS/LC", each class has a code; else, it has either a code, a reference or both.
407 The class code itself can be composed by 3 or 4 levels, given that each level follows its own
408 categorization:

```
classes: [ 'repeat(2,5)': {
  if (["PGD/LC", "TS/LC"].includes(this.legitimationSource.type)) {
    code(gen) {
      var level1 = gen.random(...gen.range(100,950,50))
      var level2 = gen.random(10,20,30,40,50)
      var level3 = gen.integer(1,999,3)
      var level4 = gen.random("01", "02")

      var class = level1 + '.' + level2 + '.' + level3
      if (Math.random() > 0.5) class += '.' + level4
      return class
    }
  }
  else {
    at_least(1) {
      code(gen) { (...) //equal to the function above },
      reference: '{{random(1,2,3,55,56)}}'
    }
  }
}
```

```

    }
  } ]

```

409 There are also year properties that must belong to the last 100 years:

```

yearStart: '{{integer(1921,2021}}}',
yearEnd(gen) {
  var year = gen.integer(1921,2021)
  while (year < this.yearStart) year = gen.integer(1921,2021)
  return year
}

```

410 Finally, there are two related fields, number of aggregations and the corresponding list,
411 where the size of the list must correspond to the number indicated:

```

numberAggregations: '{{integer(1,50}}}',
aggregations: [ 'repeat(this.numberAggregations)': {
  code: '{{pt_entity_abbrev}} - {{integer(1,200}}}',
  title: '{{lorem(3,"words"}}}',
  year: '{{integer(1921,2021}}}',
  if (["PGD/LC","TS/LC"].includes(this.legitimationSource.type)) {
    interventionNature: '{{random("PARTICIPANT","OWNER"}}}'
  }
} ]

```

412 **5 Conclusion**

413 Along the paper it was discussed the development of a multilingual data generator, with
414 built-in REST API integration. The intent behind this project was to create a versatile and
415 powerful tool that would allow for quick prototyping and testing of software applications, a
416 very important and common subject that seems to go surprisingly unnoticed, despite its vast
417 relevance.

418 Be it either small-scale projects of university students or big, complex company software,
419 every application should be thoroughly tested along its development, which requires the
420 leading team to create realistic data in order to populate the system. Even today, this process
421 is very poorly optimized, which often leads either to very time-consuming manual generation
422 or, even worse, to a scarce and inefficient testing of the system, with few records, possibly
423 leading to wrongful conclusions, unnoticed bugs and dangerous bottlenecks.

424 As such, DataGen emerges as a quick and easy to use application that allows the user to
425 swiftly prototype a data model according to their use cases and put their system to practice
426 with a newly-generated dataset, with accurate and realistic values, automating the generation
427 process and facilitating the user's role in it, ultimately enhancing the user's experience and
428 allowing more time and resources to go towards the project itself.

429 DataGen was thoroughly experimented with real-life cases and proved to be capable of
430 creating complex and sizeable datasets for third party applications. The product will very
431 soon be put in a production environment and made available for the general public, after a
432 laborious and successful development phase.

433 **6 Future work**

434 This platform will be open-source and its contents will be uploaded to GitHub. The next
 435 step for the application itself is putting it in a production environment, to be made available
 436 for anyone that may want to use it.

437 As for the grammar, the team intends to develop an user-friendly personalized grammar
 438 checker that analyzes the user's DSL model and, in the presence of errors, communicates
 439 what they are, exactly where they occur and how to fix them, in a simple and clear way, in
 440 order to enhance the user's experience and make the application easier to use.

441 Extensive documentation on all the functionalities provided is also under development,
 442 along with examples on how to use them, in order to guide the user since the application
 443 uses a DSL. Without it, the user may be overwhelmed by the amount of features they must
 444 learn by themselves. This documentation will be made available in the website and may
 445 eventually be downloaded in PDF format, if the user so wishes.

446 **References**

-
- 447 1 D.b. statistical disclosure limitation. page 461–468, 1993.
 - 448 2 General data protection regulation. In *GDPR*, 2018. Accessed: 2021-04-26. URL: <https://gdpr-info.eu/>.
 - 449 3 Artificial intelligence in health care: Benefits and challenges of machine learning in drug
 450 development (staa)-policy briefs & reports-epta network. 2020. Accessed: 2021-04-25. URL:
 451 [https://eptanetwork.org/database/policy-briefs-reports/1898-artificial-intelli](https://eptanetwork.org/database/policy-briefs-reports/1898-artificial-intelligence-in-health-care-benefits-and-challenges-of-machine-learning-in-drug-development-staa)
 452 [gence-in-health-care-benefits-and-challenges-of-machine-learning-in-drug-dev](https://eptanetwork.org/database/policy-briefs-reports/1898-artificial-intelligence-in-health-care-benefits-and-challenges-of-machine-learning-in-drug-development-staa)
 453 [elopment-staa](https://eptanetwork.org/database/policy-briefs-reports/1898-artificial-intelligence-in-health-care-benefits-and-challenges-of-machine-learning-in-drug-development-staa).
 - 454 4 Yahya Al-Hadhrami and Farookh Khadeer Hussain. Real time dataset generation framework
 455 for intrusion detection systems in iot. *Future Generation Computer Systems*, 108:414–423,
 456 2020. URL: <https://www.sciencedirect.com/science/article/pii/S0167739X19322678>,
 457 doi:<https://doi.org/10.1016/j.future.2020.02.051>.
 - 458 5 Anat Reiner Benaim, Ronit Almog, Yuri Gorelik, Irit Hochberg, Laila Nassar, Tanya Mashiach,
 459 Mogher Khamaisi, Yael Lurie, Zaher S Azzam, Johad Houry, Daniel Kurnik, and Rafael
 460 Beyar. Analyzing medical research results based on synthetic data and their relation to real
 461 data results: Systematic comparison from five observational studies. 2015. URL: [https://unece.org/fileadmin/DAM/stats/documents/ece/ces/ge.46/20150/Paper_33_Session](https://unece.org/fileadmin/DAM/stats/documents/ece/ces/ge.46/20150/Paper_33_Session_2_-_Univ._Edinburgh__Nowok_.pdf)
 462 [_2_-_Univ._Edinburgh__Nowok_.pdf](https://unece.org/fileadmin/DAM/stats/documents/ece/ces/ge.46/20150/Paper_33_Session_2_-_Univ._Edinburgh__Nowok_.pdf), doi:10.2196/18910.
 - 463 6 Maurilio Di Cicco, Ciro Potena, Giorgio Grisetti, and Alberto Pretto. Automatic model
 464 based dataset generation for fast and accurate crop and weeds detection. In *2017 IEEE/RSJ*
 465 *International Conference on Intelligent Robots and Systems (IROS)*, pages 5188–5195, 2017.
 466 doi:10.1109/IR0S.2017.8206408.
 - 467 7 Elimination records. <https://clav.dglab.gov.pt/autosEliminacaoInfo/>. Accessed:
 468 2020-05-02.
 - 469 8 Bryan Ford. Parsing Expression Grammars: A Recognition-Based Syntactic Foundation. In
 470 *Proceedings of the 31st ACM SIGPLAN-SIGACT Symposium on Principles of Programming*
 471 *Languages*, 2004. Accessed: 2021-04-20. URL: <https://bford.info/pub/lang/peg.pdf>.
 - 472 9 Georgios Gousios. The ghtorent dataset and tool suite. In *2013 10th Working Conference on*
 473 *Mining Software Repositories (MSR)*, pages 233–236, 2013. doi:10.1109/MSR.2013.6624034.
 - 474 10 Bill Howe, Julia Stoyanovich, Haoyue Ping, Bernease Herman, and Matt Gee. Synthetic data
 475 for social good. 2017.
 - 476 11 JSON Generator. <https://next.json-generator.com/4kaddUyG9/>. Accessed: 2020-05-04.
 - 477 12 Xiangjie Kong, Feng Xia, Zhaolong Ning, Azizur Rahim, Yinqiong Cai, Zhiqiang Gao, and
 478 Jianhua Ma. Mobility dataset generation for vehicular social networks based on floating car
 479
 480

- 481 data. *IEEE Transactions on Vehicular Technology*, 67(5):3874–3886, 2018. doi:10.1109/TVT.
482 2017.2788441.
- 483 13 Menno Mostert, Annelien L Bredenoord, Monique Biesart, and Johannes Delden. Big data in
484 medical research and eu data protection law: Challenges to the consent or anonymise approach.
485 2016. doi:10.1038/ejhg.2015.239.
- 486 14 Beata Nowok. Analyzing medical research results based on synthetic data and their relation
487 to real data results: Systematic comparison from five observational studies. 2020. Accessed:
488 2021-05-03.
- 489 15 PegJS. <https://pegjs.org/>. Accessed: 2021-04-20.
- 490 16 Haoyue Ping, Julia Stoyanovich, and Bill Howe. Datasynthetizer: Privacy-preserving synthetic
491 datasets. In *Proceedings of SSDBM '17*, 2017. doi:10.1145/3085504.3091117.
- 492 17 Darijo Raca, Dylan Leahy, Cormac J. Sreenan, and Jason J. Quinlan. Beyond throughput,
493 the next generation: A 5g dataset with channel and context metrics. In *Proceedings of the*
494 *11th ACM Multimedia Systems Conference, MMSys '20*, page 303–308, New York, NY, USA,
495 2020. Association for Computing Machinery. doi:10.1145/3339825.3394938.
- 496 18 Debbie Rankin, Michaela Black, Raymond Bond, Jonathan Wallace, Maurice Mulvenna, and
497 Gorka Epelde. Reliability of supervised machine learning using synthetic data in health care:
498 Model to preserve privacy for data sharing. 2020. doi:10.2196/18910.
- 499 19 REGULAMENTO NACIONAL DE INTEROPERABILIDADE DIGITAL (RNID). <https://dre.pt/application/file/a/114461891>. Accessed: 2020-04-21.
- 500 20 Design APIs fast, manage content easily. <https://strapi.io/>. Accessed: 2020-04-21.
- 501