

1 Ozymandias: A biodiversity knowledge graph

2
3 Roderic D. M. Page

4 <https://orcid.org/0000-0002-7101-9767>

5 Institute of Biodiversity, Animal Health and Comparative Medicine, College of Medical, Vet-
6 erinary and Life Sciences, Graham Kerr Building, University of Glasgow, Glasgow, UK

7 Email address: Roderic.page@glasgow.ac.uk

8 Abstract

9 Enormous quantities of biodiversity data are being made available online, but much of this
10 data remains isolated in their own silos. One approach to breaking these silos is to map local,
11 often database-specific identifiers to shared global identifiers. This mapping can then be used
12 to construct a knowledge graph, where entities such as taxa, publications, people, places,
13 specimens, sequences, and institutions are all part of a single, shared knowledge space. Moti-
14 vated by the 2018 GBIF Ebbe Nielsen Challenge I explore the feasibility of constructing a
15 “biodiversity knowledge graph” for the Australian fauna. These steps involved in constructing
16 the graph are described, and examples its application are discussed. A web interface to the
17 knowledge graph (called “Ozymandias”) is available at [https://ozymandias-](https://ozymandias-demo.herokuapp.com)
18 [demo.herokuapp.com](https://ozymandias-demo.herokuapp.com).

19
20 Keywords: knowledge graph; biodiversity informatics; linked data; identifiers;

21 Introduction

22 “Linnaeus would have been a ‘techie’” - (Godfray, 2007)

23
24 The recent announcement that the Global Biodiversity Information Facility (GBIF) has
25 reached the milestone of one billion occurrence records reflects the considerable success the
26 biodiversity community has had in mobilising data. Much of this success comes from stan-
27 dardising on a simple column-based data format (Darwin Core) (Wieczorek et al., 2012) and
28 indexing that data using three fields: taxonomic name, geographic location, and date (i.e.,
29 “what”, “where”, and “when”). By flattening the data into a single table, Darwin Core makes
30 data easy to enter and view, but at the cost of potentially obscuring relationships between enti-
31 ties, relationships that may be better represented using a network. In this paper I explore the
32 representation of biodiversity data using a network or “knowledge graph”.

33
34 A knowledge graph is a network or graph where nodes represent entities or concepts
35 (“things”) and the links or edges of the graph represent relationships between those things
36 (Bollacker et al., 2008). Each node is labelled by a unique identifier, and may have one or
37 more attributes or properties. Each edge of the graph is labelled by the name of the relation-
38 ship it represents. A common representation of a knowledge graph is the linked data triple of
39 subject, predicate, and object, where the subject (e.g., a publication) is connected to an object
40 (e.g., a person) by a predicate (e.g., “author”). Triples are not the only way knowledge graphs

41 can be modelled, but adopting triples means we can use existing technologies such as triple
42 stores and the SPARQL query language (W3C SPARQL Working Group, 2013).

43

44 Knowledge graphs are potentially global in scope, hence rely on global identifiers. Most
45 datasets will have their own local identifiers for the entities they contain, such as species, pub-
46 lications, specimens, or collectors. These identifiers are adequate for local use, but local iden-
47 tifiers also serve to keep data isolated in distinct silos. Hence we need to map identifiers for
48 the same thing between the different silos. This can be done by establishing a “broker” service
49 that asserts identify between a set of identifiers, or by mapping local identifiers to a single
50 global identifier. The case for mapping to a single global identifier (“strings to things”) is at-
51 tractive in terms of scalability (mapping each local identifier to a single global identifier is
52 easier than managing cross mappings between multiple identifiers), and is even more attrac-
53 tive if there are useful services built around that global identifier. For example, Digital Object
54 Identifiers (DOIs) are becoming the standard for identifying academic publications. Given a
55 DOI we can retrieve metadata about the work from CrossRef (“CrossRef”), we can get meas-
56 ures of attention from services such as Altmetric (“Altmetric”), and we can discover the iden-
57 tities of the work’s authors from ORCID (“ORCID”). Furthermore, by agreeing on a central-
58 ised identifier we effectively decentralise the building of the knowledge graph: given a DOI,
59 anybody that links local information to that DOI is potentially contributing to the construction
60 of the global knowledge graph.

61

62 Mapping strings to things give us a way to refer to the nodes in the knowledge graph, but
63 we also need a consistent way to label the edges of the graph. There has been an explosion in
64 vocabularies and ontologies for describing both attributes of entities and their interrelation-
65 ships. While arguments can be made that domain-specific ontologies enable us to represent
66 knowledge with greater fidelity, the existence of multiple vocabularies comes with the cogni-
67 tive overhead of having to decide which term from what vocabulary to use. In contrast to, say,
68 (Senderov et al., 2018) who use several ontologies to model taxonomic publications, the ap-
69 proach I have adopted here is to try and minimise the number of vocabularies employed, and
70 to avoid domain-specific vocabularies where ever possible. For this reason the default vo-
71 cabulary used is schema.org (“Schema.org”), being developed by a consortium of search en-
72 gine vendors including Google, Microsoft, and Yahoo. In addition to simplifying develop-
73 ment, adopting a widely used vocabulary increases the potential utility of the knowledge
74 graph. One motivation for the development of schema.org is to encourage the inclusion of
75 structured data in web pages, helping search engines interpret the contents of those pages. By
76 adopting schema.org in knowledge graphs we can make it easier for developers of biodiver-
77 sity web sites to incorporate structured data from those knowledge graphs directly into their
78 web pages.

79

80

81 There are several different categories of applications that can be built on top of a knowledge
82 graph, for example data reconciliation, data augmentation, and meta-analyses. Reconciliation
83 involves either matching strings to things, or matching entities from different data sources. An

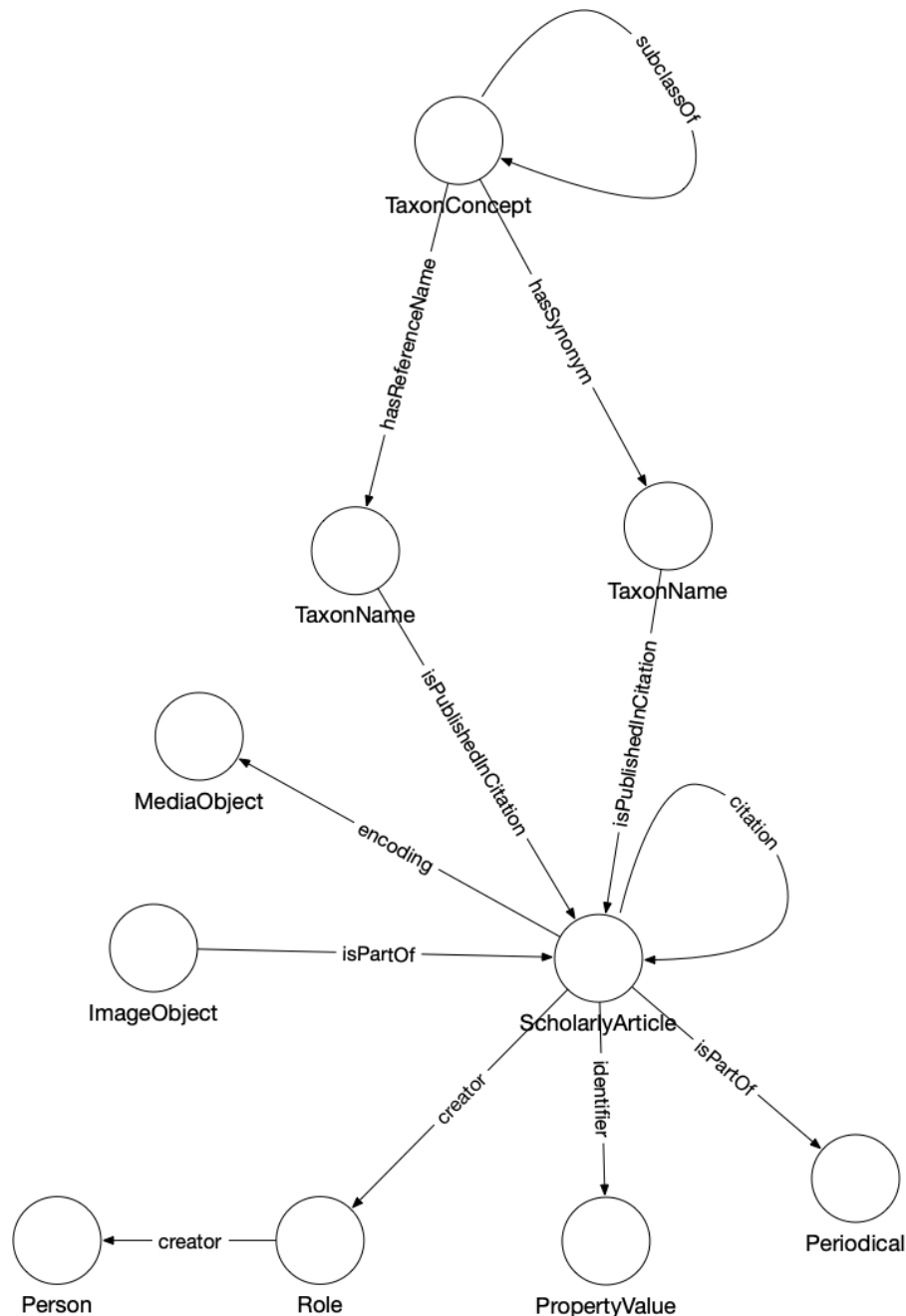
84 example of reconciliation is matching author names to identifiers. Augmentation involves
85 combining data for the same entities from different sources that individually may be incom-
86 plete, but together yield more extensive coverage of those entities. An example is supplement-
87 ing existing imagery of species with figures published in the taxonomic literature. Meta
88 analyses make use of the data aggregated in the knowledge graph to explore larger patterns.
89 There have been numerous studies of patterns of taxonomic activity (Joppa, Roberts & Pimm,
90 2011; Costello, Wilson & Houlding, 2013; Bebber et al., 2013; Grieneisen et al., 2014;
91 Sangster & Luksenburg, 2014; Tancoigne & Ollivier, 2017), typically these studies assembled
92 a custom database, and often this data is not made more widely available, or the data is not
93 actively updated. Having a biodiversity knowledge graph would enable users to ask similar
94 questions but for different taxonomic groups, or different time periods.

95
96 In response to the GBIF 2018 Ebbe Nielsen Challenge I constructed a knowledge graph for
97 the Australian fauna, based on data in the Atlas of Living Australia (ALA) (“Atlas of Living
98 Australia”) and the Australian Faunal Directory (AFD) (“Australian Faunal Directory”). This
99 regional-scale dataset was chosen to be sufficiently large to be interesting, but without being
100 too distracted by issues of scalability. The knowledge graph combines information on taxa
101 and their names, taxonomic publications, the authors of those publications together with their
102 interrelationships, such as publication, citation, and authorship. Constructing the knowledge
103 graph required extensive data cleaning and cross linking. These steps are described below,
104 and examples of the application of the knowledge graph are discussed.

105 **Materials and Methods**

106 **Knowledge graph**

107 The general structure of the knowledge graph is based on (Page, 2013, 2016a). The core enti-
108 ties are taxa, taxonomic names, publications, journals, and people. Figure 1 summarises the
109 relationships between those entities.



110

111 Figure 1. The knowledge graph model used in Ozymandias. Nodes in the graph are repre-
112 sented by circles and are labelled with the *rdf:type* of that node. Nodes are connected by
113 edges in the graph which are labelled by a term from a vocabulary, typically schema.org.

114

115 Taxa and taxonomic names were modelled using the TDWG LSID vocabulary based on
116 (Kennedy et al., 2006). Taxa are nodes in a tree representing the taxonomic classification and
117 are instances of the type *tc:TaxonConcept*. The taxonomic classification is represented by
118 *rdfs:subClassOf* relationship between parent and child taxa (a child is a *rdfs:subClassOf* its
119 parent).

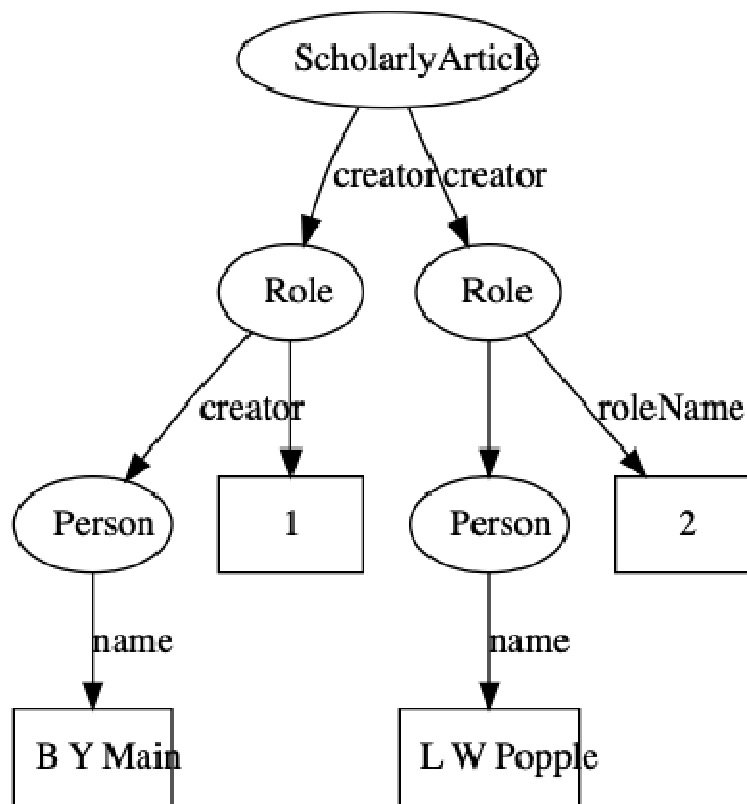
120 Taxonomic names (type *tn:TaxonName*) are connected to the corresponding taxa using
121 relations from the TAXREF vocabulary (Michel et al., 2017) and are typically either accepted
122 names or synonyms. This vocabulary was adopted to because it enables a more direct way of
123 expressing the relationship between taxa and taxonomic names than is possible using the
124 TDWG LSID vocabulary.

125

126 Taxonomic names are published in publications, which were represented using terms
127 from the schema.org vocabulary. In cases where the full text of an article is available as a
128 PDF file I make use of the *schema:encoding* property to link the publication to a
129 *schema:MediaObject* representing the PDF. Articles are linked to the journals they were pub-
130 lished in by the *schema:isPartOf* property.

131

132 Representing ordered lists in RDF is not straightforward, which presents a challenge for
133 expressing relationships such as authorship. Not only is the order of authorship an important
134 feature when formatting a published work for display, it is also useful information when try-
135 ing to reconcile author names (see below). The approach adopted here is to use the
136 *schema:Role* type (Vicki Tardif Holland & Jason Johnson, 2014) . Rather than directly con-
137 nect a publication to an author using, say, the *schema:creator* property, instead the creator of
138 a work is a Role, which in turn has the author as a creator property. The position of author in
139 the list of authors is stored using the *schema:roleName* property (e.g., “1”, “2”, etc.) (Figure
140 2).

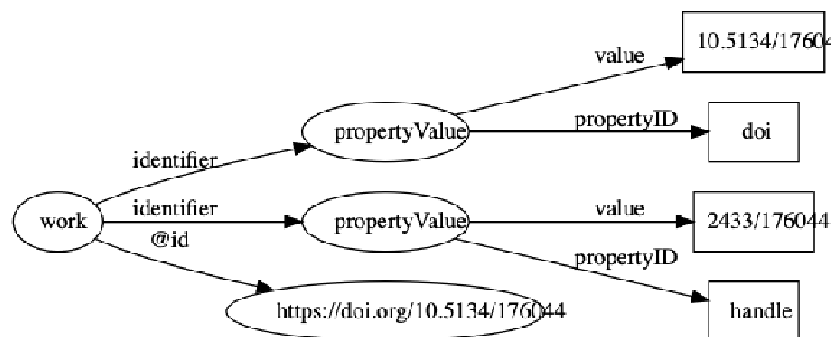


141

142 Figure 2. An example of modelling order of authorship using *schema:Role*. Each author is
143 linked to the article they authored via a *schema:Role* node, which specifies the order of au-
144 thorship for each author. In this example, “B Y Main” is the first author, “L W Pople” is the
145 second author.
146

147 Identifiers

148 Identifiers are both central to any attempt to link data together, and at the same time can be
149 one of the major obstacles to creating links. Ideally identifiers should be globally unique, per-
150 sistent, and each entity would have only a single identifier. In reality, entities may have many
151 identifiers, typically minted by different databases, and identifiers may change, or at least
152 have multiple representations. For example, DOIs may contain upper and lowercase letters,
153 but are actually case insensitive. Some databases may choose to store DOIs in lower case
154 form, others in upper case, or any combination in between. Identifiers typically require
155 dereferencing and the mechanism for this may evolve over time, often for reasons outside the
156 control of the organisation that minted the identifier. DOIs are dereferenced (“resolved”) us-
157 ing the web proxy <https://doi.org>. This proxy recently switched from the HTTP to the HTTPS
158 protocol, immediately rendering out of date any DOIs stored URLs starting with the prefix
159 “http://”. To minimise the impact of these kinds of changes, Ozymandias stores identifiers
160 both as URLs (where appropriate) but also as property-value pairs (*schema:PropertyValue*)
161 where the *schema:value* property stores the identifier string stripped of any dereferencing pre-
162 fix. For example, a DOI <https://doi.org/10.5134/176044> would be stored as a
163 *schema:PropertyValue* with *schema:propertyID* “doi” and *schema:value* “10.5134/176044”
164 (Figure 3).



165

166 Figure 3. Storing identifiers using *schema:PropertyValue*. The work has two identifiers, a
167 DOI <https://doi.org/10.5134/176044> and a Handle <https://hdl.handle.net/2433/176044>. These
168 are stored as *schema:PropertyValue* pairs.
169

170 Citations

171 One paper citing another can be represented by a direct link between two identifiers, for ex-
172 ample a link between the DOIs of the citing and the cited work. CrossRef provides lists of
173 literature cited for many of the works in its database, and many of these cited works them-
174 selves have DOIs Hence if we have a DOI for a work we can immediately populate the triple

175 store with citation links. This works well if both works have a DOI, but many taxonomically
176 relevant works do not have these identifiers. Even for those works that do have DOIs, these
177 may not have been available at the time the citing work was deposited by a publisher, for ex-
178 ample, if the cited work has only recently been assigned a DOI.

179

180 To expand the citation links beyond just those works with DOIs I also generated an iden-
181 tifier for each work modelled on the Serial Item and Contribution Identifier (SICI). This iden-
182 tifier comprised the International Standard Serial Number (ISSN) of the journal, together with
183 the volume, and the starting page. This triple uniquely identifies most articles, and is easy to
184 generate. SICIs were generated for works harvested from the Australian Faunal Directory, and
185 from the lists of literature cited obtained from CrossRef, and were stored as
186 *schema:PropertyValue* pairs in the same way as DOIs and other identifiers. By matching SI-
187 CIs it was possible to expand citation links beyond those where both works had DOIs.

188

189 **Populating the knowledge graph**

190 Perhaps the biggest challenge in constructing a knowledge graph is to map names or descrip-
191 tions of entities to one or more globally unique identifiers. In some cases the sources data will
192 already have identifiers. Taxa in the ALA each have a unique identifier (a LSID), as do taxa
193 and publications in the AFD (which use UUIDs). The ALA and AFD share the same taxon
194 identifiers, which makes linking the two databases straightforward. However, these identifiers
195 are local in the sense that they are primary keys for local databases that have been converted
196 into URLs. The knowledge graph can only grow if we use external identifiers that are shared
197 by other databases, or at least map local identifiers onto those external identifiers. For publi-
198 cations this is straightforward in the sense that a publication in a database of Australian ani-
199 mals can be unambiguously mapped onto the publication in, say, a database for Japanese
200 animals. However, a taxon as defined in the Australian Faunal Directory may not correspond
201 exactly to a taxon with the same name in another.

202

203 **Reconciling works**

204 For the works in AFD I searched for DOIs using the API provided by CrossRef. If a reference
205 was found the associated DOI was assigned to that reference. CrossRef is not the only regis-
206 tration agency for DOIs, there are several others that are used by digital libraries and publish-
207 ers, such as DataCite, the multilingual European Registration Agency (mEDRA), and Airiti
208 (華藝數位). Most of these agencies lack the discovery services provided by CrossRef, so for
209 these DOIs I harvested the article metadata using a combination of web services and screen
210 scraping, created a local MySQL database to store the metadata, and used that database to
211 match references to DOIs. This database was also used to match articles to other classes of
212 identifiers, such as Handles and URLs.

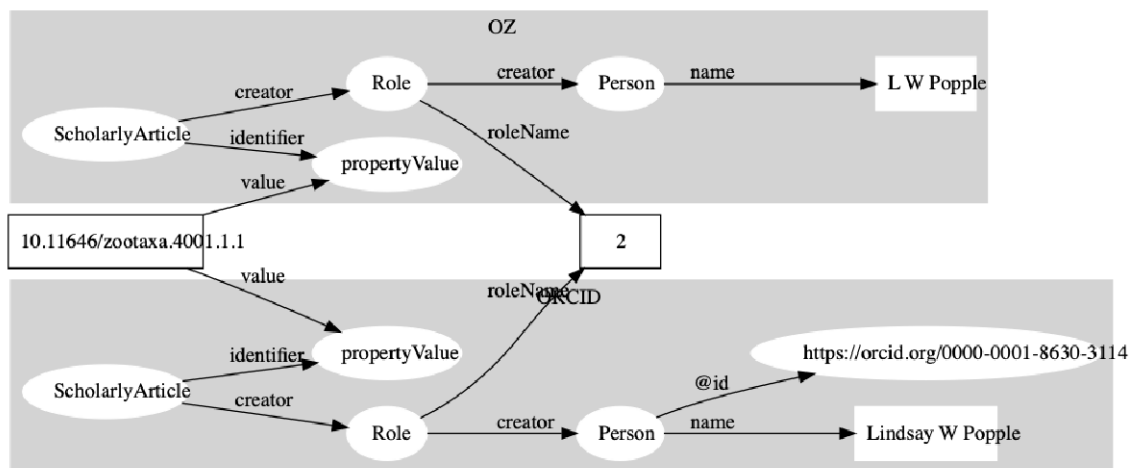
213

214 Australian natural history institutions are significant publishers of biodiversity literature,
215 and much of this has been scanned by the Biodiversity Heritage Library in Australia. As a
216 consequence many of the articles in the knowledge graph were available in my BioStor pro-

217 ject (Page, 2011). Identifiers for these articles were found by matching using the BioStor
218 OpenURL service.
219

220 Reconciling authors

221 Multiple approaches were used to match author names to external identifiers. Metadata for
222 DOIs from CrossRef would sometimes include ORCID ids for authors. The ORCID record
223 for each ORCID id was retrieved using the ORCID API and converted to a set of RDF triples
224 linking the identifiers for a work (e.g., DOI) to a person's ORCID. These triples modelled the
225 order of authorship using *schema:Role* as described above. Similarly, I parsed Wikispecies
226 pages and extracted bibliographic records for works identified by a DOI, and constructed tri-
227 ples linking the work to its authors where those authors had their own Wikispecies page.
228 Hence to match authors in the knowledge graph to authors in ORCID or Wikispecies, we can
229 ask whether the same pairing of work and author name appears in both databases. For exam-
230 ple, we can retrieve the second author of a work in the knowledge graph and in ORCID by
231 querying by DOI for the work and restricting the value of *schema:roleName* to "2" (Figure 4).
232 As a final check we can compare the author names and accept only those names whose simi-
233 larity exceeds a threshold. In this way we can automate the matching authors across data-
234 bases.



235

236 Figure 4. Matching author records in two different databases. In this example the article with
237 DOI 10.11646/zootaxa.4001.1.1 occurs in both Ozymandias (OZ) and ORCID. Using a
238 SPARQL query we retrieve the name of the second author in the two databases: "L W Pop-
239 ple" in Ozymandias and "Lindsay W Popple" in ORCID. Given the similarity in the names,
240 we conclude that the two authors are the same, and we can assign the ORCID for Lindsay W
241 Popple (<https://orcid.org/0000-0001-8630-3114>) to "L W Popple" in Ozymandias.

242

243

244 **Data sources**

245 I used several different strategies to convert data into the triples required for the knowledge
246 graph. If the source data was in the form of CSV files (e.g., the Australian Faunal Directory)
247 it was imported into a MySQL database, and PHP scripts were written to further clean the
248 data and map it to any external identifiers. Once the data was cleaned and linked, a PHP script
249 was used to export the data in N-triples format.

250

251 Several sources of data (Atlas of Living Australia, CrossRef, ORCID, Wikispecies, and
252 Biodiversity Literature Repository) were accessed via their APIs. For ALA a list of all animal
253 taxa was obtained from the ALA web site, then the JSON record for each taxon was har-
254 vested. For CrossRef, data was harvested for just those DOIs found by the bibliographic
255 string to DOI mapping process described above. These DOIs were also submitted to a custom
256 script that queried the ORCID database to discover whether any authors had works with those
257 DOIs in their ORCID profile. If this was the case, the corresponding ORCID profile was
258 downloaded. Each DOI was also used as a query term for searching Wikispecies using its API
259 with the “list” parameter set to “exturlusage” to find wiki pages that mentioned that DOI.
260 Pages found were retrieved in XML format using the API, any references on that page parsed
261 and converted into JSON. All JSON documents obtained from these sources were stored in
262 CouchDB databases and custom CouchDB views were written in Javascript to convert the
263 JSON documents into N-triples.

264

265 By default Ozymandias treats individual publications as a single, monolithic entity. How-
266 ever, some publishers such as PLoS and Pensoft provide DOIs for component parts of an arti-
267 cle, such as individual figures. (Egloff et al., 2017) have argued that even if a taxonomic arti-
268 cle itself is copyrighted, the individual figures are not eligible for copyright, and hence extract
269 and assign DOIs to large numbers of figures extracted from journals such as *Zootaxa*. These
270 figures, together with ones sourced from open access journals are available through the Bio-
271 diversity Literature Repository (“Biodiversity Literature Repository”) (BLR). The BLR is
272 hosted by Zenodo (<https://zenodo.org>) and each publication and figure has a unique identifier
273 (typically a DOI), and metadata for each publication and figure is available as JSON-LD. This
274 means data from the BLR can be directly incorporated into a triple store. However for this
275 project I wanted just a subset relevant to publications on the Australian fauna, and so I created
276 a CouchDB version of the BLR and write scripts to match publications from the AFD to the
277 corresponding record in the BLR. Metadata for each matching publication and its associated
278 figures were then retrieved directly from Zenodo.

279

280 **Knowledge graph**

281 The knowledge graph was implemented as a triple store using Blazegraph 2.1.4 running on a
282 Windows 10 server, with a nginx web server acting as a reverse proxy. N-triples for different
283 categories of data (e.g., taxa, publications, etc.) were partitioned using named graphs and up-
284 loaded to the triple store. This made it easier to manage sets of data, for example the biblio-
285 graphic data could be deleted and reloaded by simply deleting all triples in the corresponding

286 named graph, rather than having to delete the entire knowledge graph. It also facilitated some
287 queries, such as author matching across multiple data sources where distinguishing between
288 data source was an essential part of the query.
289

290 **Search**

291 Being able to simply search for relevant documents by typing in one or more terms is a fea-
292 ture users expect from almost any web site. To implement search, basic information on taxa
293 and publications was encoded into a simple JSON document (one per entity) and these JSON
294 documents were indexed using an instance of Elasticsearch 6.3.1 hosted on Google's Com-
295 pute Engine.
296

297 **Web interface**

298 Designing a semantic web browser to display a richly interconnected data set is a challenging
299 task (Quan & Karger, 2004). For Ozymandias the goal was to have a simple interface which
300 encouraged the user to explore connections between taxa, publications, and people. Apart
301 from the home page, there are two main page types in the web interface for Ozymandias. The
302 first is the search interface which is a simple list of the entities that best match the search
303 terms. Clicking on any member of that list leads to the second page type, which is a display of
304 the entity itself. This display comprises three columns. The left column displays core facts
305 about the entity. These are typically triples that have the entity as their subject, or are one
306 edge away in the knowledge graph (such as thumbnail images), and so can be retrieved from
307 the knowledge graph using either SPARQL DESCRIBE or CONSTRUCT queries. The mid-
308 dle column displays connections between the main entity on the page and related entities in
309 the knowledge graph (such as authors of a paper, taxonomic names mentioned in a work,
310 etc.), and is populated by SPARQL queries. The rightmost column is used to display the re-
311 sult of searching external sources for information relevant to the entity displayed on the page.
312 Hence, unlike columns one and two, these queries are not SPARQL queries to the local
313 knowledge graph.

314 **Results**

315
316 Ozymandias can be viewed at <https://ozymandias-demo.herokuapp.com>. Source code is avail-
317 able on GitHub <https://github.com/rdmpage/ozymandias-demo>. Below I describe the web in-
318 terface to Ozymandias, and outline some of the exploratory analyses that can be undertaken
319 using the underlying knowledge graph. Where the results are based on SPARQL queries,
320 those queries are listed in the Supplementary material.
321

322 **Web interface**

323 A screenshot of the web interface is shown in Figure 5. This shows the three-column layout
324 used to display an entity, its relationships within the knowledge graph, and any known exter-
325 nal relationships.

Ozymandias - a biodiversity knowledge graph

Search

Revision of genera of the dragonets (Pisces : Callionymidae)
Publications of the Seto Marine Biological Laboratory
1982; 27(1/3): 77-131 figs 1-30

T. Nakabo

<https://doi.org/10.5134/176044>
<https://hdl.handle.net/2433/176044>

Connections within this knowledge graph.

External knowledge graphs.

DOI in Wikidata
Q56208522
No ORCID links to this DOI
Page 3

Taxa in this work.

Pseudocallurichthys delicatulus (Smith, 1963)
Repomucenus limiceps (Ogilby, 1908)
Neesychiropus Nalbant, 1980
Bathycallionymus moretonensis (Johnson, 1971)
Minesychiropus Nakabo, 1982
Pterosychiropus splendidus (Hesse, 1927)
Neesychiropus molitorius (Schultz, 1960)
Pseudocallurichthys simplicioris (Valetonianus, 1837)
Neesychiropus ocellatus (Pallas, 1770)
Eocallionymus papilio (Günther, 1864)
Paradiplomammus Nakabo, 1982
Repomucenus meridionalis (Szwedzi, 1965)

Materials and Methods

The specimens used for making the generic revision of the Callionymidae are shown in Table 1. These specimens have been preserved in 10 % formalin, 70 % ethyl alcohol or 40 % isopropyl alcohol and are deposited at the following institutions: Australian Museum, Sydney; Academy of Natural Sciences of Philadelphia; British Museum (Natural History), London; California Academy of Sciences, San Francisco; Department of Biology, Faculty of Science, Kochi University; Department of Biology, University of the Ryukyus; Department of Fisheries, Faculty of Agriculture, Kyoto University; Department of Zoology, University Museum, University of Tokyo; Fisheries Research Station, Kyoto University; Far Sea Fisheries Research Laboratory; Laboratory of Marine Zoology, Faculty of Fisheries, Hokkaido University; Marine Science Museum, Tokai University; Museum of Tokyo University of Fisheries; Miyaaki University; Naturhistorisches Museum Wien; Department of Zoology, Natural Science Museum, Tokyo; Queensland Museum, Brisbane; J.L.B. Smith Institute of Ichthyology, Rhodes University, Grahamstown; Senckenberg

Fig. 1. Diagrammatic illustration showing the measuring methods of the portions of callionymid body. 1, total length; 2, standard length; 3, body width; 4, body depth; 5, nasal peduncle depth; 6, prenasal length; 7, caudal fin length; 8, head length; 9, eye diameter; 10, snout length; 11, upper jaw length; 12, mouthful width; 13, 1st dorsal spine length; 14, 2nd dorsal spine length; 15, 3rd dorsal spine length; 16, 4th dorsal spine length; 17, 1st dorsal eye length; 18, 2nd dorsal eye length; 19, 1st anal ray length; 20, 2nd anal ray length; 21, 3rd anal ray length; 22, 4th anal ray length; 23, 5th anal ray length; 24, 6th anal ray length; 25, 7th anal ray length.

326

327

328 Figure 5. Web interface to Ozymandias knowledge graph displaying information for an arti-
329 cle. The left column displays a summary of the article, and a PDF viewer (only available if
330 content is freely accessible). The middle column displays related content from the knowledge
331 graph, such as taxa mentioned in the article. The right column shows the result of searches in
332 external web sites for related information, in this case is displays the identifier for Wikidata
333 item that corresponds to this article. To view this page live go to [https://ozymandias-](https://ozymandias-demo.herokuapp.com/?uri=https://biodiversity.org.au/afd/publication/3e0c1402-de05-4227-9df3-803e68300623)
334 [demo.herokuapp.com/?uri=https://biodiversity.org.au/afd/publication/3e0c1402-de05-4227-](https://biodiversity.org.au/afd/publication/3e0c1402-de05-4227-9df3-803e68300623)
335 [9df3-803e68300623](https://biodiversity.org.au/afd/publication/3e0c1402-de05-4227-9df3-803e68300623).
336

337 The first example is a publication, in this case (Nakabo, 1982). The first column summa-
338 rizes basic data about the publication, and if the full text is available it is displayed using ei-
339 ther a PDF viewer, or a simple image viewer in the case of scanned images. The second col-
340 um lists taxa associated with the publication. For publications with identifiers such as DOIs
341 the third column displays whether a record with that DOI exists in external sources such as
342 Wikidata and ORCID.

The screenshot shows the Ozymandias interface for author L. W. Popple. The header includes the title "Ozymandias - a biodiversity knowledge graph" and a search bar. The main content is divided into three columns:

- Left Column:** Author profile for L. W. Popple, a list of works by this author, and a list of publications from 2000 and 2010. The 2000 publication is "A new species of Cicadetta Amyot (Hemiptera: Cicadidae) from Queensland, with notes on its calling song" from the Australian Entomologist. The 2010 publications include "A new cicada genus and redescription of Pauropsalta subolivacea Ashton (Hemiptera: Cicadidae) from eastern Australia" and "An analysis of the calling song of Burbunga mouldsi Olive (Hemiptera: Cicadidae)".
- Middle Column:** "Connections within this knowledge graph" section containing:
 - Top five coauthors:** A. Ewart (4), D. L. Emery (3), D. C. Marshall (1), K. B. R. Hill (1), N. J. Emery (1).
 - Top ten journals:** Zootaxa (7), Australian Entomologist (5), Records of the Australian Museum (1), Memoirs of the Queensland Museum - Nature (1).
 - Top 20 taxa:** A hierarchical tree diagram starting from ANIMALIA and ending with Cicadettini, listing higher taxa such as ARTHROPODA, HEXAPODA, INSECTA, Pterygota, HEMIPTERA, AUCHENORRHYNCHA, CICADOMORPHA, CICADOIDEA, CICADIDAE, and CICADETTINAE.
- Right Column:** "External knowledge graphs" section containing:
 - ORCID match:** Links to L. W. Popple, L. W. POPPLE, L.W. POPPLE, Lindsay W. Popple, and LINDSAY W. POPPLE.
 - Wikispecies match:** Links to Popple, L.W., Wikidata, and Q21393780.

343

344

345 Figure 6. Information about an author displayed in Ozymandias. The left column lists the au-
346 thor's publications, the middle column uses the knowledge graph to identify coauthors, ven-
347 ues for publication, and the taxonomic expertise of the author, the right column displays in-
348 formation from external sources. To view live go to [https://ozymandias-](https://ozymandias-demo.herokuapp.com/?uri=https://biodiversity.org.au/afd/publication/%23creator/l-w-popple)
349 [demo.herokuapp.com/?uri=https://biodiversity.org.au/afd/publication/%23creator/l-w-popple](https://biodiversity.org.au/afd/publication/%23creator/l-w-popple).
350

351

352

353

354

355

356

357

358

The second example (Figure 6) displays information for an author, including a list of publications, coauthors, journals the author publishes in, and a summary of their taxonomic expertise. This later diagram is computed by using a SPARQL query to find the top 20 taxa the author has published on. For each taxon the query uses a property path expression to retrieve the list of higher taxa each taxon belongs to, and a Javascript script assembles those lists into a tree. The third panel displays the results of matching the author to author identifiers using external web services, in this case discovering the author's ORCID id and entry in Wikidata.

Ozymandias - a biodiversity knowledge graph

Search

Acupalpa Kröber, 1912
<https://bie.ala.org.au/species/urn:lsid:biodiversity.org...>

ANIMALIA
ARTHROPODA
HEXAPODA
INSECTA
Pterygotes
DIPTERA
ORTHORRHAPHA
ASILOIDEA
THEREVIDAE
Agapophytinae

Children

Connections within this knowledge graph.

Names for this taxon.

- ✓ Acupalpa Kröber, 1912
- = Acupalpa Kröber, 1912
 - Die Thereviden der indoaustralischen Region. (Dipt)

External knowledge graphs.

Taxon in GBIF
1566276

359

360 Figure 7. Information about the genus *Acupalpa* Kröber, 1912 displayed in Ozymandias. The
361 display includes the species in the genus, details about the publication of the name *Acupalpa* ,
362 and a link to the taxon in GBIF. Live version at [https://ozymandias-](https://ozymandias-demo.herokuapp.com/?uri=https://bie.ala.org.au/species/urn:lsid:biodiversity.org.au:afd.taxon:111fc7e9-0265-453e-8e60-1761e42efc9a)
363 [demo.herokuapp.com/?uri=https://bie.ala.org.au/species/urn:lsid:biodiversity.org.au:afd.taxon](https://ozymandias-demo.herokuapp.com/?uri=https://bie.ala.org.au/species/urn:lsid:biodiversity.org.au:afd.taxon:111fc7e9-0265-453e-8e60-1761e42efc9a)
364 [:111fc7e9-0265-453e-8e60-1761e42efc9a](https://ozymandias-demo.herokuapp.com/?uri=https://bie.ala.org.au/species/urn:lsid:biodiversity.org.au:afd.taxon:111fc7e9-0265-453e-8e60-1761e42efc9a).
365

366

367 Figure 7 shows the view of a taxon, in this case genus *Acupalpa* Kröber, 1912. We see
368 the member species of the genus, the taxonomic hierarchy of the genus (generated using a
369 SPARQL property path query) and, where available, a thumbnail image from the ALA. The
370 second column lists the taxonomic names associated with the genus, together with the publi-
371 cations that made those names available. The third column shows the results of mapping the
372 taxon to one or more external taxonomic databases, in this case GBIF.
373

374

374 Wherever possible, Ozymandias uses thumbnail images from ALA to illustrate taxa.
375 However, many taxa lack images. Figure 8 shows an example where the ALA has no image
376 for a taxon (*Trigonopterus cooktownensis*). Because the taxon, its name, the publication, and
377 the figures in that publication extracted by the Biodiversity Literature Repository are all part
378 of the knowledge graph, we can traverse the graph and discover that an image for that species
379 was published in (Riedel & Tänzler, 2016) .

The screenshot shows the Ozymandias biodiversity knowledge graph interface. The title is "Ozymandias - a biodiversity knowledge graph" with "Tree SPARQL" in the top right. A search bar contains the text "Trigonopterus cooktownensis Riedel, 2016". Below the search bar, a taxonomic tree is displayed, starting with "ANIMALIA" and ending with "Trigonopterus". To the right of the search bar, there is a square with a question mark. Further right, there are sections for "Connections within this knowledge graph.", "Names for this taxon." (listing "Trigonopterus cooktownensis Riedel, 2016" and "Trigonopterus cooktownensis Riedel, 2016" with a note about a revision), and "Images of taxon" (showing two images of the weevil). On the far right, there is a section for "External knowledge graphs." with a link to "Taxon in GBIF" (8782535).

380

381

382 Figure 8. Augmenting data using knowledge graph. The Atlas of Living Australia did not
383 have an image for *Trigonopterus cooktownensis* at the time it was harvested by Ozymandias,
384 hence the “?” displayed in the square in the left column. However, the original description of
385 that species did include images which are available in the Biodiversity Literature Repository,
386 and hence are displayed by Ozymandias in the middle column. Live example

387 <https://ozymandias->

388 [demo.herokuapp.com/?uri=https://bie.ala.org.au/species/urn:lsid:biodiversity.org.au:afd.taxon:](https://ozymandias-demo.herokuapp.com/?uri=https://bie.ala.org.au/species/urn:lsid:biodiversity.org.au:afd.taxon:14fec1f-9d2a-496b-9b98-ec691289b5ce)
389 [14fec1f-9d2a-496b-9b98-ec691289b5ce](https://ozymandias-demo.herokuapp.com/?uri=https://bie.ala.org.au/species/urn:lsid:biodiversity.org.au:afd.taxon:14fec1f-9d2a-496b-9b98-ec691289b5ce).

390

391 **Strings to things**

392 Most of the work on data cleaning and linking was devoted to matching string representations
393 of publications to the corresponding digital identifiers. The result of this matching provides us
394 with an estimate of how many publications have been digitised and hence are potentially
395 available online.

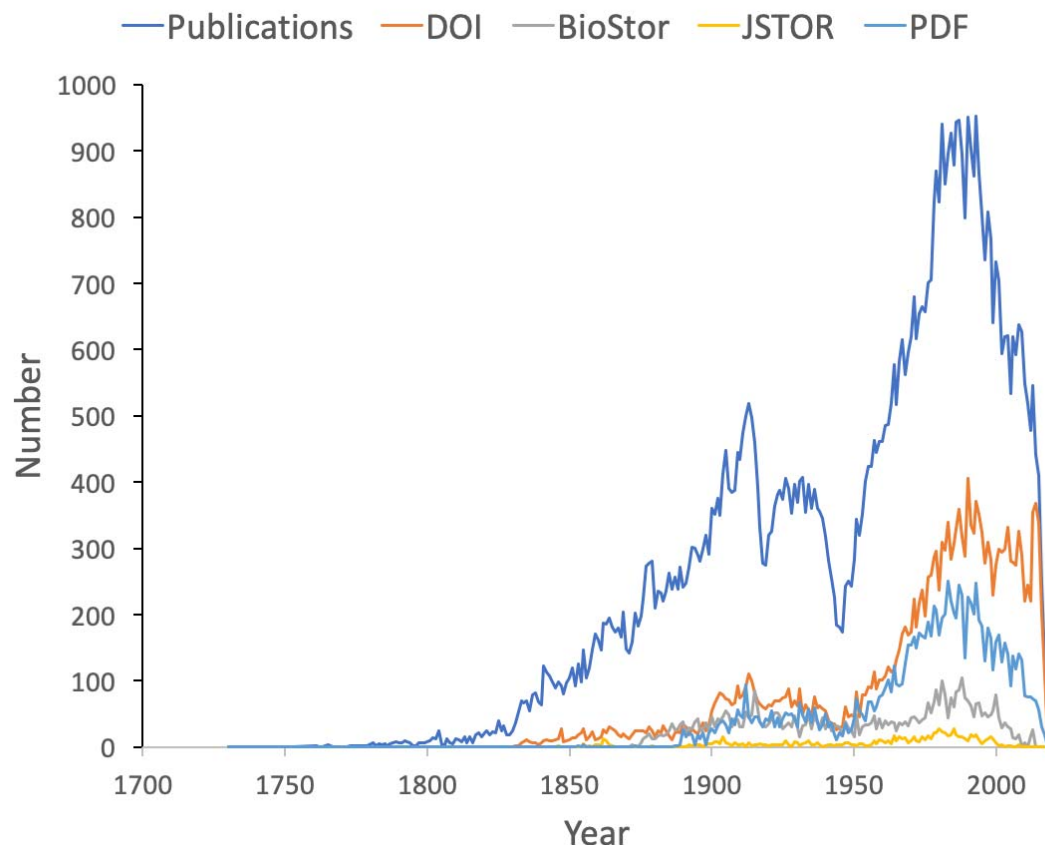
396

397 Figure 9 shows the distribution of publications over time, together with the numbers that
398 have been matched to digital identifiers. The pattern of publication shows three prominent
399 dips. The first two correspond to the two world wars in the twentieth century, the third dip
400 occurs from the mid-1990’s to the present day. Given that the AFD is retrospectively collect-

401 ing publication data, it is not clear to what extent this decline in recent publications represents
402 an actual decline in activity versus a under sampling the most recent literature.

403

404 Many publications lack a digital identifier, suggesting that a considerable amount of the
405 relevant literature has not been digitised. However, this may be overstated as the matching
406 was done by a single individual working to a deadline (in this case the Ebbe Nielsen Chal-
407 lenge submission date). As more effort is expended on matching records the gap between the
408 number of publications and the number of publications online is likely to decrease.



409

410 Figure 9. Plot of publications over time. As well as the total number of publications for each
411 year, the chart shows the numbers of publications that have a digital identifier (DOI, BioStor,
412 or JSTOR) or have a PDF available online.

413

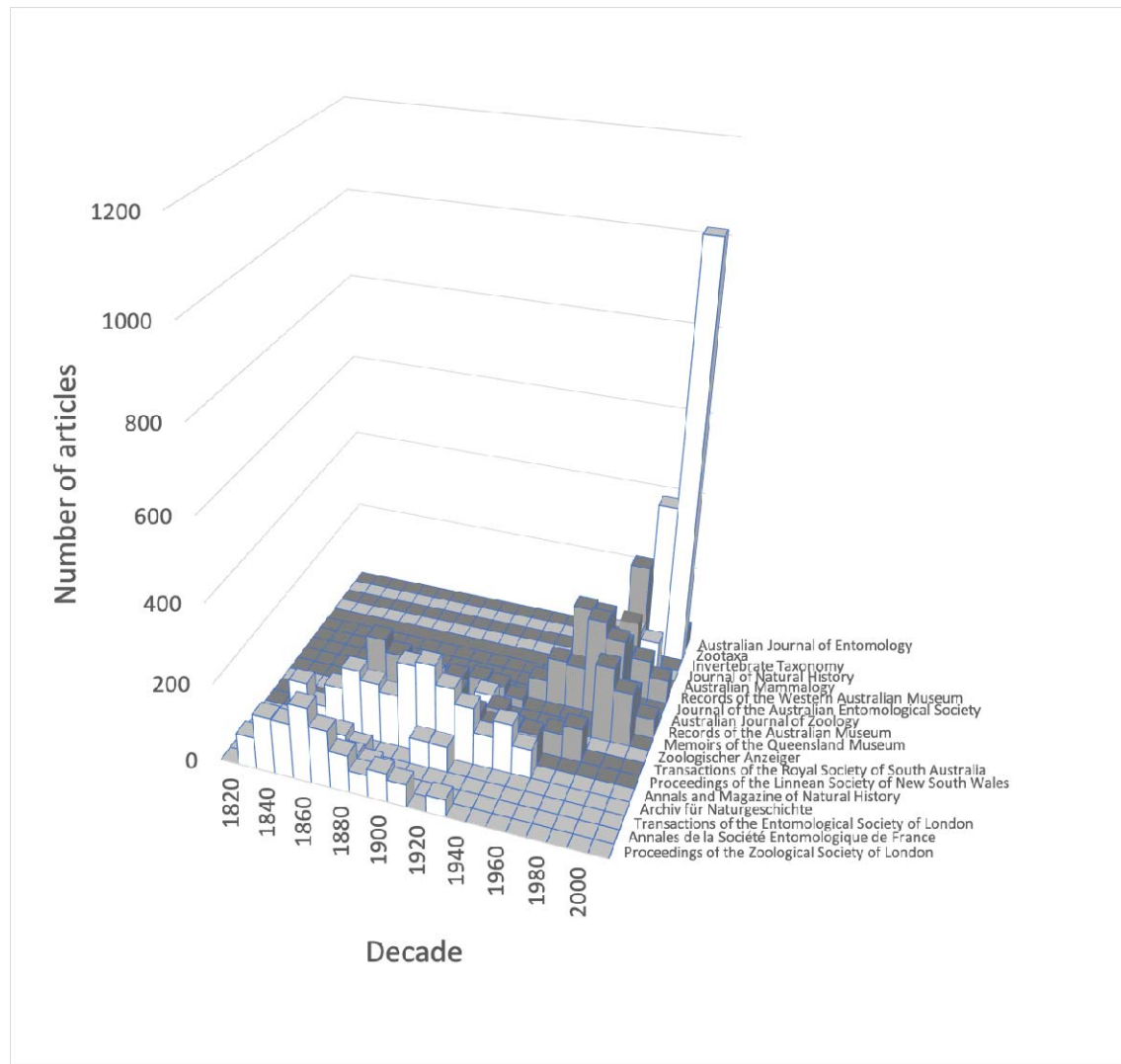
414 **Linking authors to identifiers**

415 We can measure the uptake of ORCID ids for researchers working on the Australian fauna by
416 using DOIs to match works in the knowledge graph to works in the ORCID database. ORCID
417 was launched in 2012, for period from 2012 to the present day the knowledge graph contains
418 2302 distinct author names. Matching DOIs for the works those authors published to the OR-
419 CID database shows that 346 (15%) of authors publishing in that time period have ORCID

420 ids. This number is likely to be an underestimate as not all works in ORCID have DOIs (and
421 ORCID records sometimes omit DOIs for works that have them), but it suggests limited adop-
422 tion of ORCIDs amongst taxonomists and other biodiversity researchers.

423 **Changes in taxonomic publications over time**

424 To explore the publication history of taxonomic research on Australian animals for each dec-
425 ade from 1820 to 2020 I found the ten journals that had the most articles in the knowledge
426 graph. The numbers of articles in each journal were plotted for each decade (Figure 10). Over
427 time different journals have been dominant venues for publishing taxonomic work. In the
428 18th century British or other European journals dominated, such as *Proceedings of the Zoo-*
429 *logical Society* and *Annals and Magazine of Natural History*, although the local journal *Pro-*
430 *ceedings of the Linnean Society of New South Wales* (establish 1875) was a major outlet for
431 taxonomic work. In the mid to late 20th century Australian journals, typically published by
432 museums or by the Commonwealth Scientific and Industrial Research Organisation (CSIRO)
433 were the primary venues for taxonomic papers on the Australian fauna. However, the last
434 decade has seen the spectacular rise of the “megajournal” *Zootaxa*, published in New Zealand
435 but with global coverage. Hence, taxonomic publication in Australia has gone from an early
436 colonial period where much of it was published overseas, to a national period where many
437 papers were published in local journals, culminating in the present situation where interna-
438 tional journals such as *Zootaxa* and, to a lesser extent *Zookeys*, dominate.



439

440 Figure 10. Patterns of publication of taxonomic work on Australian animals 1820-2020. Chart
441 shows the numbers of publications in the top ten journals for each decade. The 19th and early
442 20th centuries are dominated by European journals, by the mid 20th century most taxonomy
443 was published in Australian journals, more recently international journals such as *Zootaxa* are
444 increasingly important.

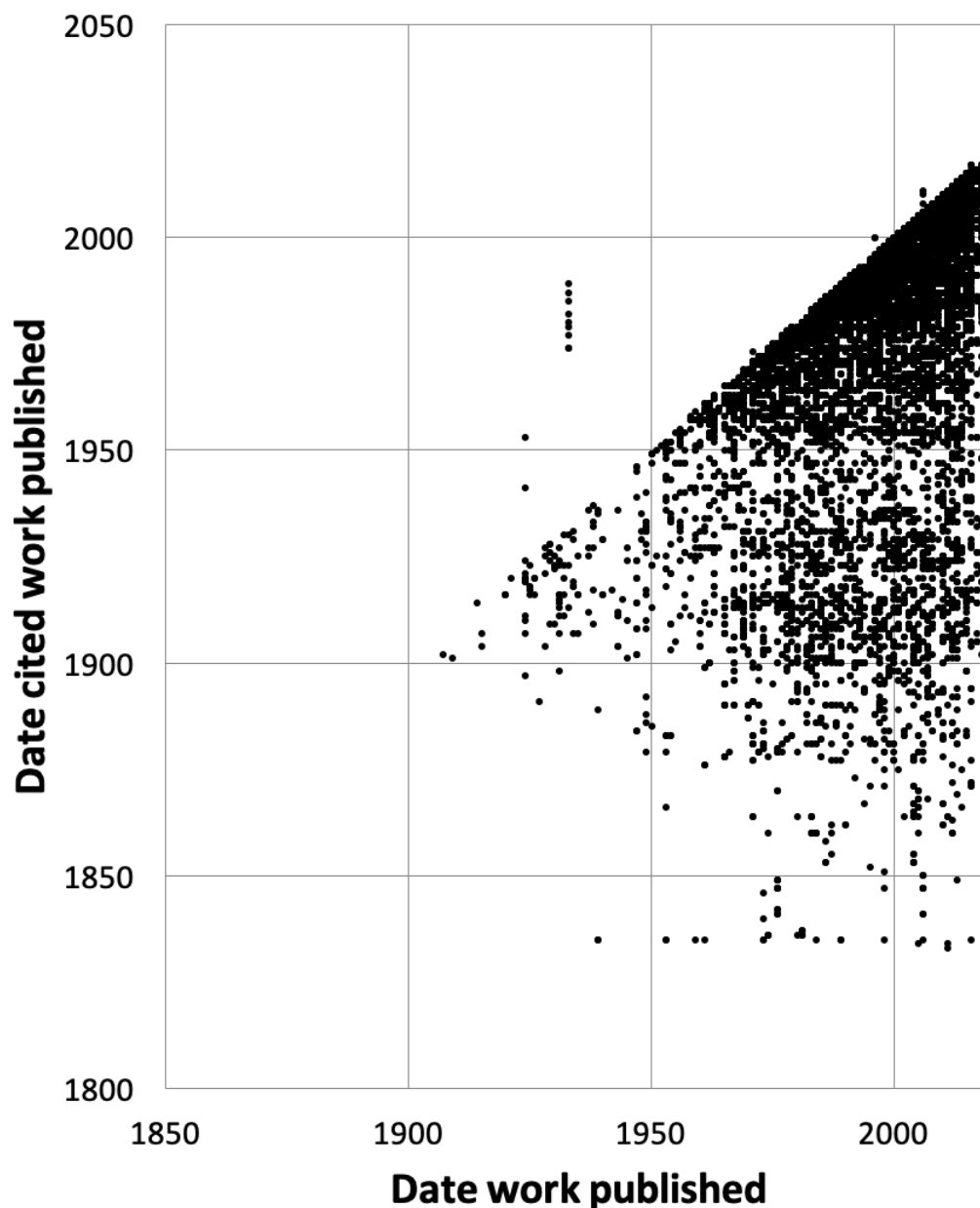
445

446

447 **Citations and taxonomy as long data**

448 Taxonomy is a “long data” discipline (Page, 2016b). In some scientific fields published pa-
449 pers have a short citation half-life and hence are relatively ephemeral, quickly losing their
450 relevance as the “research front” moves on (de Solla Price, 1965). The rise of academic
451 search engines such as Google Scholar may increase the discoverability of the older literature
452 (and hence increasing its likelihood of being cited, (Verstak et al., 2014)), but for many fields
453 the older literature fades from importance. In contrast, the taxonomic literature is essentially
454 ageless - any published work is potentially relevant. Part of this relevance reflects the impor-

455 tance of priority in biological nomenclature, given competing names for the same taxon in
456 general the oldest name wins. Another factor is the sheer number of species and the relative
457 paucity of published knowledge on many of those species. May (1988) estimated that for pub-
458 lications in the period 1978 to 1987 for insects there were on average 0.02 papers per species
459 per year, for beetles it was 0.01 papers. Hence a researcher may have to search back through a
460 hundred years of literature in order to find mention of a specific beetle species.



461

462 Figure 11. Dates of publication of works cited against the date of publication of the cited
463 work. Each point represents the (x, y) pair (publication date, cited publication date). All cited
464 works must, by definition, be published in the same year or earlier, and hence the points fall

465 on or below the diagonal. The few points that are above the diagonal represent errors in
466 CrossRef's metadata.

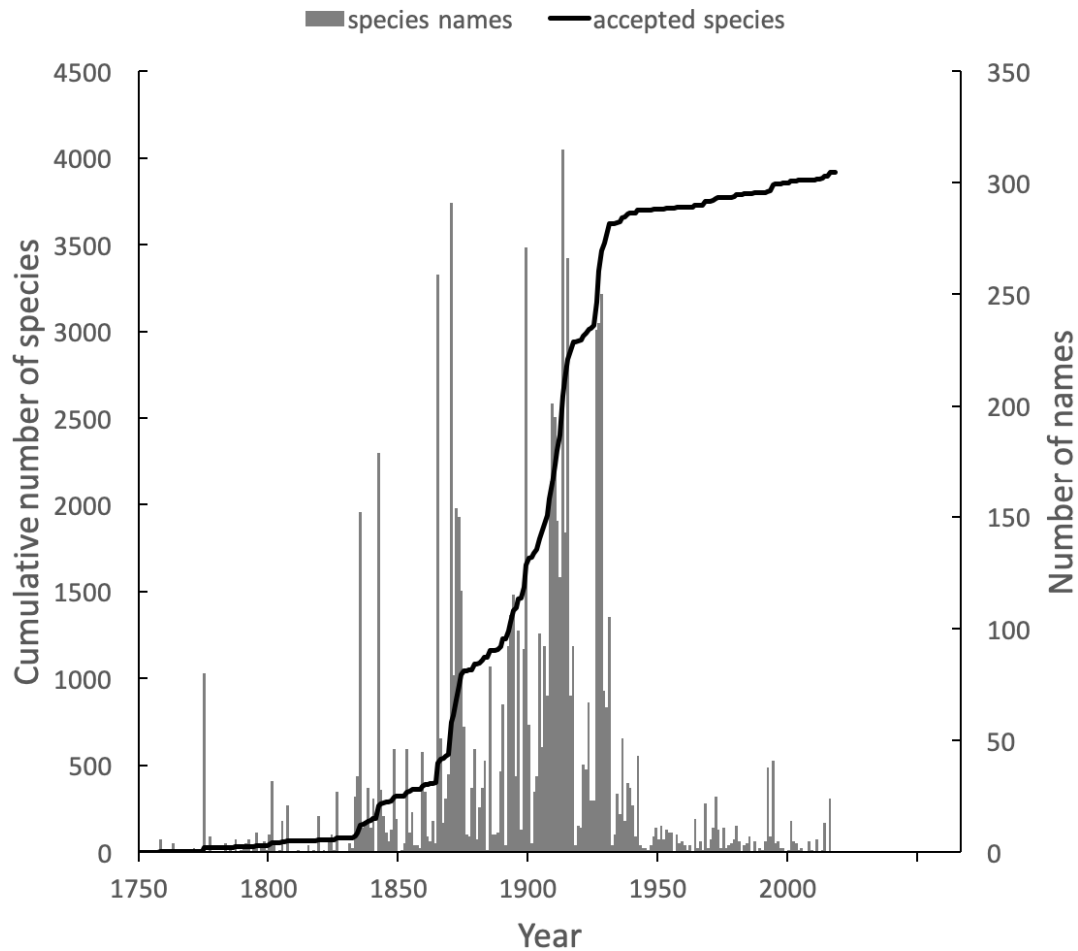
467

468 To explore the citation graph for publications on the Australian fauna I queried each cita-
469 tion relationship for the dates of publication of the citing and the cited works. The relationship
470 between these two dates (Figure 11) highlights the enduring value of the older taxonomic lit-
471 erature. If taxonomic work cited only recent publications then the points in Figure 11 would
472 fall on or close to the diagonal. However, even papers published recently (top right of the
473 chart) cite older literature (represented by the vertical columns of dots below each year), and
474 hence much of the area below the diagonal is occupied.

475

476 **History of species discovery in different taxonomic groups**

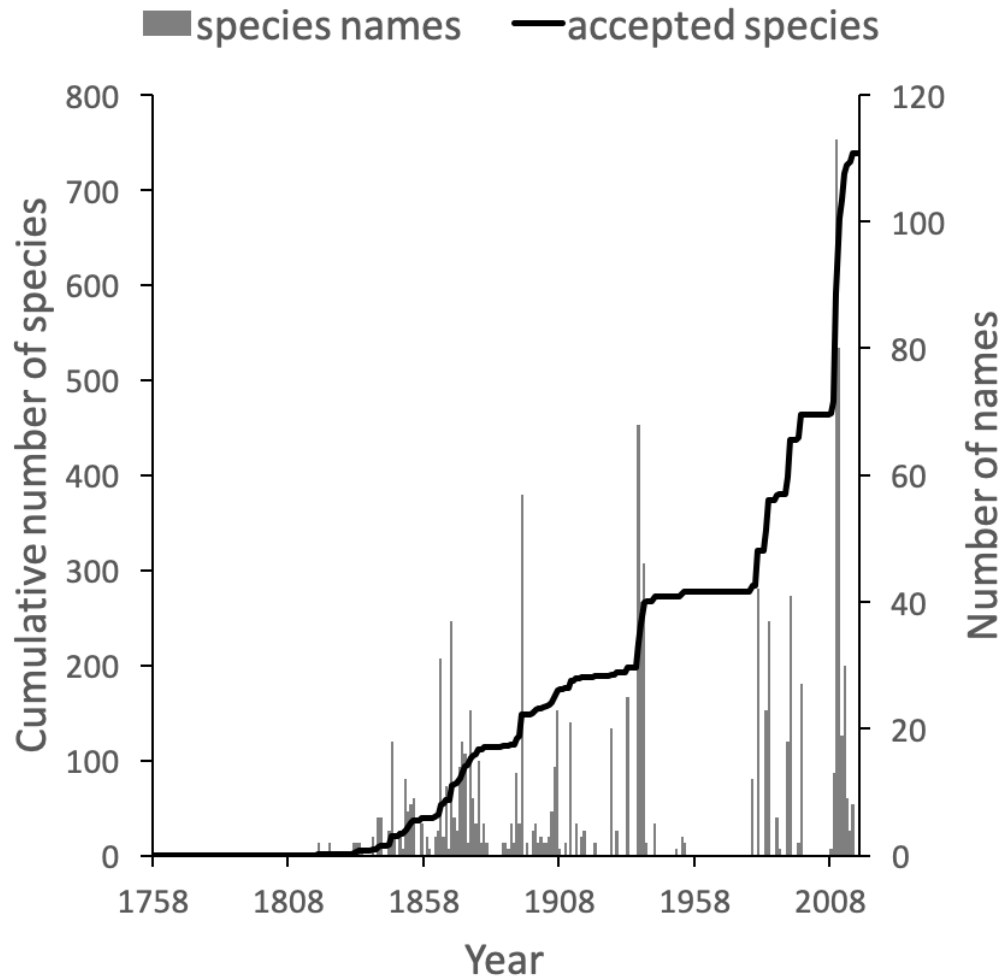
477 The knowledge graph enables exploration of the taxonomic history of any taxon of interest.
478 (Pullen, Jennings & Oberprieler, 2014) recently reviewed the history of weevil taxonomy in
479 Australia. *Ozymandias* has some 3958 accepted weevil species. For each accepted taxon in
480 the ALA classification I used a SPARQL query to retrieve the date the species was originally
481 described, and the dates were then grouped by year. The plot of cumulative numbers of ac-
482 cepted species over time (Figure 12) closely matches that reported by Pullen et al.



483

484 Figure 12. Plot of the history of species discovery for Australian weevils. The solid line is the
485 cumulative number of weevil species that are currently accepted. The vertical bars record the
486 number of new weevil species names published each year. Note the relatively modest increase
487 in names and taxa since the 1930's.
488

489 The same chart also shows the number of weevil species names published each year, in-
490 cluding synonyms. This chart shows that the bulk of weevil discovery took place in the mid-
491 19th to mid-20th centuries. The sharp drop in species discovery since the 1930's may indicate
492 that the bulk of the Australian weevil fauna has been described, but this seems unlikely given
493 that weevils are typically small and cryptic, and many species leaf-litter and other habitats
494 may remain undiscovered (Stork et al., 2008; Riedel & Tänzler, 2016)



495

496 Figure 13. Plot of the history of species discovery for Australian snails in the family Camae-
497 nidae. The solid line is the cumulative number of camaenid species that are currently ac-
498 cepted. The vertical bars record the number of new camaenid species names published each
499 year. In contrast to the weevils (Fig. weevils) new Camaenidae species are continuing to be
500 discovered.
501

502

503 These same queries can be used on other taxonomic groups, enabling us to compare the
504 state of knowledge for different taxa. For example, the land snail family Camaenidae (Figure
505 13) shows a similar pattern of discovery in the mid-19th to mid-20th centuries to that seen in
506 weevils. However, in contrast to weevils these snails have been the subject of ongoing study
507 with over 200 new species being described in the last decade (Köhler, 2010, 2011) a rate of
508 discovery that shows no sign of declining.

509 **Discussion**

510 Building a knowledge graph requires mapping textual representations of entities to identifiers
511 that are shared across data sources (“strings to things”). Creating this mapping is tedious and
512 time consuming to construct, and in a time limited project such as a challenge entry like
513 Ozymandias the mapping is likely to be incomplete before the deadline for the project. De-
514 spite its necessarily incomplete state I think the project illustrates some of the ways a network
515 approach can enrich our knowledge of a topic. The web interface exposes many more connec-
516 tions between taxa, publications and people than are evident in the Atlas of Living Australia
517 and Australian Faunal Directory that were used as source databases.

518

519 The underlying knowledge graph can be used to support queries exploring the history of
520 taxonomic publishing and discovery. Some of these queries could be used to help prioritise
521 future work. For example, the pattern of citations (Figure 11) confirms that the older taxo-
522 nomic literature is still relevant today, reinforcing the case for digitising the legacy taxonomic
523 literature. We could further explore the citation data to prioritise which journals should be
524 scanned first: for example, by focusing on those journals that have been cited the most. Given
525 that the bulk of taxonomic publications in the 20th century appeared in Australian journals,
526 initiatives such as the Biodiversity Heritage Library in Australia would seem well placed to
527 make the case that this work should be scanned and made openly available. Citation counts
528 can also be used more directly. For example, the International Institute for Species Explora-
529 tion annually issues a manually curated list of the “top 10” species discovered the previous
530 year. Such a list could be automatically generated from a knowledge graph using, for exam-
531 ple, the number of citations (or other measures of attention) that each work publishing a new
532 species has received.

533

534 Some analyses of the knowledge graph are more focussed on the state of the knowledge
535 graph itself. For example, querying for author identifiers such as ORCIDs reveals a limited
536 uptake of that identifier. This has implications for proposals to use ORCID as the basis for
537 tracking the broader activities of taxonomists, such as specimen collection and identification
538 (Shorthouse). Perhaps the development of such tools may help raise awareness of the possible
539 benefits of authors registering with ORCID.

540 **Expanding the knowledge graph**

541 The knowledge graph in Ozymandias features only a subset of the entities depicted in earlier
542 work sketching the “biodiversity knowledge graph” (Page, 2013, 2016a). There are several
543 entities that are obvious candidates to be added to Ozymandias, such as specimens and nu-
544 cleotide sequences. However, the number of specimens that could potentially be added has
545 implications for the scalability of the knowledge graph. Bearing this in mind, we could add a
546 subset of specimens, such as type specimens, or those which have been sequenced. Fontaine
547 et al. reported that the average lag time between the discovery of a specimen representing a
548 new species and the description of that species is 21 years. The generality of this observation
549 could be evaluated using a knowledge graph that contains both the taxonomic literature and
550 type specimens with dates of collection.

551 The Biodiversity Literature Repository highlights the potential of treating scientific arti-
552 cles not as monolithic entities but rather as assembles of component parts, including figures.
553 We can drill down further and start to annotate individual parts including fragments of text.
554 The idea of annotating and interlinking fragments of text has a long history, pioneered by
555 people such as Ted Nelson (Douglas R. Dechow & Daniele C. Struppa, 2015), and tools such
556 as Hypotheses.is (“Hypothes.is”) now make this possible. We could view the “micro cita-
557 tions” used by taxonomists to specify the page location of a taxonomic name as a form of an-
558 notation, hence a logical next step is to map these micro citations onto publications in the
559 knowledge graph so that we can locate these micro citations in the context of the taxonomic
560 literature that they refer to.

561 **The future of knowledge graphs**

562 To the extent that Ozymandias is judged to be a success it suggests that knowledge graphs
563 have potential to improve the way we aggregate and interface with biodiversity data. How-
564 ever, it is worth noting that the biodiversity informatics community has been aware of knowl-
565 edge graphs and semantic web technologies for a decade or more, and several taxonomic da-
566 tabases have been serving data in RDF since the mid-2000’s. Yet it is hard to point to suc-
567 cessful applications of these approaches to the study of biodiversity, and there has been lim-
568 ited uptake of linked data beyond a few databases.

569
570 There is a considerable cost involved in cross linking datasets, and to date the rewards for
571 this effort are, perhaps, unclear. At the same time, there is growing concern within biology in
572 general (McDade et al., 2011) and in taxonomy in particular, that existing measures of the
573 output of researchers, such as citations, are poor metrics of activity (cite citation papers).
574 There are also concerns that existing data aggregators do not pay enough attention to tracking
575 the provenance and authorship of information (Franz & Sterner, 2018). Researchers may do
576 much more than write papers, they may clean, prepare, and publish datasets, collect speci-
577 mens, curate collections, identify specimens, etc. Keeping track of these activities is greatly
578 facilitated by the use of stable identifiers for people and the objects they work with (e.g.,
579 specimens, collections, datasets), and a knowledge graph would be an ideal data structure to
580 quantify the work done, and trace the provenance of data and associated annotations. Projects
581 such as Scholia (Nielsen, Mietchen & Willighagen, 2017) already demonstrate the potential of
582 Wikidata to explore the output of scholars. Hence, it may be that the best way to bootstrap the
583 adoption of biodiversity knowledge graphs is to focus on the implications for being able to
584 give appropriate credit to researchers for all the activities that they undertake.

585
586 There is considerable enthusiasm for the potential of identifiers to help evaluate research
587 (Haak, Meadows & Brown, 2018) and yield insights into the behaviour of researchers
588 (Bohannon, 2017). However, the ease with which measures of research activity (such as cita-
589 tion-based measures) switch from being tools for insight into targets to be met suggests we
590 should consider the possibility that metrics developed to create incentives to build knowledge
591 graphs may ultimately harm the researchers being measured.

592

593 Beyond internal drivers, such as documenting the provenance of taxonomic information,
594 and quantifying the contributions of researchers, there are also external drivers for consider-
595 ing knowledge graphs. Wikidata (Vrandečić & Krötzsch, 2014) is an open, global knowledge
596 graph with an enthusiastic community of editors, and many of the entities taxonomists care
597 about are already included in the graph, such as taxa, people, and publications. This means
598 that we can use Wikidata to help define the scope of a knowledge graph. Anyone constructing
599 a knowledge graph rapidly runs into the problem of scope, in other words, when do you stop
600 adding entities? Once we move beyond specialist knowledge in a given field (such as speci-
601 mens, rules of nomenclature, sequences and phylogenies) and include more generic entities
602 that other communities may also be interested in (such as publications, natural history collec-
603 tions, people) we reach the point at which we can stop constructing our graph and defer to
604 Wikidata. Hence a key part of the future development of biodiversity knowledge graphs will
605 be to determine the extent to which Wikidata and its community can be responsible for man-
606 aging biodiversity-related data.

Additional Information and Declarations

Competing Interests

The author declares he has no competing interests.

Author Contributions

Roderic D.M. Page conceived and designed the experiments, performed the experiments, analysed the data, contributed reagents/materials/analysis tools, wrote the paper.

Data Availability

The Ozymandias web site is <https://ozymandias-demo.herokuapp.com>. This site includes a SPARQL interface to query the knowledge graph directly. Source code for the interface is available from GitHub <https://github.com/rdmpage/ozymandias-demo>. Source code for the scripts used to harvest and clean the data used to populate the knowledge graph is available from <https://github.com/rdmpage/oz-afd-export>, <https://github.com/rdmpage/oz-ala-harvest>, <https://github.com/rdmpage/oz-csl>, and <https://github.com/rdmpage/oz-wikispecies>.

Funding

The work described here was an entry in the Global Biodiversity Information Facility 2018 Ebbe Nielsen Challenge. GBIF had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Acknowledgements

Constructing the knowledge graph described here would have been impossible without the wealth of freely available and open source software used in the project. Furthermore, it should be obvious that none of this would have been possible without the centuries of taxonomic research by generations of researchers, and the recent efforts to make that research digitally accessible via projects such as Atlas of Living Australia and the Australian Faunal Directory. I'm also indebted to GBIF for running the 2018 Ebbe Nielsen Challenge which gave me a hard deadline to work towards. I also thank Steve Baskauf and Joel Sachs for feedback on the project, and for invitations to present Ozymandias to their colleagues.

References

- Altmetric. Available at <https://www.altmetric.com/> (accessed November 28, 2018).
- Atlas of Living Australia. Available at <https://www.ala.org.au/> (accessed November 27, 2018).
- Australian Faunal Directory. Available at <https://biodiversity.org.au/afd/home> (accessed November 27, 2018).
- Bebber DP, Wood JRI, Barker C, Scotland RW. 2013. Author inflation masks global capacity for species discovery in flowering plants. *New Phytologist* 201:700–706. DOI: 10.1111/nph.12522.
- Biodiversity Literature Repository. Available at <http://plazi.org/resources/bibliography-of-life-bol/biodiversity-literature-repository-blr/> (accessed November 28, 2018).
- Bohannon J. 2017. Vast set of public CVs reveals the world's most migratory scientists. *Science*. DOI: 10.1126/science.aal1189.
- Bollacker K, Evans C, Paritosh P, Sturge T, Taylor J. 2008. Freebase: A Collaboratively Created Graph Database for Structuring Human Knowledge. In: *Proceedings of the 2008 ACM SIGMOD International Conference on Management of Data*. SIGMOD '08. New York, NY, USA: ACM, 1247–1250. DOI: 10.1145/1376616.1376746.
- Costello MJ, Wilson S, Houlding B. 2013. More Taxonomists Describing Significantly Fewer Species per Unit Effort May Indicate That Most Species Have Been Discovered. *Systematic Biology* 62:616–624. DOI: 10.1093/sysbio/syt024.
- CrossRef. Available at <https://www.crossref.org/> (accessed November 27, 2018).
- Douglas R, Dechow, Daniele C, Struppa (eds.). 2015. Intertwined: The work and influence of Ted Nelson. *History of Computing*. DOI: 10.1007/978-3-319-16925-5.
- Egloff W, Agosti D, Kishor P, Patterson D, Miller J. 2017. Copyright and the Use of Images as Biodiversity Data. *Research Ideas and Outcomes* 3:e12502. DOI: 10.3897/rio.3.e12502.
- Franz NM, Sterner BW. 2018. To increase trust, change the social design behind aggregated biodiversity data. *Database* 2018. DOI: 10.1093/database/bax100.
- Godfray HCJ. 2007. Linnaeus in the information age. *Nature* 446:259–260. DOI: 10.1038/446259a.

- Grieneisen ML, Zhan Y, Potter D, Zhang M. 2014. Biodiversity, Taxonomic Infrastructure, International Collaboration, and New Species Discovery. *BioScience* 64:322–332. DOI: 10.1093/biosci/biu035.
- Haak LL, Meadows A, Brown J. 2018. Using ORCID, DOI, and Other Open Identifiers in Research Evaluation. *Frontiers in Research Metrics and Analytics* 3. DOI: 10.3389/frma.2018.00028.
- Hypothes.is. Available at <https://web.hypothes.is/> (accessed November 28, 2018).
- Joppa LN, Roberts DL, Pimm SL. 2011. The population ecology and social behaviour of taxonomists. *Trends in Ecology & Evolution* 26:551–553. DOI: 10.1016/j.tree.2011.07.010.
- Kennedy J, Hyam R, Kukla R, Paterson T. 2006. Standard Data Model Representation for Taxonomic Information. *OMICS: A Journal of Integrative Biology* 10:220–230. DOI: 10.1089/omi.2006.10.220.
- Köhler F. 2010. Uncovering local endemism in the Kimberley, Western Australia: description of new species of the genus *Amplirhagada* Iredale, 1933 (Pulmonata: Camaenidae). *Records of the Australian Museum* 62:217–284. DOI: 10.3853/j.0067-1975.62.2010.1554.
- Köhler F. 2011. Descriptions of new species of the diverse and endemic land snail *Amplirhagada* Iredale, 1933 from rainforest patches across the Kimberley, Western Australia (Pulmonata: Camaenidae). *Records of the Australian Museum* 63:167–202. DOI: 10.3853/j.0067-1975.63.2011.1581.
- May RM. 1988. How Many Species Are There on Earth? *Science* 241:1441–1449. DOI: 10.1126/science.241.4872.1441.
- McDade LA, Maddison DR, Guralnick R, Piwowar HA, Jameson ML, Helgen KM, Herendeen PS, Hill A, Vis ML. 2011. Biology Needs a Modern Assessment System for Professional Productivity. *BioScience* 61:619–625. DOI: 10.1525/bio.2011.61.8.8.
- Michel F, Gargominy O, Terceire S, Faron Zucker C. 2017. A Model to Represent Nomenclatural and Taxonomic Information as Linked Data. Application to the French Taxonomic Register, TAXREF. In: *S4Biodiv 2017 - 2nd International Workshop on Semantics for Biodiversity co-located with ISWC 2017*. Vienna, Austria, 1–12.
- Nakabo T. 1982. REVISION OF GENERA OF THE DRAGONETS (PISCES: CALLIONYMIDAE). *Publications of the Seto Marine Biological Laboratory* 27:77–131. DOI: 10.5134/176044.
- Nielsen FÅ, Mietchen D, Willighagen E. 2017. Scholia, Scientometrics and Wikidata. *Lecture Notes in Computer Science*:237–259. DOI: 10.1007/978-3-319-70407-4_36.
- ORCID. Available at <https://orcid.org/> (accessed November 27, 2018).
- Page RD. 2011. Extracting scientific articles from a large digital archive: BioStor and the Biodiversity Heritage Library. *BMC Bioinformatics* 12:187. DOI: 10.1186/1471-2105-12-187.
- Page RDM. 2013. BioNames: linking taxonomy, texts, and trees. *PeerJ* 1:e190. DOI: 10.7717/peerj.190.
- Page RDM. 2016a. Towards a biodiversity knowledge graph. *Research Ideas and Outcomes* 2:e8767. DOI: 10.3897/rio.2.e8767.

- Page RDM. 2016b. DNA barcoding and taxonomy: dark taxa and dark texts. *Phil. Trans. R. Soc. B* 371:20150334. DOI: 10.1098/rstb.2015.0334.
- Pullen KR, Jennings D, Oberprieler RG. 2014. Annotated catalogue of Australian weevils (Coleoptera: Curculionoidea). *Zootaxa* 3896:1. DOI: 10.11646/zootaxa.3896.1.1.
- Quan DA, Karger R. 2004. How to make a semantic web browser. In: ACM Press,. DOI: 10.1145/988672.988707.
- Riedel A, Tänzler R. 2016. Revision of the Australian species of the weevil genus *Trigonopterus* Fauvel. *ZooKeys* 556:97–162. DOI: 10.3897/zookeys.556.6126.
- Sangster G, Luksenburg JA. 2014. Declining Rates of Species Described per Taxonomist: Slowdown of Progress or a Side-effect of Improved Quality in Taxonomy? *Systematic Biology* 64:144–151. DOI: 10.1093/sysbio/syu069.
- Schema.org. Available at <https://schema.org/> (accessed December 3, 2018).
- Senderov V, Simov K, Franz N, Stoev P, Catapano T, Agosti D, Sautter G, Morris RA, Penev L. 2018. OpenBiodiv-O: ontology of the OpenBiodiv knowledge management system. *Journal of Biomedical Semantics* 9. DOI: 10.1186/s13326-017-0174-5.
- Shorthouse DP. Bloodhound. Available at <https://bloodhound.shorthouse.net> (accessed November 28, 2018).
- de Solla Price DJ. 1965. Networks of Scientific Papers. *Science* 149:510–515. DOI: 10.1126/science.149.3683.510.
- Stork NE, Grimbacher PS, Storey R, Oberprieler RG, Reid C, Slipinski SA. 2008. What determines whether a species of insect is described? Evidence from a study of tropical forest beetles. *Insect Conservation and Diversity* 1:114–119. DOI: 10.1111/j.1752-4598.2008.00016.x.
- Tancoigne E, Ollivier G. 2017. Evaluating the progress and needs of taxonomy since the Convention on Biological Diversity: going beyond the rate of species description. *Australian Systematic Botany* 30:326. DOI: 10.1071/sb16017.
- Verstak A, Acharya A, Suzuki H, Henderson S, Iakhiaev M, Lin CCY, Shetty N. 2014. On the Shoulders of Giants: The Growing Impact of Older Articles. *arXiv:1411.0275 [cs]*.
- Vicki Tardif Holland, Jason Johnson. 2014. Introducing “Role.” Available at <http://blog.schema.org/2014/06/introducing-role.html> (accessed November 29, 2018).
- Vrandečić D, Krötzsch M. 2014. Wikidata. *Communications of the ACM* 57:78–85. DOI: 10.1145/2629489.
- W3C SPARQL Working Group. 2013. SPARQL 1.1 Overview. Available at <https://www.w3.org/TR/sparql11-overview/> (accessed November 27, 2018).
- Wieczorek J, Bloom D, Guralnick R, Blum S, Döring M, Giovanni R, Robertson T, Vieglais D. 2012. Darwin Core: An Evolving Community-Developed Biodiversity Data Standard. *PLoS ONE* 7:e29715. DOI: 10.1371/journal.pone.0029715.

Supplementary Information

Publications and identifiers

Get count of number of published works for each year, and number of works with identifiers.

```
PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
SELECT ?work_date (COUNT(?w) as ?c) (COUNT(?doi) as ?c_doi)
(COUNT(?biostor) as ?c_biostor) (COUNT(?jstor) as ?c_jstor)
(COUNT(?pdf) as ?c_pdf)
WHERE
{
  ?w <http://schema.org/datePublished> ?work_date .

  # just articles
  ?w <http://www.w3.org/1999/02/22-rdf-syntax-ns#type>
<http://schema.org/ScholarlyArticle> .

  # DOI?
  OPTIONAL {
    ?w <http://schema.org/identifier> ?doi .
    ?doi <http://schema.org/propertyID> "doi" .
  }

  # BioStor?
  OPTIONAL {
    ?w <http://schema.org/identifier> ?biostor .
    ?biostor <http://schema.org/propertyID> "biostor" .
  }

  # JSTOR?
  OPTIONAL {
    ?w <http://schema.org/identifier> ?jstor .
    ?jstor <http://schema.org/propertyID> "jstor" .
  }

  # PDF?
  OPTIONAL {
    ?w <http://schema.org/encoding> ?pdf .
  }
}
```



```
?pdf <http://schema.org/fileFormat> "application/pdf" .
}

FILTER regex(?work_date, "^[0-9]{4}$")

#FILTER (xsd:integer(?work_date) > 1980)
}
GROUP BY ?work_date
ORDER BY ?work_date
```

Data in publications.tsv

Journal ranks

Query to retrieve top 10 journals for a given decade (in this case 1910)

```
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
PREFIX tc: <http://rs.tdwg.org/ontology/voc/TaxonConcept#>
SELECT ?journal ?issn (COUNT(?journal) AS ?count) WHERE
{
  ?work <http://www.w3.org/1999/02/22-rdf-syntax-ns#type>
  <http://schema.org/ScholarlyArticle> .
  ?work <http://schema.org/isPartOf> ?container .
  ?container <http://schema.org/name> ?journal .
  ?work <http://schema.org/datePublished> ?year .

  OPTIONAL {
    ?container <http://schema.org/issn> ?issn .
  }
  FILTER ((xsd:integer(?year) >= 1910) && (xsd:integer(?year)
  < " . ($year + 9) . "))
}
GROUP BY ?journal ?issn
ORDER BY DESC(?count)
LIMIT 10
```

Repeat this query for all decades, aggregate results, then filter for journals with > 200 articles.

Data in journals.tsv

Citation patterns

Find all pairs of citing articles and get dates they were published.

```
PREFIX xsd: <http://www.w3.org/2001/XMLSchema#>
SELECT ?cited_identifier_type (xsd:integer(?w_year) as ?from)
(xsd:integer(?work_year) as ?to)
WHERE
{

?w <http://schema.org/identifier> ?identifier .
  ?w <http://schema.org/name> ?w_name .
?w <http://schema.org/datePublished> ?w_year .
# Identifier (e.g., DOI) for work we are displaying
?identifier <http://schema.org/value> ?identifier_value .

?citing_identifier <http://schema.org/value> ?identifier_value
.
?citing <http://schema.org/identifier> ?citing_identifier .

# What does this work cite (typically from CrossRef data)
?citing <http://schema.org/citation> ?cited .

# Translate the citing work\'s DOI (or other identifier) into
AFD identifier
# Get identifier (typically a DOI) for citing work
?cited <http://schema.org/identifier> ?cited_identifier .
?cited_identifier <http://schema.org/value>
?cited_identifier_value .
?cited_identifier <http://schema.org/propertyID>
?cited_identifier_type .

# Get work(s) with this identifier (may have > 1 if we have
CrossRef record in our triple store
?work_identifier <http://schema.org/value>
?cited_identifier_value .
?work <http://schema.org/identifier> ?work_identifier .
?work <http://schema.org/name> ?name .
?work <http://schema.org/datePublished> ?work_year .

# Just include citing records that are also in ALA
```

```
FILTER regex(str(?work), \'biodiversity.org.au\') .  
FILTER regex(str(?w), \'biodiversity.org.au\') .
```

```
FILTER regex(?w_year, "^[0-9]{4}$")  
FILTER regex(?work_year, "^[0-9]{4}$")  
}
```

Data in cites.tsv

Weevils

Number of accepted taxon names per year.

```
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>  
SELECT ?year (COUNT(?taxonName) AS ?count)  
WHERE  
{  
VALUES ?root_name {"CURCULIONOIDEA"}  
?root <http://schema.org/name> ?root_name .  
?child rdfs:subClassOf+ ?root .  
?child rdfs:subClassOf ?parent .  
?child <http://schema.org/name> ?child_name .  
?parent <http://schema.org/name> ?parent_name .  
  
  ?child  
<http://taxref.mnhn.fr/lod/property/hasReferenceName> ?taxon-  
Name .  
  
  ?taxonName  
<http://rs.tdwg.org/ontology/voc/TaxonName#rankString> "spe-  
cies" .  
  ?taxonName <http://rs.tdwg.org/ontology/voc/TaxonName#year>  
?year .  
}  
GROUP BY ?year  
ORDER BY ?year
```

Sum these to generate cumulative total.

Number of weevil names published each year:

```
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
SELECT ?year (COUNT(DISTINCT ?name) AS ?c)
WHERE
{
VALUES ?root_name {"CURCULIONOIDEA"}
?root <http://schema.org/name> ?root_name .
?child rdfs:subClassOf+ ?root .
?child rdfs:subClassOf ?parent .
?child <http://schema.org/name> ?child_name .
?parent <http://schema.org/name> ?parent_name .

    ?child
<http://taxref.mnhn.fr/lod/property/hasReferenceName>|<http://
taxref.mnhn.fr/lod/property/hasSynonym> ?taxonName .
    ?taxonName
<http://rs.tdwg.org/ontology/voc/TaxonName#rankString> "spe-
cies" .
    ?taxonName <http://schema.org/name> ?name .
    ?taxonName <http://rs.tdwg.org/ontology/voc/TaxonName#year>
?year .
}

GROUP BY ?year
ORDER BY ?year
```

Combined data in weevils.tsv

Snails

Number of accepted taxon names per year

```
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
SELECT ?year (COUNT(?taxonName) AS ?count)
WHERE
{
VALUES ?root_name {"CAMAENIDAE"}
?root <http://schema.org/name> ?root_name .
?child rdfs:subClassOf+ ?root .
?child rdfs:subClassOf ?parent .
?child <http://schema.org/name> ?child_name .
?parent <http://schema.org/name> ?parent_name .
```

```
    ?child
    <http://taxref.mnhn.fr/lod/property/hasReferenceName> ?taxon-
    Name .

    ?taxonName
    <http://rs.tdwg.org/ontology/voc/TaxonName#rankString> "spe-
    cies" .
    ?taxonName <http://rs.tdwg.org/ontology/voc/TaxonName#year>
    ?year .

}
GROUP BY ?year
ORDER BY ?year
```

Sum these to generate cumulative total.

Number of snail names published each year:

```
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
SELECT ?year (COUNT(DISTINCT ?name) AS ?c)
WHERE
{
VALUES ?root_name {"CAMAENIDAE"}
?root <http://schema.org/name> ?root_name .
?child rdfs:subClassOf+ ?root .
?child rdfs:subClassOf ?parent .
?child <http://schema.org/name> ?child_name .
?parent <http://schema.org/name> ?parent_name .

    ?child
    <http://taxref.mnhn.fr/lod/property/hasReferenceName>|<http://
    taxref.mnhn.fr/lod/property/hasSynonym> ?taxonName .
    ?taxonName
    <http://rs.tdwg.org/ontology/voc/TaxonName#rankString> "spe-
    cies" .
    ?taxonName <http://schema.org/name> ?name .
    ?taxonName <http://rs.tdwg.org/ontology/voc/TaxonName#year>
    ?year .

}

GROUP BY ?year
ORDER BY ?year
```

Combined data in snails.tsv

Authors and ORCIDs

How many authors of works with DOIs post 2011?

```
SELECT (COUNT(DISTINCT ?creator) as ?c)
WHERE
{
  GRAPH <https://biodiversity.org.au/afd/publication> {

    ?work <http://schema.org/identifier> ?identifier .
    ?identifier <http://schema.org/propertyID> "doi" .
    ?identifier <http://schema.org/value> ?doi .

    ?work <http://schema.org/datePublished> ?datePublished .

    ?work <http://schema.org/creator> ?role .
    ?role <http://schema.org/roleName> ?roleName .
    ?role <http://schema.org/creator> ?creator .
    ?creator <http://schema.org/name> ?name .
  }

  FILTER (xsd:integer(?datePublished) > 2011)
}
```

How many authors of works with DOIs post 2011 had an ORCID?

```
SELECT DISTINCT ?orcid_creator
WHERE
{
  GRAPH <https://biodiversity.org.au/afd/publication> {

    ?work <http://schema.org/identifier> ?identifier .
    ?identifier <http://schema.org/propertyID> "doi" .
    ?identifier <http://schema.org/value> ?doi .

    ?work <http://schema.org/datePublished> ?datePublished .
```

```
?work <http://schema.org/creator> ?role .
?role <http://schema.org/roleName> ?roleName .
?role <http://schema.org/creator> ?creator .
?creator <http://schema.org/name> ?name .
}

GRAPH <https://orcid.org>
{
  ?orcid_identifier <http://schema.org/value> ?doi .
  ?orcid_work <http://schema.org/identifier> ?or-
cid_identifier .

  ?orcid_work <http://schema.org/creator> ?orcid_role .
  ?orcid_role <http://schema.org/roleName> ?orcid_roleName
.

  ?orcid_role <http://schema.org/creator> ?orcid_creator .

  ?orcid_creator <http://schema.org/name> ?orcid_name .
}

FILTER(?roleName = ?orcid_roleName)
FILTER (xsd:integer(?datePublished) > 2011)
}
```