

# Bioschemas data harvesting project report

Alasdair Gray<sup>1</sup>, Petros Papadopoulos<sup>1</sup>, Alban Gaignard<sup>2</sup>, Thomas Rosnet<sup>3</sup>, Ivan Mičetić<sup>4</sup>, and Sébastien Moretti<sup>5</sup>

**BioHackathon series:**

[BioHackathon Europe 2021](#)

Barcelona, Spain, 2021

[Bioschemas team](#)

**Submitted:** 28 Apr 2022

**License:**

Authors retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC-BY](#)).

Published by [BioHackrXiv.org](#)

<sup>1</sup> Heriot-Watt University, Edinburgh, UK <sup>2</sup> Nantes Université, CNRS, INSERM, l'institut du thorax, F-44000 Nantes, France <sup>3</sup> Institut Français de Bioinformatique, CNRS UAR 3601, France <sup>4</sup> University of Padua, Padova, Italy <sup>5</sup> Université de Lausanne, Lausanne, Switzerland

## Introduction

The promise of Bioschemas is that it makes consuming data from multiple resources more straightforward. However, this hypothesis has not been tested by conducting a large scale harvest of deployed markup and making this available for others to reuse. Therefore, the goal of this hackathon project is to harvest a collection of Bioschemas markup from a number of different sites listed on the Bioschemas [live deploys page](#) using the Bioschemas Markup Scraper and Extractor ([BMUSE](#)). The harvested data will be made available for others and loaded into a triplestore to allow for further exploration.

## Data Harvesting

Prior to the BioHackathon, we set about harvesting data from as many of the Bioschemas [live deploy sites](#) as possible. At the time of the BioHackathon, there were 70 sites listed, and 137 profile deployments (a site can deploy multiple profiles, e.g. Dataset and DataCatalog). Not all deployments could be harvested since they do not provide sitemaps listing the pages within the site. At the time of the BioHackathon there were 25 sites with sitemaps. Several of these do not list the pages containing data, limiting the amount that could be harvested.

The list of sites to be harvested were gathered in a GitHub [project board](#) so that progress could be tracked. The cards in this board were annotated to state whether the source was known to use a static site deployment (i.e. the markup is embedded in the page source by the server) or dynamic single page application (i.e. the page content is generated client side using Javascript), and also whether they were known to have data content or limited content of Dataset and DataCatalog.

The Bioschemas Markup Scraper and Extractor ([BMUSE](#)) was used for the harvesting of the data. During the harvesting we found a two key issues with BMUSE which arose due to the scale of the data harvest. The first was that errors in the JSON-LD were not correctly identified and logged. The second was a memory limit relating to JSoup which meant that only about 24,104 pages were scraped out of the 50,000 in the sitemap file ([BMUSE #82](#)). Fixes to these issues were applied resulting in BMUSE v0.5.2 being used for most of the harvesting.

The data harvesting workflow consisted of the following steps:

1. Pick one of the sites to be harvested: priority was given to static sites with data content since these could be harvested more quickly and went beyond Dataset/DataCatalog markup.
2. For each sitemap in the sitemap index, harvest the content from the source pages.
3. Merge the individual nquad files for each page in the (sub)sitemap into a single nquad file.
4. Load the merged nquad file into the triplestore.
5. Make the merged nquad file available on the web.
6. Update the project [README](#) with details of what had been harvested.

Where issues were found with the source site, these were feedback to the data provider to allow them to revise their markup. For example, it was found that MassBank were including characters in their string values such as " that need to be escaped to generate valid JSON-LD ([MassBank-Web #316](#)).

In total, six sites were found to be unscrapable. These were

- InterPro: a dynamic site providing a sitemap. However, the sitemap did not conform with the sitemap standard.
- Scholia COVID-19 URL list: a site generated via SPARQL queries over the Wikidata endpoint. Unable to scrape due to timeout being reached before the data being available.
- SwissModel: the provided sitemap did not conform with the sitemap standard.
- WikiPathways: sitemap was empty.
- IPPIDB: a dynamic site with data content corresponding to MolecularEntities. However, the pages exhibit inconsistent rendering when tested in the browser and could not be harvested with BMUSE.
- OrphaNet: a static site with disease markup. The sitemap conforms with the older Google proposal for sitemaps rather than the widely used 0.9 version expected by BMUSE

## Data Analysis

We reused the notebook originally developed at BioHackathon 2020 (Gray et al., 2021) and since evolved for the Intrinsically Disordered Protein Knowledge Graph (IDP-KG) (Gray et al., 2022). We include the HCLS Dataset Description profile statistics queries<sup>1</sup> (Dumontier et al., 2016), read in from an existing [repository](#). We also include [queries](#) developed specifically for the analysis of the Bioschemas harvested data.

To use the [notebook](#) ([MyBinder launcher](#)), you simply need to run all cells and then select the query you would like to execute from the resulting dropdown menu.

We now present the results of the queries obtained during the hackathon, i.e. the data values are as they were on 11 November 2021. Running the notebook in March 2022 obtains different results due to more harvested data having been added.

## HCLS Dataset Statistics

We include here a selection of results from some of the HCLS statistics queries. We focus on those providing the most interesting statistics for the available data. For the full set of queries and results, please run the notebook.

### Number of triples

This is the raw count of the number of triples contained in the triplestore repository.

triples
10,610,743

### Number of named graphs

The result presented here is equivalent to number of pages harvested since BMUSE generates a named graph for each page harvested.

<sup>1</sup>[Dataset Descriptions: HCLS Community Profile §6](#) accessed March 2022

---

graphs

---

413,748

---

### Number of instance per class

There are many different types included in the markup. BMUSE extracts all markup, not just Bioschemas profiles.

The results are ordered by the Class IRI; in the notebook you can edit the query and change the ordering of results.

(57 results)

Class	distinctInstances
<a href="http://rdfs.org/sioc/ns#Item">http://rdfs.org/sioc/ns#Item</a>	57
<a href="http://xmlns.com/foaf/0.1/Document">http://xmlns.com/foaf/0.1/Document</a>	89
<a href="http://xmlns.com/foaf/0.1/Image">http://xmlns.com/foaf/0.1/Image</a>	219
<a href="https://bioschemas.org/Gene">https://bioschemas.org/Gene</a>	238,079
<a href="https://bioschemas.org/Protein">https://bioschemas.org/Protein</a>	1,262
<a href="https://bioschemas.org/Taxon">https://bioschemas.org/Taxon</a>	55,884
<a href="https://schema.org/AboutPage">https://schema.org/AboutPage</a>	1
<a href="https://schema.org/Action">https://schema.org/Action</a>	3
<a href="https://schema.org/Answer">https://schema.org/Answer</a>	8
<a href="https://schema.org/BioChemEntity">https://schema.org/BioChemEntity</a>	49,823
<a href="https://schema.org/BreadcrumbList">https://schema.org/BreadcrumbList</a>	14,037
<a href="https://schema.org/ChemicalSubstance">https://schema.org/ChemicalSubstance</a>	29
<a href="https://schema.org/CollectionPage">https://schema.org/CollectionPage</a>	187
<a href="https://schema.org/CollegeOrUniversity">https://schema.org/CollegeOrUniversity</a>	2
<a href="https://schema.org/ContactPoint">https://schema.org/ContactPoint</a>	148
<a href="https://schema.org/CreativeWork">https://schema.org/CreativeWork</a>	14,299
<a href="https://schema.org/DataCatalog">https://schema.org/DataCatalog</a>	7,439
<a href="https://schema.org/DataDownload">https://schema.org/DataDownload</a>	1,497
<a href="https://schema.org/Dataset">https://schema.org/Dataset</a>	201,302
<a href="https://schema.org/DefinedTerm">https://schema.org/DefinedTerm</a>	4,261
<a href="https://schema.org/DefinedTermSet">https://schema.org/DefinedTermSet</a>	4,112
<a href="https://schema.org/DigitalDocument">https://schema.org/DigitalDocument</a>	1
<a href="https://schema.org/EducationalOrganization">https://schema.org/EducationalOrganization</a>	3
<a href="https://schema.org/Event">https://schema.org/Event</a>	12,818
<a href="https://schema.org/FAQPage">https://schema.org/FAQPage</a>	1
<a href="https://schema.org/Gene">https://schema.org/Gene</a>	39
<a href="https://schema.org/GeoShape">https://schema.org/GeoShape</a>	19,398
<a href="https://schema.org/GovernmentOrganization">https://schema.org/GovernmentOrganization</a>	1
<a href="https://schema.org/ItemList">https://schema.org/ItemList</a>	187
<a href="https://schema.org/ListItem">https://schema.org/ListItem</a>	28,137
<a href="https://schema.org/MolecularEntity">https://schema.org/MolecularEntity</a>	199,350
<a href="https://schema.org/NGO">https://schema.org/NGO</a>	11,717
<a href="https://schema.org/Offer">https://schema.org/Offer</a>	5
<a href="https://schema.org/Organization">https://schema.org/Organization</a>	206,715
<a href="https://schema.org/PeopleAudience">https://schema.org/PeopleAudience</a>	2,475
<a href="https://schema.org/Person">https://schema.org/Person</a>	326,935
<a href="https://schema.org/Place">https://schema.org/Place</a>	19,438
<a href="https://schema.org/PostalAddress">https://schema.org/PostalAddress</a>	307,406
<a href="https://schema.org/PropertyValue">https://schema.org/PropertyValue</a>	144,002
<a href="https://schema.org/Protein">https://schema.org/Protein</a>	4,462

Class	distinctInstances
<a href="https://schema.org/QAPage">https://schema.org/QAPage</a>	1
<a href="https://schema.org/Question">https://schema.org/Question</a>	8
<a href="https://schema.org/ScholarlyArticle">https://schema.org/ScholarlyArticle</a>	9,350
<a href="https://schema.org/SearchAction">https://schema.org/SearchAction</a>	5
<a href="https://schema.org/SequenceAnnotation">https://schema.org/SequenceAnnotation</a>	15,786
<a href="https://schema.org/SequenceRange">https://schema.org/SequenceRange</a>	15,786
<a href="https://schema.org/SoftwareApplication">https://schema.org/SoftwareApplication</a>	4
<a href="https://schema.org/SoftwareSourceCode">https://schema.org/SoftwareSourceCode</a>	4
<a href="https://schema.org/Study">https://schema.org/Study</a>	4,328
<a href="https://schema.org/Thing">https://schema.org/Thing</a>	27,872
<a href="https://schema.org/URL">https://schema.org/URL</a>	1
<a href="https://schema.org/WebApplication">https://schema.org/WebApplication</a>	3
<a href="https://schema.org/WebPage">https://schema.org/WebPage</a>	55,114
<a href="https://schema.org/WebSite">https://schema.org/WebSite</a>	5
<a href="https://schema.org/contact">https://schema.org/contact</a>	40
<a href="https://schema.org/hostInstitution">https://schema.org/hostInstitution</a>	40
<a href="https://schema.org/url">https://schema.org/url</a>	10,360

## Bioschemas Queries

The following queries focus on features of interest to the Bioschemas community.

### Instances per Bioschemas Class

Note that due to the data content we need to include some properties with both a Bioschemas namespace and a Schema.org namespace.

The results are ordered by the count of the number of instances; in the notebook you can edit the query and change the ordering of results.

(18 results)

Class	instances
<a href="https://schema.org/Person">https://schema.org/Person</a>	326,935
<a href="https://bioschemas.org/Gene">https://bioschemas.org/Gene</a>	238,079
<a href="https://schema.org/Organization">https://schema.org/Organization</a>	206,715
<a href="https://schema.org/Dataset">https://schema.org/Dataset</a>	201,302
<a href="https://schema.org/MolecularEntity">https://schema.org/MolecularEntity</a>	199,350
<a href="https://bioschemas.org/Taxon">https://bioschemas.org/Taxon</a>	55,884
<a href="https://schema.org/BioChemEntity">https://schema.org/BioChemEntity</a>	49,823
<a href="https://schema.org/SequenceAnnotation">https://schema.org/SequenceAnnotation</a>	15,786
<a href="https://schema.org/SequenceRange">https://schema.org/SequenceRange</a>	15,786
<a href="https://schema.org/Event">https://schema.org/Event</a>	12,818
<a href="https://schema.org/ScholarlyArticle">https://schema.org/ScholarlyArticle</a>	9,350
<a href="https://schema.org/DataCatalog">https://schema.org/DataCatalog</a>	7,439
<a href="https://schema.org/Protein">https://schema.org/Protein</a>	4,462
<a href="https://schema.org/Study">https://schema.org/Study</a>	4,328
<a href="https://bioschemas.org/Protein">https://bioschemas.org/Protein</a>	1,262
<a href="https://schema.org/Gene">https://schema.org/Gene</a>	39
<a href="https://schema.org/ChemicalSubstance">https://schema.org/ChemicalSubstance</a>	29
<a href="https://schema.org/SoftwareApplication">https://schema.org/SoftwareApplication</a>	4

### Number of Domains

This result informs us how many web domains were harvested. This is approximately equal to the number of datasets, but some sites may host more than one dataset so not necessarily an exact correspondence.

count
25

### Number of Pages per Domain

We now report the number of pages that have been harvested from each domain. Note that we do not understand the empty domain as all markup was extracted from a web domain.

(25 results)

domain	count
massbank.eu	76,253
scholia.toolforge.org	74,319
www.gbif.org	68,167
test.intermine.org	49,959
bgee.org	49,022
www.metanetx.org	49,012
tess.elixir-europe.org	13,939
ega-archive.org	11,833
fairsharing.org	6,351
prosite.expasy.org	5,858
ippidb.pasteur.fr	2,433
mobidb.org	2,082
disprot.org	2,043
pcddb.cryst.bbk.ac.uk	1,402
www.ebi.ac.uk	672
proteinensemble.org	187
www.france-bioinformatique.fr	86
pairedomicsdata.bioinformatics.nl	78
www.covid19dataportal.org	19
	12
www.alliancegenome.org	11
biopragnomics.github.io	3
nanocommons.github.io	3
bridgedb.github.io	2
www.uniprot.org	2

### Count of Types per Domain

We now report the number of instances of each type on each domain. What is interesting here is the fact that Bgee has many proteins listed on their pages.

The results are ordered by the count of the number of instances; in the notebook you can edit the query and change the ordering of results.

(146 results)

domain	type	count
www.gbif.org	<a href="https://schema.org/PostalAddress">https://schema.org/PostalAddress</a>	297,090
www.gbif.org	<a href="https://schema.org/Person">https://schema.org/Person</a>	291,260
bgee.org	<a href="https://bioschemas.org/Gene">https://bioschemas.org/Gene</a>	263,793
www.gbif.org	<a href="https://schema.org/Organization">https://schema.org/Organization</a>	186,688
www.gbif.org	<a href="https://schema.org/PropertyValue">https://schema.org/PropertyValue</a>	126,268
massbank.eu	<a href="https://schema.org/Dataset">https://schema.org/Dataset</a>	76,249
massbank.eu	<a href="https://schema.org/MolecularEntity">https://schema.org/MolecularEntity</a>	76,249
scholia.toolforge.org	<a href="https://schema.org/CreativeWork">https://schema.org/CreativeWork</a>	74,310
scholia.toolforge.org	<a href="https://schema.org/MolecularEntity">https://schema.org/MolecularEntity</a>	74,310
www.gbif.org	<a href="https://schema.org/Dataset">https://schema.org/Dataset</a>	63,134
test.intermine.org	<a href="https://schema.org/Dataset">https://schema.org/Dataset</a>	49,959
test.intermine.org	<a href="https://schema.org/BioChemEntity">https://schema.org/BioChemEntity</a>	49,823
bgee.org	<a href="https://bioschemas.org/Taxon">https://bioschemas.org/Taxon</a>	49,059
bgee.org	<a href="https://schema.org/WebPage">https://schema.org/WebPage</a>	49,009
www.metanetx.org	<a href="https://schema.org/CreativeWork">https://schema.org/CreativeWork</a>	49,002
www.metanetx.org	<a href="https://schema.org/MolecularEntity">https://schema.org/MolecularEntity</a>	49,001
prosite.expasy.org	<a href="https://schema.org/Person">https://schema.org/Person</a>	31,364
tess.elixir-europe.org	<a href="https://schema.org/ListItem">https://schema.org/ListItem</a>	27,872
tess.elixir-europe.org	<a href="https://schema.org/Thing">https://schema.org/Thing</a>	27,872
www.gbif.org	<a href="https://schema.org/GeoShape">https://schema.org/GeoShape</a>	19,398
www.gbif.org	<a href="https://schema.org/Place">https://schema.org/Place</a>	19,398
tess.elixir-europe.org	<a href="https://schema.org/BreadcrumbList">https://schema.org/BreadcrumbList</a>	13,938
tess.elixir-europe.org	<a href="https://schema.org/Event">https://schema.org/Event</a>	12,778
prosite.expasy.org	<a href="https://schema.org/Organization">https://schema.org/Organization</a>	11,715
prosite.expasy.org	<a href="https://schema.org/NGO">https://schema.org/NGO</a>	11,714
disprot.org	<a href="https://schema.org/PropertyValue">https://schema.org/PropertyValue</a>	11,046
disprot.org	<a href="https://schema.org/SequenceAnnotation">https://schema.org/SequenceAnnotation</a>	11,046
disprot.org	<a href="https://schema.org/SequenceRange">https://schema.org/SequenceRange</a>	11,046
prosite.expasy.org	<a href="https://schema.org/url">https://schema.org/url</a>	10,360
tess.elixir-europe.org	<a href="https://schema.org/PostalAddress">https://schema.org/PostalAddress</a>	10,316
ega-archive.org	<a href="https://schema.org/DataCatalog">https://schema.org/DataCatalog</a>	7,431
ega-archive.org	<a href="https://schema.org/Dataset">https://schema.org/Dataset</a>	7,431
tess.elixir-europe.org	<a href="https://schema.org/Organization">https://schema.org/Organization</a>	7,110
prosite.expasy.org	<a href="https://bioschemas.org/Taxon">https://bioschemas.org/Taxon</a>	6,796
prosite.expasy.org	<a href="https://schema.org/ScholarlyArticle">https://schema.org/ScholarlyArticle</a>	6,681
disprot.org	<a href="https://schema.org/DefinedTerm">https://schema.org/DefinedTerm</a>	6,599
fairsharing.org	<a href="https://schema.org/Dataset">https://schema.org/Dataset</a>	6,328
prosite.expasy.org	<a href="https://schema.org/WebPage">https://schema.org/WebPage</a>	6,093
prosite.expasy.org	<a href="https://schema.org/CreativeWork">https://schema.org/CreativeWork</a>	5,857
fairsharing.org	<a href="https://schema.org/CreativeWork">https://schema.org/CreativeWork</a>	5,542
mobidb.org	<a href="https://schema.org/PropertyValue">https://schema.org/PropertyValue</a>	4,486
mobidb.org	<a href="https://schema.org/SequenceAnnotation">https://schema.org/SequenceAnnotation</a>	4,486
mobidb.org	<a href="https://schema.org/SequenceRange">https://schema.org/SequenceRange</a>	4,486
ega-archive.org	<a href="https://schema.org/Study">https://schema.org/Study</a>	4,328
disprot.org	<a href="https://schema.org/DefinedTermSet">https://schema.org/DefinedTermSet</a>	4,076
mobidb.org	<a href="https://schema.org/DefinedTerm">https://schema.org/DefinedTerm</a>	3,400
mobidb.org	<a href="https://schema.org/DefinedTermSet">https://schema.org/DefinedTermSet</a>	3,400
tess.elixir-europe.org	<a href="https://schema.org/Person">https://schema.org/Person</a>	3,298
tess.elixir-europe.org	<a href="https://schema.org/CreativeWork">https://schema.org/CreativeWork</a>	2,876
disprot.org	<a href="https://schema.org/ScholarlyArticle">https://schema.org/ScholarlyArticle</a>	2,857
tess.elixir-europe.org	<a href="https://schema.org/PeopleAudience">https://schema.org/PeopleAudience</a>	2,475
proteinensemble.org	<a href="https://schema.org/PropertyValue">https://schema.org/PropertyValue</a>	2,202
mobidb.org	<a href="https://schema.org/CreativeWork">https://schema.org/CreativeWork</a>	2,073

domain	type	count
mobidb.org	<a href="https://schema.org/Protein">https://schema.org/Protein</a>	2,073
disprot.org	<a href="https://schema.org/CreativeWork">https://schema.org/CreativeWork</a>	2,038
disprot.org	<a href="https://schema.org/Protein">https://schema.org/Protein</a>	2,038
proteinensemble.org	<a href="https://schema.org/DefinedTerm">https://schema.org/DefinedTerm</a>	1,626
pcddb.cryst.bbk.ac.uk	<a href="https://schema.org/Organization">https://schema.org/Organization</a>	1,402
pcddb.cryst.bbk.ac.uk	<a href="https://schema.org/DataDownload">https://schema.org/DataDownload</a>	1,394
prosite.expasy.org	<a href="https://bioschemas.org/Protein">https://bioschemas.org/Protein</a>	1,376
pcddb.cryst.bbk.ac.uk	<a href="https://schema.org/Dataset">https://schema.org/Dataset</a>	697
pcddb.cryst.bbk.ac.uk	<a href="https://schema.org/Person">https://schema.org/Person</a>	697
biopragmatics.github.io	<a href="https://schema.org/Dataset">https://schema.org/Dataset</a>	287
www.france-bioinformatique.fr	<a href="https://schema.org/ListItem">https://schema.org/ListItem</a>	265
proteinensemble.org	<a href="https://schema.org/Protein">https://schema.org/Protein</a>	254
proteinensemble.org	<a href="https://schema.org/SequenceAnnotation">https://schema.org/SequenceAnnotation</a>	254
proteinensemble.org	<a href="https://schema.org/SequenceRange">https://schema.org/SequenceRange</a>	254
pairedomicsdata.bioinformatics.nl	<a href="https://schema.org/Person">https://schema.org/Person</a>	222
www.ebi.ac.uk	<a href="http://xmlns.com/foaf/0.1/Image">http://xmlns.com/foaf/0.1/Image</a>	222
proteinensemble.org	<a href="https://schema.org/CollectionPage">https://schema.org/CollectionPage</a>	187
proteinensemble.org	<a href="https://schema.org/CreativeWork">https://schema.org/CreativeWork</a>	187
proteinensemble.org	<a href="https://schema.org/DefinedTermSet">https://schema.org/DefinedTermSet</a>	187
proteinensemble.org	<a href="https://schema.org/ItemList">https://schema.org/ItemList</a>	187
proteinensemble.org	<a href="https://schema.org/ScholarlyArticle">https://schema.org/ScholarlyArticle</a>	181
pairedomicsdata.bioinformatics.nl	<a href="https://schema.org/ContactPoint">https://schema.org/ContactPoint</a>	148
www.covid19dataportal.org	<a href="https://schema.org/Organization">https://schema.org/Organization</a>	148
www.covid19dataportal.org	<a href="https://schema.org/Dataset">https://schema.org/Dataset</a>	110
www.france-bioinformatique.fr	<a href="https://schema.org/BreadcrumbList">https://schema.org/BreadcrumbList</a>	99
	<a href="https://schema.org/Dataset">https://schema.org/Dataset</a>	97
test.intermine.org	<a href="https://schema.org/Protein">https://schema.org/Protein</a>	97
www.ebi.ac.uk	<a href="http://xmlns.com/foaf/0.1/Document">http://xmlns.com/foaf/0.1/Document</a>	91
	<a href="https://schema.org/Person">https://schema.org/Person</a>	90
bgee.org	<a href="https://schema.org/Dataset">https://schema.org/Dataset</a>	87
pairedomicsdata.bioinformatics.nl	<a href="https://schema.org/CreativeWork">https://schema.org/CreativeWork</a>	74
pairedomicsdata.bioinformatics.nl	<a href="https://schema.org/DataCatalog">https://schema.org/DataCatalog</a>	74
pairedomicsdata.bioinformatics.nl	<a href="https://schema.org/DataDownload">https://schema.org/DataDownload</a>	74
pairedomicsdata.bioinformatics.nl	<a href="https://schema.org/Dataset">https://schema.org/Dataset</a>	74
pairedomicsdata.bioinformatics.nl	<a href="https://schema.org/Organization">https://schema.org/Organization</a>	74
www.ebi.ac.uk	<a href="http://rdfs.org/sioc/ns#Item">http://rdfs.org/sioc/ns#Item</a>	59
www.france-bioinformatique.fr	<a href="https://schema.org/Event">https://schema.org/Event</a>	40
www.france-bioinformatique.fr	<a href="https://schema.org/Place">https://schema.org/Place</a>	40
www.france-bioinformatique.fr	<a href="https://schema.org/contact">https://schema.org/contact</a>	40
www.france-bioinformatique.fr	<a href="https://schema.org/hostInstitution">https://schema.org/hostInstitution</a>	40
test.intermine.org	<a href="https://schema.org/Gene">https://schema.org/Gene</a>	39
	<a href="https://bioschemas.org/Taxon">https://bioschemas.org/Taxon</a>	29
nanocommons.github.io	<a href="https://schema.org/ChemicalSubstance">https://schema.org/ChemicalSubstance</a>	29
	<a href="https://schema.org/Organization">https://schema.org/Organization</a>	27
bridgedb.github.io	<a href="https://schema.org/DataDownload">https://schema.org/DataDownload</a>	23
bridgedb.github.io	<a href="https://schema.org/Dataset">https://schema.org/Dataset</a>	23
www.covid19dataportal.org	<a href="https://schema.org/DataCatalog">https://schema.org/DataCatalog</a>	19
www.uniprot.org	<a href="https://schema.org/Organization">https://schema.org/Organization</a>	14
	<a href="https://schema.org/CreativeWork">https://schema.org/CreativeWork</a>	9
nanocommons.github.io	<a href="https://schema.org/Organization">https://schema.org/Organization</a>	9
bgee.org	<a href="https://schema.org/Answer">https://schema.org/Answer</a>	8
bgee.org	<a href="https://schema.org/CreativeWork">https://schema.org/CreativeWork</a>	8
bgee.org	<a href="https://schema.org/Question">https://schema.org/Question</a>	8



domain	type	count
	<a href="https://schema.org/WebPage">https://schema.org/WebPage</a>	7
nanocommons.github.io	<a href="https://schema.org/Dataset">https://schema.org/Dataset</a>	7
nanocommons.github.io	<a href="https://schema.org/DataDownload">https://schema.org/DataDownload</a>	6
	<a href="https://schema.org/DataCatalog">https://schema.org/DataCatalog</a>	5
bgee.org	<a href="https://schema.org/Offer">https://schema.org/Offer</a>	5
	<a href="https://schema.org/SearchAction">https://schema.org/SearchAction</a>	4
bgee.org	<a href="https://schema.org/SoftwareSourceCode">https://schema.org/SoftwareSourceCode</a>	4
nanocommons.github.io	<a href="https://schema.org/CreativeWork">https://schema.org/CreativeWork</a>	4
	<a href="https://schema.org/EducationalOrganization">https://schema.org/EducationalOrganization</a>	3
	<a href="https://schema.org/NGO">https://schema.org/NGO</a>	3
	<a href="https://schema.org/ScholarlyArticle">https://schema.org/ScholarlyArticle</a>	3
bgee.org	<a href="https://schema.org/WebApplication">https://schema.org/WebApplication</a>	3
biopragmatics.github.io	<a href="https://schema.org/Person">https://schema.org/Person</a>	3
prosite.expasy.org	<a href="https://schema.org/Action">https://schema.org/Action</a>	3
	<a href="https://schema.org/CollegeOrUniversity">https://schema.org/CollegeOrUniversity</a>	2
	<a href="https://schema.org/WebSite">https://schema.org/WebSite</a>	2
bgee.org	<a href="https://schema.org/SoftwareApplication">https://schema.org/SoftwareApplication</a>	2
biopragmatics.github.io	<a href="https://schema.org/WebPage">https://schema.org/WebPage</a>	2
bridgedb.github.io	<a href="https://schema.org/CreativeWork">https://schema.org/CreativeWork</a>	2
www.uniprot.org	<a href="https://schema.org/GovernmentOrganization">https://schema.org/GovernmentOrganization</a>	2
www.uniprot.org	<a href="https://schema.org/NGO">https://schema.org/NGO</a>	2
www.uniprot.org	<a href="https://schema.org/WebPage">https://schema.org/WebPage</a>	2
	<a href="https://schema.org/GovernmentOrganization">https://schema.org/GovernmentOrganization</a>	1
bgee.org	<a href="https://schema.org/AboutPage">https://schema.org/AboutPage</a>	1
bgee.org	<a href="https://schema.org/FAQPage">https://schema.org/FAQPage</a>	1
biopragmatics.github.io	<a href="https://schema.org/WebSite">https://schema.org/WebSite</a>	1
bridgedb.github.io	<a href="https://schema.org/SoftwareApplication">https://schema.org/SoftwareApplication</a>	1
bridgedb.github.io	<a href="https://schema.org/WebPage">https://schema.org/WebPage</a>	1
bridgedb.github.io	<a href="https://schema.org/WebSite">https://schema.org/WebSite</a>	1
massbank.eu	<a href="https://schema.org/DataCatalog">https://schema.org/DataCatalog</a>	1
massbank.eu	<a href="https://schema.org/ScholarlyArticle">https://schema.org/ScholarlyArticle</a>	1
nanocommons.github.io	<a href="https://schema.org/DataCatalog">https://schema.org/DataCatalog</a>	1
nanocommons.github.io	<a href="https://schema.org/Person">https://schema.org/Person</a>	1
nanocommons.github.io	<a href="https://schema.org/URL">https://schema.org/URL</a>	1
nanocommons.github.io	<a href="https://schema.org/WebSite">https://schema.org/WebSite</a>	1
prosite.expasy.org	<a href="https://schema.org/DigitalDocument">https://schema.org/DigitalDocument</a>	1
prosite.expasy.org	<a href="https://schema.org/SearchAction">https://schema.org/SearchAction</a>	1
www.metanetx.org	<a href="https://schema.org/SoftwareApplication">https://schema.org/SoftwareApplication</a>	1
www.uniprot.org	<a href="https://schema.org/Dataset">https://schema.org/Dataset</a>	1
www.uniprot.org	<a href="https://schema.org/QAPage">https://schema.org/QAPage</a>	1

## Connectivity of the Data

We were interested to gain some insight as to how connected the data was both internally, and how many points where it would link up with other knowledge graphs. The queries in this section focus on the connectedness of the data.

We first investigated the number of nodes that only contained incoming edges. We report the total number of object nodes there are (excluding literals), and the number of edge IRIs, i.e. those that only have incoming properties. Only 4.65% of the nodes only contain incoming edges.



Object IRIs	Edge IRIs
2,057,094	95,610

We then investigated the number of outgoing links per class. We report here the top 20 results.

s	class	nb_out_edges
<a href="https://bioschemas.org/crawl/v1/bgee/?page=gene&amp;gene_id=ENSBTAG00000027937/20211110/90020/bgee.org/?page=gene&amp;gene_id=ENSBTAG00000027937/1779564251">https://bioschemas.org/crawl/v1/bgee/?page=gene&amp;gene_id=ENSBTAG00000027937/20211110/90020/bgee.org/?page=gene&amp;gene_id=ENSBTAG00000027937/1779564251</a>	<a href="https://bioschemas.org/Gene">https://bioschemas.org/Gene</a>	856
<a href="https://www.metanetx.org/chem_info/MNXM1944">https://www.metanetx.org/chem_info/MNXM1944</a>	<a href="https://schema.org/MolecularEntity">https://schema.org/MolecularEntity</a>	654
<a href="https://doi.org/10.15468/hb9rjv">https://doi.org/10.15468/hb9rjv</a>	<a href="https://schema.org/Dataset">https://schema.org/Dataset</a>	594
<a href="https://bioschemas.org/crawl/v1/bgee/?page=gene&amp;gene_id=ENSBTAG00000043564/20211110/94715/bgee.org/?page=gene&amp;gene_id=ENSBTAG00000043564/1772156424">https://bioschemas.org/crawl/v1/bgee/?page=gene&amp;gene_id=ENSBTAG00000043564/20211110/94715/bgee.org/?page=gene&amp;gene_id=ENSBTAG00000043564/1772156424</a>	<a href="https://bioschemas.org/Gene">https://bioschemas.org/Gene</a>	519
<a href="https://bioschemas.org/crawl/v1/bgee/?page=gene&amp;gene_id=ENSBTAG00000043584/20211110/94734/bgee.org/?page=gene&amp;gene_id=ENSBTAG00000043584/2022662406">https://bioschemas.org/crawl/v1/bgee/?page=gene&amp;gene_id=ENSBTAG00000043584/20211110/94734/bgee.org/?page=gene&amp;gene_id=ENSBTAG00000043584/2022662406</a>	<a href="https://bioschemas.org/Gene">https://bioschemas.org/Gene</a>	474
<a href="https://doi.org/10.15468/m5vrza">https://doi.org/10.15468/m5vrza</a>	<a href="https://schema.org/Dataset">https://schema.org/Dataset</a>	406
<a href="http://www.ebi.ac.uk/pdbe/about/past-events">http://www.ebi.ac.uk/pdbe/about/past-events</a>	<a href="http://rdfs.org/sioc/ns#Item">http://rdfs.org/sioc/ns#Item</a>	346
<a href="http://www.ebi.ac.uk/pdbe/about/past-events">http://www.ebi.ac.uk/pdbe/about/past-events</a>	<a href="http://xmlns.com/foaf/0.1/Document">http://xmlns.com/foaf/0.1/Document</a>	346
<a href="https://doi.org/10.15468/vmf5ye">https://doi.org/10.15468/vmf5ye</a>	<a href="https://schema.org/Dataset">https://schema.org/Dataset</a>	296
<a href="https://bioschemas.org/crawl/v1/bgee/?page=gene&amp;gene_id=ENSBTAG00000043559/20211110/94710/bgee.org/?page=gene&amp;gene_id=ENSBTAG00000043559/1377128066">https://bioschemas.org/crawl/v1/bgee/?page=gene&amp;gene_id=ENSBTAG00000043559/20211110/94710/bgee.org/?page=gene&amp;gene_id=ENSBTAG00000043559/1377128066</a>	<a href="https://bioschemas.org/Gene">https://bioschemas.org/Gene</a>	292
<a href="https://doi.org/10.5281/zenodo.291971">https://doi.org/10.5281/zenodo.291971</a>	<a href="https://schema.org/Dataset">https://schema.org/Dataset</a>	289
<a href="https://bioschemas.org/crawl/v1/bgee/?page=gene&amp;gene_id=ENSBTAG00000043546/20211110/94697/bgee.org/?page=gene&amp;gene_id=ENSBTAG00000043546/1804549344">https://bioschemas.org/crawl/v1/bgee/?page=gene&amp;gene_id=ENSBTAG00000043546/20211110/94697/bgee.org/?page=gene&amp;gene_id=ENSBTAG00000043546/1804549344</a>	<a href="https://bioschemas.org/Gene">https://bioschemas.org/Gene</a>	284
<a href="https://doi.org/10.15472/hy9nif">https://doi.org/10.15472/hy9nif</a>	<a href="https://schema.org/Dataset">https://schema.org/Dataset</a>	282
<a href="https://bioschemas.org/crawl/v1/bgee/?page=gene&amp;gene_id=ENSBTAG00000043550/20211110/94701/bgee.org/?page=gene&amp;gene_id=ENSBTAG00000043550/1476242225">https://bioschemas.org/crawl/v1/bgee/?page=gene&amp;gene_id=ENSBTAG00000043550/20211110/94701/bgee.org/?page=gene&amp;gene_id=ENSBTAG00000043550/1476242225</a>	<a href="https://bioschemas.org/Gene">https://bioschemas.org/Gene</a>	269
<a href="https://www.metanetx.org/chem_info/MNXM383">https://www.metanetx.org/chem_info/MNXM383</a>	<a href="https://schema.org/MolecularEntity">https://schema.org/MolecularEntity</a>	264

<a href="https://bioschemas.org/crawl/v1/bgee/?page=gene&amp;gene_id=ENSBTAG000000043577/20211110/94727/bgee.org/?page=gene&amp;gene_id=ENSBTAG000000043577/1610681495">https://bioschemas.org/crawl/v1/bgee/?page=gene&amp;gene_id=ENSBTAG000000043577/20211110/94727/bgee.org/?page=gene&amp;gene_id=ENSBTAG000000043577/1610681495</a>	<a href="https://bioschemas.org/Gene">https://bioschemas.org/Gene</a>	261
<a href="https://bioschemas.org/crawl/v1/bgee/?page=gene&amp;gene_id=ENSBTAG000000043556/20211110/94707/bgee.org/?page=gene&amp;gene_id=ENSBTAG000000043556/1277162978">https://bioschemas.org/crawl/v1/bgee/?page=gene&amp;gene_id=ENSBTAG000000043556/20211110/94707/bgee.org/?page=gene&amp;gene_id=ENSBTAG000000043556/1277162978</a>	<a href="https://bioschemas.org/Gene">https://bioschemas.org/Gene</a>	240
<a href="https://bioschemas.org/crawl/v1/bgee/?page=gene&amp;gene_id=ENSBTAG000000043568/20211110/94719/bgee.org/?page=gene&amp;gene_id=ENSBTAG000000043568/2065005542">https://bioschemas.org/crawl/v1/bgee/?page=gene&amp;gene_id=ENSBTAG000000043568/20211110/94719/bgee.org/?page=gene&amp;gene_id=ENSBTAG000000043568/2065005542</a>	<a href="https://bioschemas.org/Gene">https://bioschemas.org/Gene</a>	235
<a href="https://bioschemas.org/crawl/v1/bgee/?page=gene&amp;gene_id=ENSBTAG000000043560/20211110/94711/bgee.org/?page=gene&amp;gene_id=ENSBTAG000000043560/1818154049">https://bioschemas.org/crawl/v1/bgee/?page=gene&amp;gene_id=ENSBTAG000000043560/20211110/94711/bgee.org/?page=gene&amp;gene_id=ENSBTAG000000043560/1818154049</a>	<a href="https://bioschemas.org/Gene">https://bioschemas.org/Gene</a>	229

## Discussion

Through this Hackathon project we have demonstrated that it is possible to harvest Bioschemas markup from a number of sites and load them into a triplestore. However, this has revealed a number of challenges that need to be resolved.

Harvesting content from a whole site is very time consuming, particularly for dynamic sites. Harvesting requires visiting each page with markup in turn and extracting the markup. In the case of dynamic sites the content needs to be rendered before it can be extracted. Most of the sites that could be completed harvested did not contain data content beyond Dataset and DataCatalog.

The quality of deployed markup is very problematic. As reported above, a number of sites could not be harvested due to issues with their sitemap. Of those that could be harvested, a number of pages were not harvested due to issues with the markup contained within them, e.g. inclusion of non-escaped characters within strings was a common error. Even the markup that could be extracted contained errors, e.g. the use of different namespaces for the declaration of the Gene type as identified in the Instances per Bioschemas Class query. This highlights the need for a Bioschemas validator capable of both syntactic and semantic checking (see the data validation outputs of [Project 29](#)).

## Future work

The next steps for this work would be to improve the robustness of the data harvesting pipeline, including automating the manual steps of iterating over an index of sitemap files; merging individual files for each harvested page into a single file per subsitemap; and loading the harvested data into the triplestore.

The use of data dumps for sites should be considered to eliminate the need for data harvesting, which is a costly process for both data producers – who have to have sufficient bandwidth and compute resources to serve the content – and data consumers – who have to have sufficient compute resources to retrieve, render, and extract the content from each page.

## Jupyter notebooks, GitHub repositories and data repositories

- GitHub repository: <https://github.com/BioSchemas/bioschemas-data-harvesting>
- Jupyter Notebook: <https://github.com/BioSchemas/bioschemas-data-harvesting/blob/main/AnalysisQueries.ipynb>
  - MyBinder launch: <https://mybinder.org/v2/gh/BioSchemas/bioschemas-data-harvesting/HEAD?labpath=AnalysisQueries.ipynb>
- SPARQL Endpoint: <https://swel.macs.hw.ac.uk/data/repositories/bioschemas>
  - Snorql Extended Interface: <https://swel.macs.hw.ac.uk/bioschemas/>
- Data download director: <https://swel.macs.hw.ac.uk/bioschemas-data/>

## Acknowledgements

This work was done during the BioHackathon Europe 2021 organised by ELIXIR and run in November 2021. We wish to thank the organizers and supporters of the Biohackathon Europe 2021 for offering the venue for improving Bioschemas community efforts.

## References

- Dumontier, M., Gray, A. J., Marshall, M. S., Alexiev, V., Ansell, P., Bader, G., Baran, J., Bolleman, J. T., Callahan, A., Cruz-Toledo, J. others. (2016). The health care and life sciences community profile for dataset descriptions. *PeerJ*, 4, e2331. <https://doi.org/10.7717/peerj.2331>
- Gray, A. J. G., Papadopoulos, P., Asif, I., Mičetić, I., & Hatos, A. (2022, January 11). Creating and exploiting the intrinsically disordered protein knowledge graph (IDP-KG). *13th International Semantic Web Applications and Tools for Health Care and Life Sciences Conference (Swat4hcls 2022)*. <http://www.swat4hcls.org/>
- Gray, A. J. G., Papadopoulos, P., Mičetić, I., & Hatos, A. (2021). *Exploiting bioschemas markup to populate IDPcentral*. BioHackrXiv. <https://doi.org/10.37044/osf.io/v3jct>