

**A peer-reviewed version of this preprint was published in PeerJ on 17 September 2018.**

[View the peer-reviewed version](https://peerj.com/articles/cs-164) (peerj.com/articles/cs-164), which is the preferred citable publication unless you specifically need to cite this preprint.

Thessen AE, Poelen JH, Collins M, Hammock J. 2018. 20 GB in 10 minutes: a case for linking major biodiversity databases using an open socio-technical infrastructure and a pragmatic, cross-institutional collaboration. PeerJ Computer Science 4:e164  
<https://doi.org/10.7717/peerj-cs.164>

# 20 GB in 10 minutes: A case for linking major biodiversity databases using an open socio-technical infrastructure and a pragmatic, cross-institutional collaboration

Anne E Thessen <sup>Corresp., 1</sup>, Jorrit H Poelen <sup>2</sup>, Matthew Collins <sup>3</sup>, Jen Hammock <sup>4</sup>

<sup>1</sup> Ronin Institute for Independent Scholarship, Montclair, New Jersey, USA

<sup>2</sup> Self-employed, Oakland, California, USA

<sup>3</sup> University of Florida, Gainesville, Florida, USA

<sup>4</sup> National Museum of Natural History, Washington DC, USA

Corresponding Author: Anne E Thessen

Email address: annethessen@gmail.com

Biodiversity information is made available through numerous databases that each have their own data models, web services, and data types. Combining data across databases leads to new insights, but is not easy because each database uses its own system of identifiers. In the absence of stable and interoperable identifiers, databases are often linked using taxonomic names. This labor intensive, error prone, and lengthy process relies on accessible versions of nomenclatural authorities and fuzzy-matching algorithms.

To approach the challenge of linking diverse data, more than technology is needed. New social collaborations like the Global Unified Open Data Architecture (GUODA) that combine skills from diverse groups of computer engineers from iDigBio, server resources from the Advanced Computing and Information Systems (ACIS) Lab, global-scale data presentation from EOL, and independent developers and researchers are what is needed to make concrete progress on finding relationships between biodiversity datasets.

This paper will discuss a technical solution developed by the GUODA collaboration for faster linking across databases with a use case linking Wikidata and the Global Biodiversity Interactions database (GloBI). The GUODA infrastructure is a 12-node, high performance computing cluster made up of about 192 threads with 12 TB of storage and 288 GB memory. Using GUODA, 20GB of compressed JSON from Wikidata was processed and linked to GloBI in about 10-11 minutes. Instead of comparing name strings or relying on a single identifier, Wikidata and GloBI were linked by comparing graphs of biodiversity identifiers external to each system. This method resulted in adding 119,957 Wikidata links in GloBI, an increase of 13.7% of all outgoing name links in GloBI. Wikidata and GloBI were compared to Open Tree Taxonomy to examine consistency and coverage. The process of parsing Wikidata, Open Tree Taxonomy and GloBI archives and calculating consistency metrics was done in minutes on the GUODA platform. As a model collaboration, GUODA has the potential to revolutionize biodiversity science by bringing diverse technically minded people together with high performance computing resources that are accessible from a laptop or desktop. However, participating in such a collaboration still requires basic programming skills.

1 20 GB in 10 minutes: A case for linking  
2 major biodiversity databases using an  
3 open socio-technical infrastructure and  
4 a pragmatic, cross-institutional  
5 collaboration

6  
7 Anne E Thessen<sup>1\*</sup>, Jorrit H. Poelen<sup>2</sup>, Matthew Collins<sup>3</sup>, Jen Hammock<sup>4</sup>

8  
9 <sup>1</sup>The Ronin Institute for Independent Scholarship, Montclair, NJ 07043, USA

10 <sup>2</sup>400 Perkins Street, Apt. 104, Oakland, CA 94610, USA

11 <sup>3</sup>University of Florida, 968 Center Drive, Gainesville, FL, 32611, USA

12 <sup>4</sup>National Museum of Natural History, 10th and Constitution Ave., Washington, DC

13 \*Corresponding author [annethessen@gmail.com](mailto:annethessen@gmail.com)

14

## 15 Abstract

16 Biodiversity information is made available through numerous databases that each have their own  
17 data models, web services, and data types. Combining data across databases leads to new  
18 insights, but is not easy because each database uses its own system of identifiers. In the absence  
19 of stable and interoperable identifiers, databases are often linked using taxonomic names. This  
20 labor intensive, error prone, and lengthy process relies on accessible versions of nomenclatural  
21 authorities and fuzzy-matching algorithms.

22

23 To approach the challenge of linking diverse data, more than technology is needed. New social  
24 collaborations like the Global Unified Open Data Architecture (GUODA) that combine skills  
25 from diverse groups of computer engineers from iDigBio, server resources from the Advanced  
26 Computing and Information Systems (ACIS) Lab, global-scale data presentation from EOL, and  
27 independent developers and researchers are what is needed to make concrete progress on finding  
28 relationships between biodiversity datasets.

29

30 This paper will discuss a technical solution developed by the GUODA collaboration for faster  
31 linking across databases with a use case linking Wikidata and the Global Biodiversity  
32 Interactions database (GloBI). The GUODA infrastructure is a 12-node, high performance  
33 computing cluster made up of about 192 threads with 12 TB of storage and 288 GB memory.  
34 Using GUODA, 20GB of compressed JSON from Wikidata was processed and linked to GloBI  
35 in about 10-11 minutes. Instead of comparing name strings or relying on a single identifier,  
36 Wikidata and GloBI were linked by comparing graphs of biodiversity identifiers external to each  
37 system. This method resulted in adding 119,957 Wikidata links in GloBI, an increase of 13.7%  
38 of all outgoing name links in GloBI. Wikidata and GloBI were compared to Open Tree  
39 Taxonomy to examine consistency and coverage. The process of parsing Wikidata, Open Tree  
40 Taxonomy and GloBI archives and calculating consistency metrics was done in minutes on the  
41 GUODA platform. As a model collaboration, GUODA has the potential to revolutionize  
42 biodiversity science by bringing diverse technically minded people together with high  
43 performance computing resources that are accessible from a laptop or desktop. However,  
44 participating in such a collaboration still requires basic programming skills.

45

## 46 Introduction

47 Biodiversity databases provide global access to information about species via the Web. These  
48 databases contain information as varied as observation records, text descriptions, images, maps,  
49 genetic sequences, phylogenetic trees, and trait data (Table 1). All of these data become much  
50 more useful if they can be linked. Many biodiversity databases share information with each other  
51 (Bingham et al. 2017), but creating the links can be very difficult for several reasons including  
52 the size of the databases, the heterogeneous nature of the data, and the heterogeneous nature of  
53 the identifiers used by the different resources (Page 2008).

54

55 The more popular methods for linking biodiversity databases include taxonomic names, lsid, and  
56 doi. The Encyclopedia of Life uses taxonomic names to automatically aggregate data from  
57 hundreds of providers (Parr et al. 2014). BioNames links data using lsid, doi, handles,  
58 bibliographic citations, and taxonomic names (Page 2013). The iPhylo LinkOut service mapped  
59 identifiers used by the NCBI taxonomy database (which provides the taxonomic backbone for  
60 GenBank) to Wikipedia pages using taxonomic names, including synonyms (Page 2011).  
61 TBMMap provides links from TreeBase across several taxonomic databases, such as ITIS and  
62 NCBI (Page 2007). This mapping was also achieved using taxonomic names, but in some cases  
63 GenBank Accession numbers and museum specimen codes were available for supplement. The  
64 use of taxonomic names to aggregate data can lead to errors and requires significant *a priori*  
65 knowledge either in the form of curators or an authoritative nomenclature.

66  
67 Many databases expose their own internal identifiers, such as the WoRMS Aphia ID, so others  
68 can link their data to those resources within their own systems, often by providing a URL.  
69 Databases like WoRMS provide web services that allow users to look up an identifier for a taxon  
70 in question, one at a time. While this makes linking easier, it is still difficult to scale across all  
71 databases. For example, a list of all the taxon identifiers in EOL is 300 MB compressed. No  
72 system of identifiers is universal across biodiversity databases and none of them are easy to  
73 implement at scale.

74  
75 While the data would be much more useful if linked, there is a lack of tools for linking data  
76 across databases at scale. Most mappings are done at great expense and then are made available  
77 as a separate file or incorporated into the resources themselves. LinkOut, BioNames, GBIF, and  
78 EOL take more than a day to link across their entire body of aggregated content. This paper  
79 discusses links made between GloBI and Wikidata (WD) in 10 minutes using GUODA, a high  
80 performance computing system available for analysis of large biodiversity data sets.

## 81 Methods

### 82 Description of Resources

#### 83 GUODA

84 Following an iDigBio hack-a-thon in June 2015, GUODA was created as a pragmatic way to  
85 compute over multiple large biodiversity databases in a mutually beneficial collaboration  
86 between iDigBio, EOL, Kew Garden, and independent developers. Catalyzed by various  
87 presentations at conferences, hardware provided by ACIS, 20+ meetings, and several prototypes  
88 (e.g., <http://effechecka.org>, <https://gimmefreshdata.github.io>), a general access biodiversity data  
89 integration and analysis environment was created. This environment, with the aggregated  
90 experience and perspectives of all the collaborators, was used to produce the results of this paper.

91

92 Housed at the ACIS Lab at the University of Florida, the GUODA infrastructure consists of 12  
93 IBM HS22 blades each with 8 cores, 24 GB of memory, and 1 TB of storage each. This makes a  
94 total of 192 threads, 288 GB of memory and 12 TB of disk space available for processing jobs  
95 using Apache Spark (Fig 1; Zaharia et al. 2016). The cluster is managed under Apache Mesos  
96 (Hindman et al. 2011) which is a distributed scheduling system for periodic jobs. For long  
97 running processes, such as web APIs or databases, the Marathon  
98 (<https://github.com/mesosphere/marathon>) framework is run within Mesos. Marathon facilitates  
99 running always-up services with monitoring, automatic deployment of code, re-scaling to  
100 multiple nodes, and other management features. Mesos is responsible for accepting requests to  
101 start Spark frameworks, processes which do the actual computation and may span multiple  
102 servers, and allocation of resources requested by the framework.

103

104 Hadoop HDFS (Shvachko et al. 2010) is installed outside of Mesos directly on all 12 nodes of  
105 the cluster and provides redundant parallel shared storage to all nodes as well as the Jupyter  
106 notebook (Kluyver et al. 2016) server that provides a programming interface to end users. Each  
107 node has 1 TB of local disk storage for a total of about 3.5 TB of usable storage space for data  
108 files in Apache Parquet format. Spark is aware of the placement of data on an HDFS cluster and  
109 will divide processing among nodes in a way that prefers to read and write data that is local to  
110 the node to minimize network traffic.

## 111 Wikidata

112 Wikidata (WD) is a free and open knowledge base that provides structured data for Wikimedia  
113 projects ([www.wikidata.org](http://www.wikidata.org); Vrandečić & Krötzsch 2014). Similar to Wikipedia, anyone can  
114 read or edit the resource. Information, including links to other resources, can be added to  
115 Wikidata using bots and batch imports through their Data Import Hub  
116 ([https://www.wikidata.org/wiki/Wikidata:Data\\_Import\\_Hub](https://www.wikidata.org/wiki/Wikidata:Data_Import_Hub)). Wikidata information about taxa  
117 can be conceptualized as a graph linking related taxa to each other and identifiers from other  
118 databases to the taxa they represent (Fig 2). Every taxon in Wikidata is issued a Wikidata  
119 identifier. While a public Wikidata SPARQL endpoint and associated tools (Voß 2016) exist,  
120 these APIs are not suitable for batch processing. For example, when attempting to retrieve all  
121 taxa using the public SPARQL endpoint, a query timeout error was reported. In addition, the  
122 APIs are expected to return different results over time, so reproducing results is difficult if not  
123 impossible. This is why we used a json archive to access Wikidata (Wikidata 2018).

## 124 GloBI

125 GloBI is a database of biotic interactions recorded as Organism\_1:has\_relationship:Organism\_2  
126 (Poelen et al. 2014). GloBI uses a combination of web APIs, taxon archives, and name  
127 correction/parsing methods in an attempt to link names from species interaction datasets to  
128 existing sources. Spatial, temporal, and taxonomic coverage in GloBI is sparse and unevenly  
129 distributed (see Eltonian shortfall, Hortal et al. 2015), with spatial concentrations in Europe and  
130 North America and taxonomically concentrated in Arthropods, Fungi, and Plants. Only 8% of

131 taxa in ITIS are also in GloBI. A detailed technical description of the GloBI data model and  
132 services has been published elsewhere (Poelen et al. 2014). GloBI maintains a graph of related  
133 taxa and their identifiers from different databases (Poelen et al. 2014). GloBI does not introduce  
134 its own taxon ids. Instead, it records how names were mapped from a source name into an  
135 external taxonomic database using a taxon graph (see  
136 <https://globalbioticinteractions.org/references>). We used GloBI Taxon Graph v0.4.2 (Poelen  
137 2018b).

## 138 Open Tree Taxonomy

139 To assess taxonomic id coverage, the taxa in Wikidata and GloBI were compared to Open Tree  
140 Taxonomy (OTT 3.0; <http://files.opentreeoflife.org/ott/ott3.0/ott3.0.tgz>; Rees & Cranston 2017).  
141 OTT was built using an automated algorithm with informed choices to aggregate and link  
142 existing naming authorities into a reasonably comprehensive, artificial, taxonomy. OTT contains  
143 4,385,000 external links for 3,594,550 taxa aggregated and linked over 5 authorities (i.e., GBIF,  
144 IF, SILVA, WoRMS, NCBI).  
145

## 146 Linking Wikidata And GloBI

147 Both Wikidata and GloBI have taxon graphs that map to identifiers from external databases (e.g.  
148 NCBI, ITIS, GBIF, EOL, Index Fungorum (IF), Fishbase and WoRMS). A Wikidata dump was  
149 loaded into GUODA and processed to extract taxon items (about 2.3 million) and their links to  
150 NCBI, ITIS, GBIF, EOL, IF, Fishbase and WoRMS. This was the Wikidata taxon graph. This  
151 taxon graph was loaded into a lookup table where each row contained an NCBI, ITIS, GBIF,  
152 EOL, IF, Fishbase or WoRMS identifier and the corresponding Wikidata identifier. The GloBI  
153 taxon graph was already in a similarly formatted lookup table. The taxon graphs in GloBI and  
154 Wikidata were mapped to each other with a join of the NCBI, ITIS, GBIF, EOL, IF, Fishbase or  
155 WoRMS identifiers of the respective lookup tables. So, for each external identifier that occurred  
156 in both Wikidata and GloBI, the corresponding Wikidata identifier inserted in the GloBI lookup  
157 table. For instance, consider Wikidata taxon item Q140 (<https://www.wikidata.org/wiki/Q140>  
158 accessed on 30 March 2018; *Panthera leo*) points to ITIS:183803. With the matching algorithm  
159 used, GloBI now considers WD:Q140 to be linked to all taxon entries that are considered the  
160 same as, or synonymous to, ITIS:183803.  
161

162 This final joined graph was saved into HDFS as a Parquet file and linked entries were appended  
163 to GloBI Taxon Graph from v0.3.0 onward (Poelen 2018c). In addition, the GloBI ingestion  
164 engine was updated to automatically perform the taxon graph matching for future updates. This  
165 linkage enabled lookups of diet items of lions by Wikidata identifier via  
166 <https://www.globalbioticinteractions.org/?interactionType=eats&sourceTaxon=WD%3AQ140>  
167 and facilitates future integration of species interaction data with Wikidata.



## 168 Taxon Graph Overlap and Consistency

169 OTT, Wikidata, and GloBI taxon graphs maintain links to GBIF, IF, NCBI and WoRMS  
170 identifiers (referred to as external identifiers). The taxon graphs are considered to (partially)  
171 overlap if individual taxon ids from different graphs have at least one external identifier in  
172 common. In addition, a taxon graph is inconsistent if a taxon id links to multiple external  
173 identifiers from the same identifier scheme. Similarly, overlapping taxon ids are said to be  
174 inconsistent if they link to multiple external identifiers from the same identifier scheme. Where  
175 overlap is a measure for taxon graph similarity, consistency can be seen as a way to measure the  
176 relative quality of (overlapping) taxon graphs.

177

178 For instance, let's say that OTT:1087695 is linked to NCBI:191633, WoRMS:156905, and  
179 GBIF:1449280. In addition, WD:Q7247420 (<https://www.wikidata.org/wiki/Q7247420>) points to  
180 WORMS:156905, GBIF:1449280, and NCBI:191633. This would mean that links of these OTT  
181 and WD ids overlap and are consistent, because they do not point to different names in same  
182 naming schemes. However, when considering the GloBI taxon "id" "GLOBI:null@Procladius  
183 sp1 M\_PL\_014", multiple links to external ids were found (e.g., NCBI:1981571, NCBI:1981569,  
184 NCBI:1981572, NCBI:1981573, NCBI:1981574, NCBI:1981570). In this case, the GloBI taxon  
185 id is inconsistent.

## 186 Data Access

187 All of the input data sets can be found at:

188 <https://doi.org/10.5281/zenodo.755513> (GloBI Taxon Graph),

189 <http://files.opentreeoflife.org/ott/ott3.0/ott3.0.tgz> (Open Tree of Life Taxonomy)

190 <http://doi.org/10.5281/zenodo.1211767> (Wikidata)

191

192 A selection of intermediary and result datasets are available online (Poelen 2018d; Poelen  
193 2018a).

194

195 All of the scripts used to make the statements in the results can be found here

196 (<https://github.com/bio-guoda/guoda-datasets/tree/master/wikidata>) with instructions on how to

197 duplicate the analysis.

## 198 Results

199 After 10 minutes of processing, GloBI was linked to Wikidata using pre-existing identifier  
200 mappings. The Wikidata dump was 20 GB of compressed JSON with 40-50 million data items. It  
201 took about 10 minutes for GUODA to extract taxa (about 2.3 million) and their links in Wikidata  
202 and then less than one minute to map the Wikidata taxon graph to the GloBI taxon graph. The  
203 119,957 WikiData links that were added to GloBI increased its outgoing name links by 13.7%  
204 (Poelen 2018d). Eighty-seven percent (86.7%) of the external identifiers in Wikidata overlap  
205 with the external identifiers in OTT (Fig 3). Eighty-six percent (86.1%) of the external identifiers



206 in GloBI overlap with the external identifiers in OTT (Fig. 3). Wikidata provided mappings for  
207 65.2% of the external identifiers in GloBI (Fig. 3). Out of the 77,000 external identifiers that  
208 occurred only in OTT and GloBI, only 56 were inconsistent ([https://github.com/bio-](https://github.com/bio-guoda/guoda-datasets/blob/master/wikidata/inconsistentNameIdsGloBI_OTT.tsv)  
209 [guoda/guoda-datasets/blob/master/wikidata/inconsistentNameIdsGloBI\\_OTT.tsv](https://github.com/bio-guoda/guoda-datasets/blob/master/wikidata/inconsistentNameIdsGloBI_OTT.tsv)). These 56  
210 links pointed to seven OTT “taxa”. No inconsistent links were found between WD and GloBI.  
211 Out of the 38,000 links only found in GloBI, 9,000 were inconsistent ([https://github.com/bio-](https://github.com/bio-guoda/guoda-datasets/blob/master/wikidata/inconsistentNameIdsGloBIOnly.tsv)  
212 [guoda/guoda-datasets/blob/master/wikidata/inconsistentNameIdsGloBIOnly.tsv](https://github.com/bio-guoda/guoda-datasets/blob/master/wikidata/inconsistentNameIdsGloBIOnly.tsv)). The OTT,  
213 Wikidata, and GloBI identifier graphs related to this coverage analysis is a 74 MB compressed  
214 tab separated values file consisting of about 12 million identifier mapping records (see  
215 <https://zenodo.org/record/1213477/files/links-globi-wd-ott.tsv.gz>). The resulting Wikidata taxon  
216 objects were merged into GloBI’s Taxon Graph (Poelen 2018d).  
217  
218 In order for a mapping to be considered consistent, there can only be one identifier per resource  
219 included in each local graph. Thus, after removing the inconsistent identifiers, the external id  
220 overlap can be interpreted as an estimate of the number of shared taxon names between two  
221 databases (Table 2). This cannot be interpreted as total taxa in each resource.

## 222 Discussion

223 GUODA is a high performance computing resource for biodiversity science that provides  
224 scalable solutions for working with large data sets in a collaborative, online environment. The 10  
225 minute processing time for 20 GB of compressed JSON is far faster than any current mapping  
226 method used in biodiversity; however, it does benefit from the mapping already completed inside  
227 Wikidata. For example, the Wikidata entry for *Panthera leo*  
228 (<https://www.wikidata.org/wiki/Q140>) has 25 links to external databases, not all of them  
229 biodiversity-related. Other efforts using name-string-matching to link biodiversity databases take  
230 much longer to map resources together. For instance, EOL takes more than a day to map the  
231 content it receives from providers to a unified classification (Rice pers. comm.). Similarly, the  
232 taxon matching in BioNames and LinkOut took days to complete (Page pers. comm.). Projects  
233 like OTT, Wikidata, and GloBI that keep identifier-based taxonomic graphs make it easier to link  
234 databases at scale.

235  
236 Despite the notoriously poor nature of taxon names as identifiers, they are still commonly used to  
237 link biodiversity data. A much-discussed solution has been the use of universal, unique,  
238 persistent, resolvable identifiers across the biodiversity data landscape, but the social barrier to a  
239 universal identifier system has, thus far, proven insurmountable. Rather than rely on name strings  
240 or a universal identifier system, this method uses the graph of identifiers to map taxa across two  
241 databases. This identifier-based method has the potential to be faster and easier than name-string  
242 matching without some of the social difficulties of a single identifier system.

243  
244 Most biodiversity databases and nomenclatural authorities expose their data in idiosyncratic  
245 ways that are not suitable for batch processing. If data sources published their taxon identifier  
246 graph as a lookup table (as described in this paper) integrating across databases would be much

247 easier. Now, users have to learn a unique format for every data source. These lookup tables have  
248 the advantage of being easy to version and integrate.

249

250 In addition to fast linking of biodiversity databases, comparison of identifier graphs may be a  
251 scalable way to find inconsistencies, especially when multiple biodiversity databases/identifiers  
252 are included. By linking GloBI to OTT and WD, inconsistent names or false positive name  
253 matches were detected by considering the (lack of) overlap of GloBI names with OTT and WD  
254 external identifier schemes. These inconsistencies might be introduced by a dataset or a name  
255 resolution method that produces ambiguous results. In addition, inconsistencies can indicate a  
256 disputed / outdated name like “GLOBI:null@Senecio pectinatus” which maps to GBIF:8317096  
257 and GBIF:8414746. This would be considered an inconsistent mapping and suggests that  
258 *Senecio pectinatus* is an outdated name. Combining the speediness with the promise of  
259 scalability, a near-real-time name consistency check can be implemented to detect  
260 inconsistencies across various systems in the biodiversity data-ecosystem introduced by  
261 integration bugs, taxonomy updates or differences of interpretation.

262

263 GUODA has been available since 2015 and contains data dumps from GBIF, EOL TraitBank,  
264 iNaturalist, iDigBio, and BHL which are all accessible via a Jupyter notebook, web services, or  
265 Apache Spark shell on the command line. Despite its computing power and successful  
266 demonstrations at major conferences, GUODA has not been used to its full potential. The barrier  
267 of learning new programming and computing paradigms as well as developing an understanding  
268 of large dataset work flows seems to be a barrier to many in the biodiversity community. Despite  
269 this, GUODA is being used in several capacities. The Effechecka application generates  
270 taxonomic checklists using a web interface that allows a user to draw a polygon on a map and  
271 returns a deduplicated list of taxa aggregated from observation data held in GBIF, iNaturalist,  
272 etc. The EOL Freshdata project uses it to enable the detection of new occurrence records given  
273 geospatial and taxonomic and data source constraints and notifies interested users via email.  
274 Several workshops have used it to teach Spark programming skills to students at the University  
275 of Florida.

276

277 Future work on the GUODA infrastructure includes training and evaluating neural network  
278 models on image data, containerization of the GUODA components to allow the system to be run  
279 in additional data centers, and refinement of the end-user interface to integrate programming,  
280 source code, and publication to make research more reproducible. GUODA’s most impactful  
281 contribution has likely been the availability of readily formatted biodiversity data and new data  
282 sets will continue to be added to the collaboration platform, enabling domain experts and  
283 technical experts to answer new questions in the future.

284

285 GUODA, and hosted data analytics infrastructure in general, has the potential to drastically  
286 improve biodiversity science by making multiple biodiversity databases accessible to scientists  
287 for analysis on their laptop or desktop. Users still need to have some programming skills, which  
288 have now become an essential skill in biodiversity science.

289

## 290 Conclusions

291 Sharing information between biodiversity databases can be difficult because of the amount and  
292 heterogeneity of the data and the identifiers. Most mappings are done using taxonomic name  
293 strings at great expense. We were able to map Wikidata to GloBI in 10 minutes using identifier  
294 graphs and GUODA, a high performance computing infrastructure developed through  
295 collaboration between diverse players. The mapping increased GloBI's outgoing name links by  
296 13.7%. This method of mapping across databases using identifier graphs is faster than comparing  
297 name strings and can help find inconsistencies that point to a disputed or outdated name.  
298 GUODA, and systems like it, have the potential to revolutionize biodiversity science by bringing  
299 diverse technically minded people together with high performance computing resources that are  
300 accessible from a laptop or desktop.

## 301 Acknowledgements

302 The authors would like to acknowledge support and resources provided by the ACIS lab. The  
303 authors would like to acknowledge José A.B. Fortes for providing infrastructure and creating  
304 room for collaboration. Funding was provided by David Rubenstein and the Encyclopedia of  
305 Life, which also helped establish an informal yet pragmatic cross-institutional collaboration and  
306 by iDigBio, NSF award 1547229, which provided for staff, travel, and publication fees for this  
307 collaboration.

## 308 References

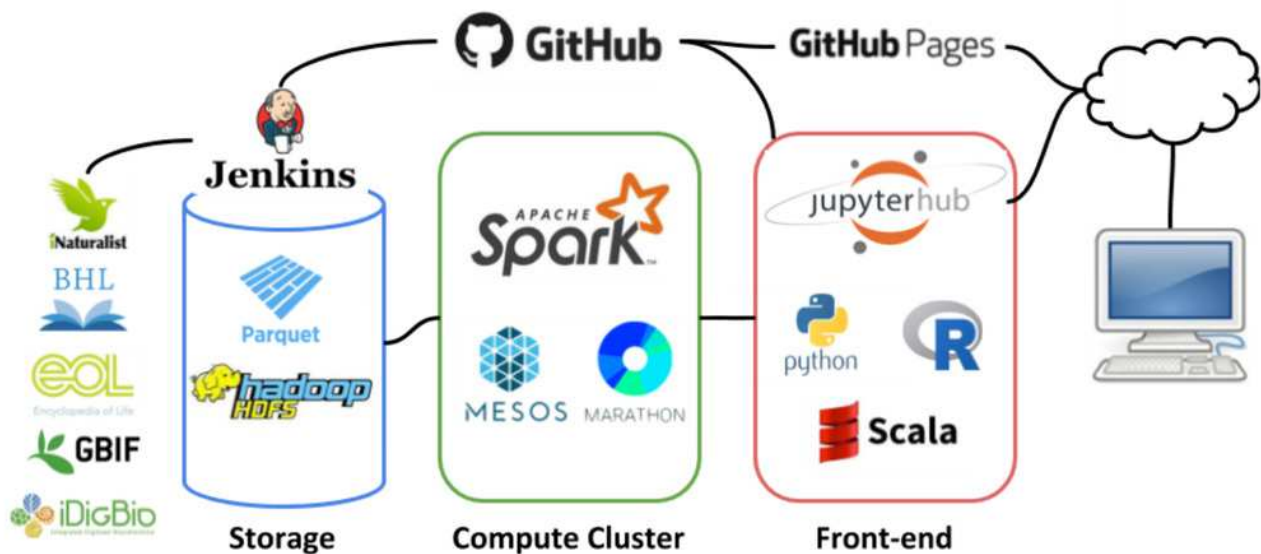
- 309 Bingham, H. et al., 2017. The Biodiversity Informatics Landscape: Elements, Connections and  
310 Opportunities. *RIO*, 3, e14059. Available at: <https://riojournal.com/article/14059/>.
- 311 Hindman, B. et al., 2011. Mesos: a platform for fine-grained resource sharing in the data center.  
312 In *Proceedings of the 8th USENIX conference on Networked systems design and*  
313 *implementation*. pp. 295–308.
- 314 Hortal, J. et al., 2015. Seven Shortfalls that Beset Large-Scale Knowledge of Biodiversity.  
315 *Annual Review of Ecology, Evolution, and Systematics*, 46(1), pp. 523–549. Available at:  
316 <http://www.annualreviews.org/doi/10.1146/annurev-ecolsys-112414-054400>.
- 317 Kluyver, T. et al., 2016. Jupyter Notebooks -- a publishing format for reproducible  
318 computational workflows. In F. Loizides & B. Schmidt, eds. *Positioning and Power in*  
319 *Academic Publishing: Players, Agents and Agendas*. Amsterdam: IOS PRes, pp. 87–90.
- 320 Page, R., 2008. Biodiversity informatics: the challenge of linking data and the role of shared  
321 identifiers. *Briefings in Bioinformatics*, 9(5), pp. 345–354. Available at:  
322 <https://academic.oup.com/bib/article/9/5/345/267216>.
- 323 Page, R., 2013. BioNames: linking taxonomy, texts, and trees. *PeerJ*, 1, p.e190. Available at:  
324 <https://peerj.com/articles/190/>.
- 325 Page, R., 2011. Linking NCBI to Wikipedia: a wiki-based approach. *PLoS Currents*, 3,  
326 RRN1228.
- 327 Page, R., 2007. Tbmap: a taxonomic perspective on the phylogenetic database treebase. *BMC*

- 328 *Bioinformatics*, 8(1), p. 158.
- 329 Parr, C. et al., 2014. The encyclopedia of life v2: providing global access to knowledge about life  
330 on earth. *Biodiversity Data Journal*, 2, e1079.
- 331 Poelen, J., 2018a. 20 GB in 10 minutes: Data linking across major biodiversity databases: Data  
332 supplements. Version 0.1 Zenodo <http://doi.org/10.5281/zenodo.1213477>
- 333 Poelen, J., 2018b. Global Biotic Interactions: Taxon Graph. Version 0.4.2 Zenodo  
334 <http://doi.org/10.5281/zenodo.1210315>
- 335 Poelen, J., 2018c. Global Biotic Interactions: Taxon Graph. Version 0.3.0 Zenodo  
336 <http://doi.org/10.5281/zenodo.1210308>
- 337 Poelen, J., 2018d. Global Biotic Interactions: Taxon Graph. Version 0.3.1 Zenodo  
338 <http://doi.org/10.5281/zenodo.1213465>
- 339 Poelen, J.H., Simons, J.D. & Mungall, C.J., 2014. Global biotic interactions: An open  
340 infrastructure to share and analyze species-interaction datasets. *Ecological Informatics*, 24,  
341 pp. 148–159.
- 342 Rees, J.A. & Cranston, K., 2017. Automated assembly of a reference taxonomy for phylogenetic  
343 data synthesis. *Biodiversity Data Journal*, 5, e12581.
- 344 Shvachko, K. et al., 2010. The Hadoop distributed file system. In *MSST '10 Proceedings of the*  
345 *2010 IEEE 26th Symposium on Mass Storage Systems and Technologies*. pp. 1–10.
- 346 Voß, J., 2016. wikidata-taxonomy 0.2.7. Version 0.2.7 Zenodo  
347 <http://doi.org/10.5281/zenodo.60708>
- 348 Vrandečić, D. & Krötzsch, M., 2014. Wikidata: A free collaborative knowledgebase.  
349 *Communications of the ACM*, 57(10), pp. 78–85. Available at:  
350 <https://cacm.acm.org/magazines/2014/10/178785-wikidata/fulltext>.
- 351 Wikidata, 2018. Wikidata dump 2017-12-27. Zenodo <http://doi.org/10.5281/zenodo.1211767>
- 352 Zaharia, M. et al., 2016. Apache Spark: a unified engine for big data processing.  
353 *Communications of the ACM*, 59(11), pp. 56–65.
- 354

# Figure 1

## GUODA Infrastructure

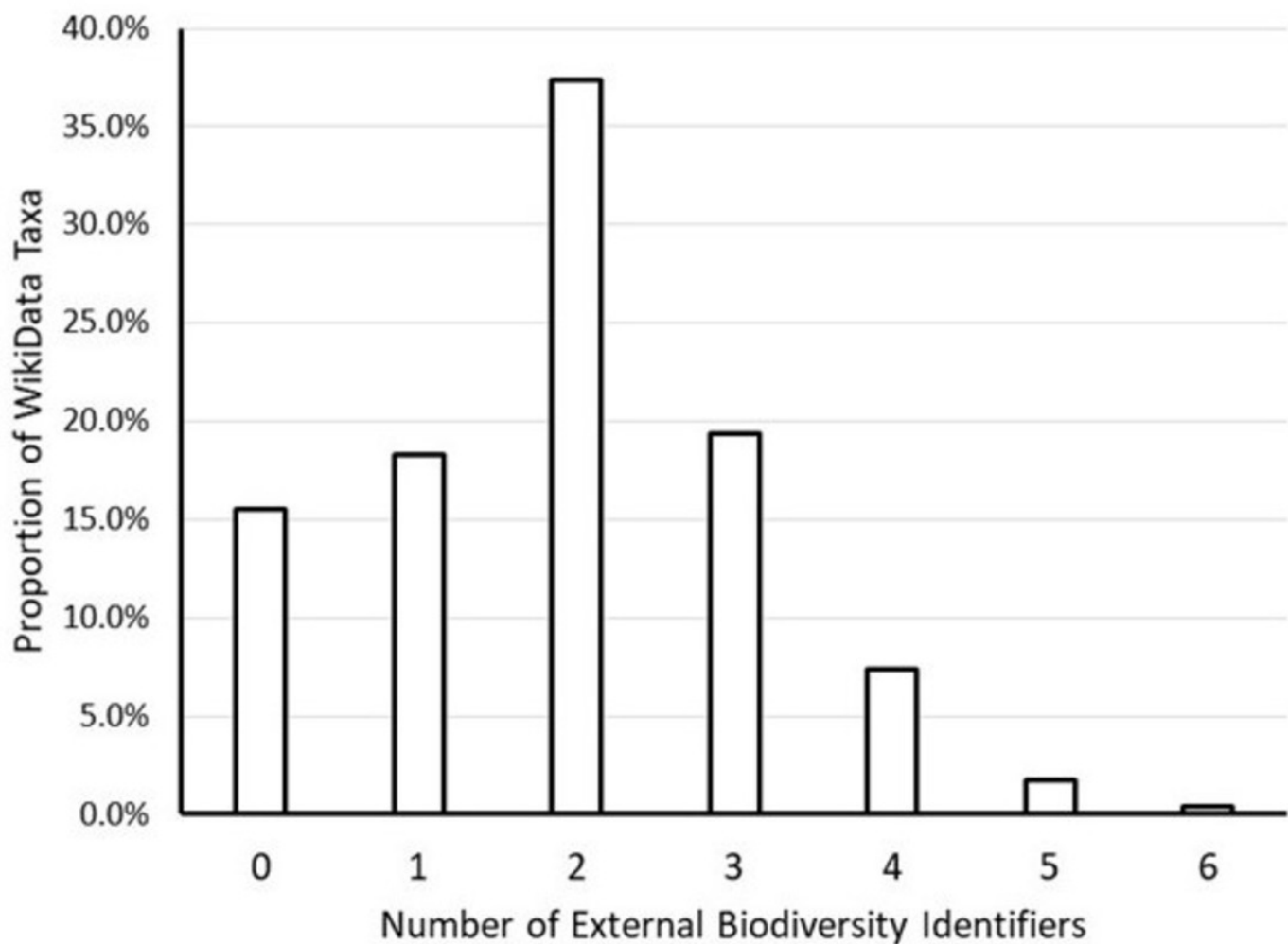
Data from biodiversity databases is loaded into GUODA as Parquet files (Storage). When a user working in a Jupyter Notebook (Front-end Server) triggers a job interactively or via GitHub and Jenkins, the data are analyzed using Apache Spark (Compute Cluster). This infrastructure allows a user working from a laptop or desktop to compute over multiple biodiversity databases at once.



## Figure 2

### Frequency of Wikidata taxa linked to biodiversity databases

This graph shows the proportion of the approximately 2.3 million Wikidata taxa with 0,1,2, etc. links to external biodiversity databases (NCBI, ITIS, GBIF, EOL, FishBase, Index Fungorum and iNaturalist). The majority of Wikidata taxa had at least two links. A little more than 15% of Wikidata taxa had no links to external biodiversity databases.

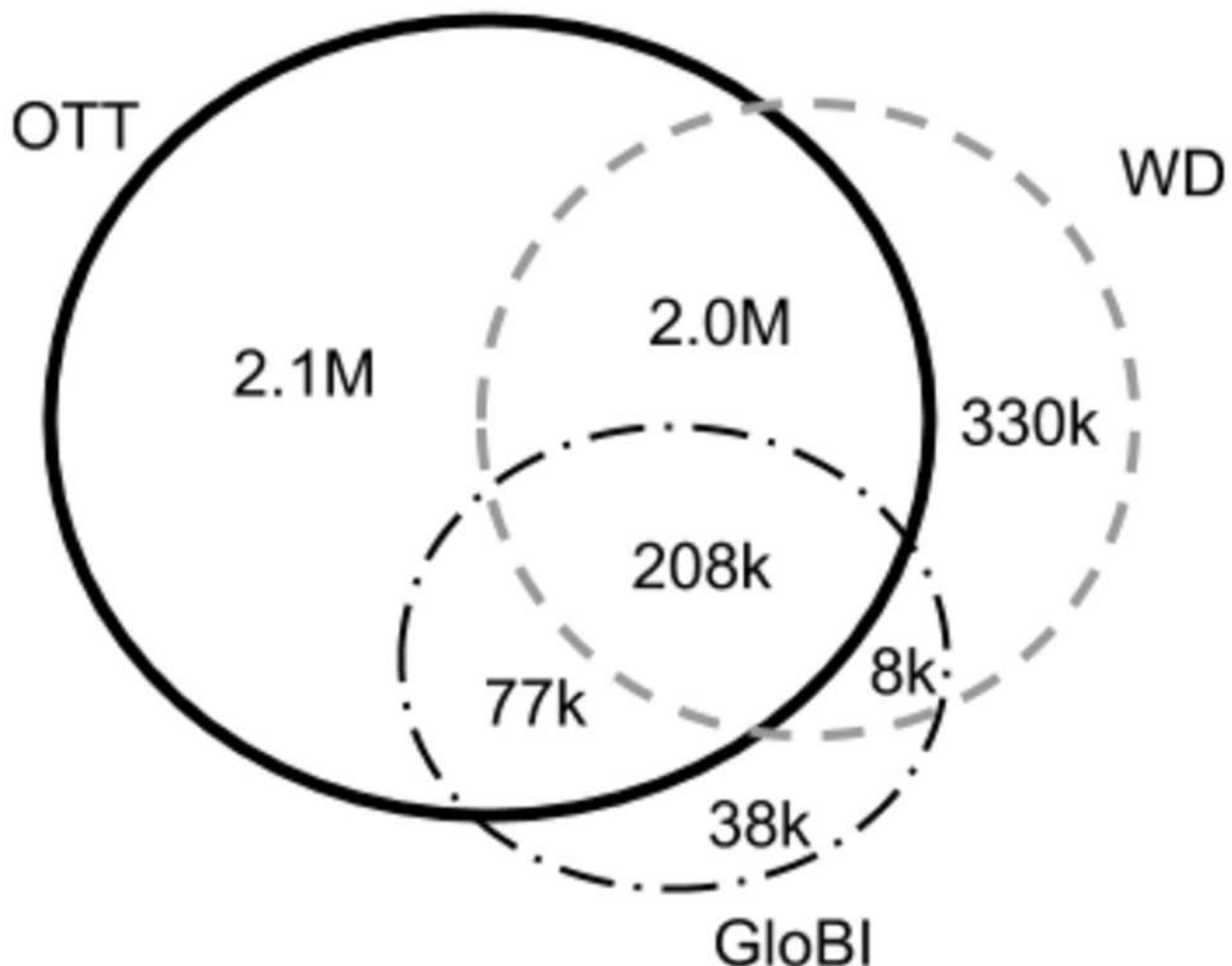




## Figure 3

Identifier overlap between Wikidata (WD), OTT, and GloBI

This Venn Diagram shows the number of overlapping external identifiers that can be found in one of three databases. Only 208,000 external ids can be found in all three. These consisted of 23,000 WoRMS links, 72,000 NCBI links, 103,000 GBIF links and 10,000 IF links. Over two million ids are only known to one of the three databases. OTT contains more than half of the external ids in Wikidata and in GloBI, but neither contain half of the external ids in OTT. Mapping Wikidata to GloBI matched 65.2% of the external ids in GloBI.



**Table 1** (on next page)

Selected biodiversity databases and their size

1

<b>Database</b>	<b>Data Quantity (Jan 2018)</b>	<b>Size (compressed)</b>
GBIF	964,547,793 occurrence records	139 GB
Catalogue of Life/ITIS	1.7 million taxa	2.9 GB
GloBI	3,363,528 interactions	206 MB
iDigBio	106,922,498 specimen records	35.5 GB
GenBank	206,293,625 sequences	3 TB
Biodiversity Heritage Library	53,739,062 pages	2.7 GB
WoRMS	243,323 marine species	71 MB
OpenTree	2,722,024 taxa and 6,810 trees	189 MB
EOL TraitBank	Over 11 million records	46 GB uncompressed
EOL	7,705,748 data objects (May 2017)	10 TB uncompressed
Wikidata	42,648,426 data items	20 GB

2

**Table 2** (on next page)

Absolute and relative link counts from OTT, WD, and GloBI compared to WoRMS, GBIF, Index Fungorum (IF), and NCBI

1

	<b>WoRMS</b>	<b>GBIF</b>	<b>IF</b>	<b>NCBI</b>	<b>combined</b>
<b>OTT</b>	327929 (100%)*	2451566 (100%)	276262 (100%)	1355207 (100%)	4410964 (100%)
<b>WD</b>	288110 (88%)	1779789 (73%)	76497 (28%)	410092 (30%)	2554488 (58%)
<b>GloBI</b>	68565 (21%)	315173 (13%)	33400 (12%)	704361 (52%)	1121499 (25%)

2 \*Overlap between each resource and OTT is set at 100%. The other percentages give a relative  
3 estimate of size and scale and should not be interpreted as overlapping ids.

4