

Liberating host-virus knowledge from biological dark data

Non Peer-reviewed author version

Upham, Nathan; Poelen, Jorrit H.; Paul, Deborah Leo; Groom, Quentin John; Simmons, Nancy B.; VANHOVE, Maarten; Bertolino, Sandro; Reeder, DeeAnn M.; Bastos-Silveira, Cristiane; Sen, Atriya; Sterner, Beckett; Franz, Nico; Guidoti, Marcus; Penev, Lyubomir & Agosti, Donat (2021) Liberating host-virus knowledge from biological dark data.

DOI: [10.32942/osf.io/tzekq](https://doi.org/10.32942/osf.io/tzekq)

Handle: <http://hdl.handle.net/1942/36893>

Title: Liberating host-virus knowledge from biological dark data

Authors: Nathan S. Upham, Ph.D.^{1,*}, Jorrit H. Poelen, M.S.², Deborah Paul, M.S.³, Quentin J. Groom, Ph.D.⁴, Nancy B. Simmons, Ph.D.⁵, Prof. Maarten P. M. Vanhove, Ph.D.⁶, Sandro Bertolino, Ph.D.⁷, Prof. DeeAnn M. Reeder, Ph.D.⁸, Cristiane Bastos-Silveira, Ph.D.⁹, Atriya Sen, Ph.D.¹⁰, Beckett Sterner, Ph.D.¹, Prof. Nico M. Franz, Ph.D.¹, Marcus Guidoti, Ph.D.¹¹, Prof. Lyubomir Penev, Ph.D.¹², and Donat Agosti, Ph.D.¹³

Affiliations:

¹School of Life Sciences, Arizona State University, Tempe, AZ 85782, USA. <https://orcid.org/0000-0001-5412-9342> (NSU); <https://orcid.org/0000-0001-5219-7616> (BS); <https://orcid.org/0000-0001-7089-7018> (NMF)

²Ronin Institute for Independent Scholarship, Montclair, NJ 07043, USA; Cheadle Center for Biodiversity and Ecological Restoration, UC Santa Barbara, Santa Barbara, CA 93106, USA. <https://orcid.org/0000-0003-3138-4118>

³Illinois Natural History Survey, University of Illinois, Urbana-Champaign, IL 61820, USA. <https://orcid.org/0000-0003-2639-7520>

⁴Meise Botanic Garden, 1860 Meise, Belgium. <https://orcid.org/0000-0002-0596-5376>

⁵Department of Mammalogy, Division of Vertebrate Zoology, American Museum of Natural History, New York, NY 10024, USA. <https://orcid.org/0000-0001-8807-7499>

⁶Research Group Zoology: Biodiversity and Toxicology, Centre for Environmental Sciences, Hasselt University, 3590 Diepenbeek, Belgium. <https://orcid.org/0000-0003-3100-7566>

⁷Department of Life Sciences and Systems Biology, University of Turin, 10123 Torino, Italy. <https://orcid.org/0000-0002-1063-8281>

⁸Department of Biology, Bucknell University, Lewisburg, PA 17837, USA. <https://orcid.org/0000-0001-8651-2012>

⁹Centre for Ecology, Evolution and Environmental Changes (CE3C), Universidade de Lisboa, Lisbon, Portugal. <https://orcid.org/0000-0001-8249-9383>

¹⁰Department of Computer Science, University of New Orleans, New Orleans, LA 70148, USA.

¹¹Plazi, Porto Alegre, Brazil. <https://orcid.org/0000-0003-1744-6191>

¹²Pensoft Publishers, Sofia, Bulgaria. <https://orcid.org/0000-0002-2186-5033>

¹³Plazi, Bern, Switzerland. <https://orcid.org/0000-0001-9286-1200>

*Correspondence to: Nathan S. Upham nathan.upham@asu.edu, Arizona State University, Natural History Collections, P.O. Box 874108, Tempe, AZ, 85287-4108, USA

Abstract: 191 words; Main text (with figure legends): 2,102 words; References: 27 total.

Abstract:

Connecting basic data about bats and other potential hosts of SARS-CoV-2 with their ecological context is critical for understanding the emergence and spread of COVID-19. However, when
40 global lockdown started in March 2020, the world's bat experts were locked out of their research laboratories, which, in turn, locked up large volumes of offline ecological and taxonomic data. Pandemic lockdowns have put a magnifying glass on the long-standing problem of biological 'dark data': data which are published, but disconnected from digital knowledge resources, and thus unavailable for high-throughput analysis. Knowledge of host-to-virus ecological interactions will
45 be biased until this challenge is addressed. Here we outline two viable solutions: (i) how to interconnect published data about host organisms, viruses, and other pathogens in the short term; and (ii) how to shift the publishing paradigm beyond unstructured text ('PDF prison') to labeled networks of digital knowledge. Biological taxonomy is foundational to both solutions as the indexing system for biodiversity data. Building digitally connected 'knowledge graphs' of host-
50 pathogen interactions will establish the needed agility for quickly identifying reservoir hosts of novel zoonoses, allow for more robust predictions of emergence, and thereby strengthen planetary health systems.

Main text:

55 An irony of COVID-19 likely originating from a bat-borne coronavirus (1) is that the global lockdown to quell the pandemic also locked up physical access to much-needed knowledge about bats. Basic data about bat diversity, ecology, and geography, as well as that of other potential mammal hosts (1,2), was suddenly critical for understanding SARS-CoV-2's emergence and spread. However, with the world's bat experts unable to access their research laboratories, any
60 undigitized or offline data was also locked down. In a matter of days, lockdowns around the world had dramatically reduced the accessibility of scientific knowledge. Why, in this digitally-connected age, was basic knowledge about species and their ecological interactions not already digitized,

online, and openly accessible to all? What must be done to improve global access to public health-related biodiversity knowledge?

65

Understanding why biodiversity science was unprepared—and how to fix it before the next crisis—has been a hot topic, spawning multiple taskforces in the biodiversity research community since the pandemic began (e.g., (3–5)). Of key interest has been mending the chasm in knowledge transfer from the physical biocollections, which contain the preserved specimens, tissues, and associated material used to describe biodiversity, to biomedical scientists in health-related fields like infectious disease, epidemiology, and virology. Most biodiversity knowledge ever published remains effectively locked in textual, unstructured articles, and is thus isolated from efforts to synthesize global ecological interactions. These data are ‘known’ in publications but are digitally disconnected—revealing a striking ‘knowledge frontier’ that is preventing scientists from digitally discovering their existence. With human activities like land conversion hastening the emergence of zoonoses (6), it is increasingly urgent to build interconnected networks of digital knowledge.

70

75

ILLUMINATING BIODIVERSITY DARK DATA

Physicists accept that dark matter exists, but they have difficulty measuring it. In the same way, biodiversity scientists are aware of large quantities of ‘dark data’ in publications, but have difficulty synthesizing it, either because such data are old and rare (e.g., inside archival or gray literature) or new and locked (e.g., behind paywalls, in digitally unreadable formats, or unlinked to other data). Traditionally, a particular research project might manually synthesize information from hundreds or thousands of articles in disparate formats over the course of years, yielding a comprehensive ‘snapshot’ of written knowledge. Still today, gathering the widely scattered biodiversity data relevant to mammal host-virus interactions would take years instead of the needed weeks for responding to a crisis like the SARS-CoV-2 outbreak. Remarkably, new articles continue to worsen the dark data dilemma, since the ubiquitous ‘portable data format’ (PDF)

80

85

requires substantial efforts to make ecological phenomena like host-virus interactions extractable
for re-use (hence the term 'PDF prison' (7)). To address deeply interconnected global problems
like COVID-19, it is imperative to implement new solutions rooted in building expansive digital
knowledge (Fig. 1).

For data to form digital knowledge, they must first be published in datasets that are open access
and [FAIR](#) — Findable on the web, digitally Accessible, Interoperable among different computing
systems, and thus Reusable for later analyses. Satisfying all of these criteria opens the door for
creating highly useful 'knowledge graphs' (8,9), in which digital open data are meaningfully linked
together on massive scales, forming knowledge that is collectively greater than its sum. As Tim
Berners-Lee presciently wrote in 2006, "it is the unexpected re-use of information which is the
value added by the web" (10). Illuminating the zoonotic origins of COVID-19 is exactly the kind of
unexpected re-use of data that biodiversity science was ill-prepared to address at the start of the
pandemic. Building a comprehensive host-virus knowledge graph will furthermore enable rapidly
improving artificial intelligence algorithms (e.g., in the fields of natural language processing, NLP
(11), and knowledge reasoning (12)) to flexibly learn from the structure of digital knowledge.

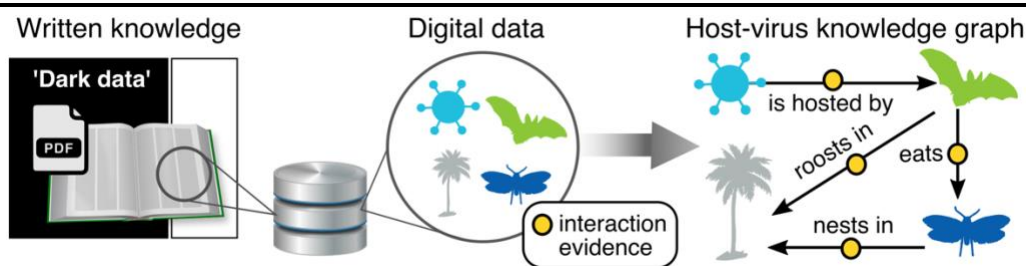


Fig. 1. The evolution of biodiversity knowledge from analog to digital. Extracting written knowledge from publications into databases is only the first step toward creating forms of digital, structured knowledge in which ecological interactions and evidence thereof are additionally annotated. Such 'knowledge graphs' include levels of confidence in each annotation as derived from evidence sources, which enable high-throughput integrative modeling of complex ecological dynamics like viral spillover. Much written knowledge is undigitized and digitally disconnected, forming 'dark data' from the perspective of synthetic knowledge graphs.

TAXONOMY AS THE KEY TO HOST-VIRUS KNOWLEDGE

115 Linking viruses to animal hosts, hosts to environments, and hosts to other hosts is the raw material
needed to build a host-virus knowledge graph (Fig. 2A). However, meaningfully connecting host
species, viral species, and their ecological traits requires mastery of a fundamental but undersold
discipline: biological taxonomy. For at least three centuries, mainstream science has used the
names of species—most often the ‘genus & species’ pair of Linnaean taxonomy—to index
120 research findings. Virtually all observations about organismal behaviors and functions, habitats,
genomics, and pathogens are linked to species names via sections of publications called
‘taxonomic treatments’, in which authors describe the boundaries of species (and other taxa)
based on physical evidence. Because that evidence—especially from preserved specimens and
derived data like DNA sequences—has improved along with the science of taxonomy through
125 time, multiple names may have been used to refer to similar sets of organisms. Thus, making
sense of biodiversity data requires keeping track of how the meaning of taxonomic names has
changed historically (e.g., synonyms, varying name usages).

How species names have been used by different authors over time is the ‘taxonomic passkey’ for
130 opening otherwise locked host-to-virus interactions in publications. By linking species names,
evidence, and taxonomic treatments through time, it is possible to create ‘taxonomic intelligence’
services (13) that allow for flexible conversion of named species data across taxonomies. For
example, SARS-like coronaviruses observed in horseshoe bats identified as *Rhinolophus sinicus*
in 2013 (14) need to be resolved relative to the 2019 re-classification of portions of this species
135 as *R. thomasi* and *R. rouxii* (15). However, updating the taxonomy of named data when taxonomic
concepts have been split is not yet possible aside from manually on small scales. Existing
taxonomic infrastructures like the [Catalogue of Life](#) have not prioritized building large scale
solutions to this problem, primarily because taxonomic changes are often very rapid. Even in a
relatively well known group like mammals, the global number of species recognized has changed

140 by >40% in the last 25 years (16) over which time the number of described viruses has increased
 by a staggering 400% (17). Keeping track of mammal-to-virus interactions relative to that
 taxonomic flux has not been incentivized in proportion to its importance for understanding
 zoonotic emergence. Therefore, we must make efforts to prioritize the building of taxonomic
 intelligence services, which will then enable the extraction and meaningful linking of named host-
 145 to-virus interaction data on planetary scales.

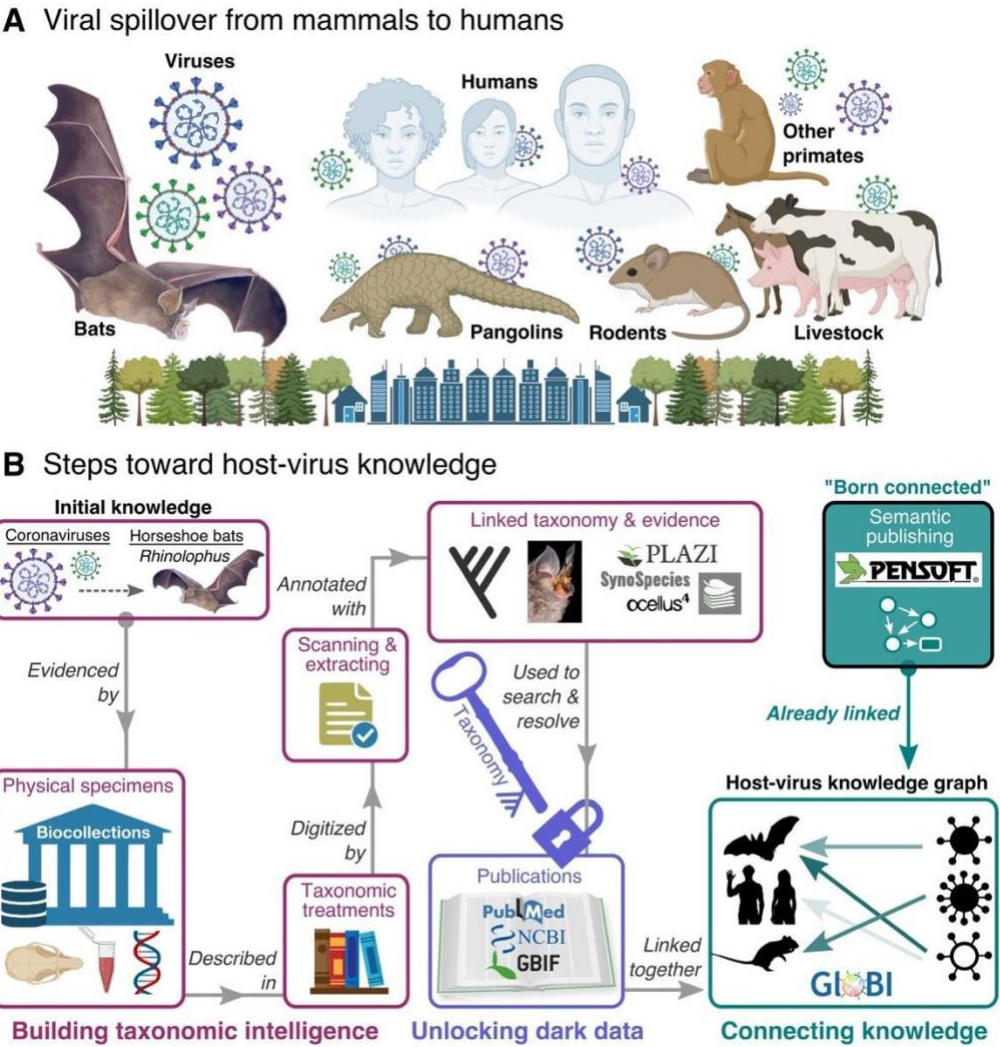


Fig. 2. Connecting digital knowledge of host-virus interactions. (A) The sharing of viruses among humans and other mammals is remarkably common, yet the ecological circumstances under which spillover occurs are poorly understood. (B) Digitally liberating ecological knowledge from locked publications requires building taxonomic intelligence—i.e., how and why species names have been used through time—and then using that taxonomic ‘passkey’ to liberate and

connect ‘dark’ interaction data hidden in publications. Alternatively, data can be ‘born connected’ if new articles are published using computer-readable (semantic) tags for ecological interactions like ‘has host’ or ‘pathogen of.’ Both pathways will enable newly comprehensive knowledge graphs that connect host-virus interactions with underlying evidence.

TOWARD A HOST-VIRUS KNOWLEDGE GRAPH

Thankfully, two decades of work in the digital knowledge arena (e.g., (8,10,13,18,19)) has established foundations for a two-pronged approach to building host-virus knowledge (Fig. 2B).

First, dark data needs to be liberated from existing publications. These efforts are being led by [Plazi](#) (18)—a pioneering platform for literature digitization, extraction, and linking—to create new flows of digital data from printed books, archives, and otherwise locked publications. For example, the Plazi services [Synospecies](#) and [Ocellus](#) have recently indexed taxonomic names and images, respectively, from taxonomic treatments spanning from Linnaeus’ initial 1758 publication *Systema Naturae* to the recent *Handbook of the Mammals of the World* series (20), making them available on the [Biodiversity Literature Repository](#) (21). Once digitally indexed, taxonomic data can be annotated and connected to biocollection-based evidence to formally align taxonomic names with their biological meanings. This liberated taxonomic knowledge allows for more robust literature searches and subsequent name translation of host-virus interaction data. Such efforts have already discovered reliable data on 1,146 host-virus interactions from selected publications ([Coronavirus-Host Community](#) on Zenodo). Second, new articles need to be published without creating more dark data. Exemplary in this area are efforts being led by [Pensoft](#)—publisher of biodiversity journals such as *ZooKeys*—to publish using computer-readable semantic annotations during the normal publishing process (22), allowing immediate indexing afterwards (23). For example, Pensoft responded to COVID-19 by beginning to index parts of speech such as ‘[has host](#)’ and ‘[pathogen of](#)’ to assist with mining biotic interactions from article texts and tables, which has netted over 2,000 biotic interactions now annotated as article metadata (24). Such digital enhancements greatly streamline the process of data extraction, because new articles already

contain digital text, linking terms, and thus a native form of digital knowledge. These data are
180 'born connected' relative to the post-processing steps needed with traditional PDF publishing.

To build a singular host-virus knowledge graph requires a central hub for discovering relevant
data, resolving disparate taxonomies, and connecting the resulting insights. Promising progress
by the Global Biotic Interactions database [GloBI](#)—an open-access ecological network across all
185 of life (19)—has led to new pipelines for ecological data to flow from sources of both 'old' (20,21)
and 'new' literature (24). From April to October 2020 alone, these pipelines resulted in adding
>53,000 host-virus data points to GloBI (see dataset on Zenodo (25)). These associations involve
19% more valid species of mammals than were identified in a recent host-virus synthesis (897 vs.
754 species in (26)). Such a dramatic initial effort illustrates the potential for broad-scale data
190 linking to yield new insights. Yet these are small steps relative to what could eventually comprise
a comprehensive and taxonomically nimble graph of not only host-virus but host-pathogen and
broader ecological knowledge. What interconnected phenomena might be illuminated when such
knowledge is freely available to the world's scientists and public health specialists?

195 **BEYOND THE PDF: KNOWLEDGE THAT IS BORN CONNECTED**

We have outlined ways to interconnect, and thus liberate, previously 'dark' host-pathogen
interactions from publications. However, doing so is expensive and so is infeasible at scale if
publishers continue to publish under the same paradigm. Therefore, we recommend three
immediate policy changes: (i) major journals should switch to publishing formats that are not only
200 open access and FAIR, but also semantically tagged with terms relevant to broad-scale ecological
interactions (especially host-pathogen and host-host relationships); (ii) academic institutions
should incentivize (e.g., via tenure evaluations, paying open-access fees) publishing in such 'born
connected' journals; and (iii) investments in data generation should be balanced with
infrastructure enhancements for data reuse, incentivizing the construction of increasingly

205 complete biodiversity knowledge graphs. Taxonomists, ecologists, data scientists, and
policymakers have essential roles to play in this paradigm shift toward digital knowledge.

The value-added by digitally connected knowledge is tremendous, both for its potential to build
nonlinear insights and to expand the capacity of biodiversity researchers around the world,
210 especially in the Global South (27). Limitations to accessing biodiversity information in developing
countries are diverse, including gaps in geographical knowledge; lack of data sharing among and
between scientists and policymakers; inaccessible presentations of information; and limited
financial resources. Efforts are hence needed not only to increase, as is often called for,
biodiversity monitoring, but also to support the capacity of local scientific and citizen communities
215 to mobilize the resulting data into digital knowledge infrastructures. Pandemics demonstrate that
human societies are inextricably linked regardless of wealth, so that building a biodiversity
knowledge commons will benefit all.

We cannot continue to waste resources to rediscover biodiversity a second time. Unprecedented
220 reliability of knowledge about biological interactions is now required to address multiple socio-
ecological challenges, from COVID-19 to biodiversity loss and runaway climate change, each of
which exists on scales too massive and too detailed for any one individual to observe alone. The
COVID-19 pandemic teaches us that siloed science does not serve society as well as its
alternative. Multiple novel solutions, including vaccines and treatments, are beginning to free us
225 from this pandemic. However, the solution for our limited access to ecological knowledge is
already here. We already have much of the technology needed to liberate and connect
biodiversity data across the entire tree of life—what is most lacking is the collective will to do so.

230

REFERENCES

- 235 1. Boni MF, Lemey P, Jiang X, Lam TT-Y, Perry BW, Castoe TA, et al. Evolutionary origins of the SARS-CoV-2 sarbecovirus lineage responsible for the COVID-19 pandemic. *Nat Microbiol* [Internet]. 2020 Jul 28 [cited 2020 Jul 28];1–10. Available from: <https://www.nature.com/articles/s41564-020-0771-4>
- 240 2. Xia X. Extreme Genomic CpG Deficiency in SARS-CoV-2 and Evasion of Host Antiviral Defense. *Mol Biol Evol* [Internet]. 2020 Sep 1 [cited 2020 Sep 30];37(9):2699–705. Available from: <https://academic.oup.com/mbe/article/37/9/2699/5819559>
- 245 3. CETAF-DiSSCo COVID-19 Taskforce. Communities Taking Action | CETAF – Consortium of European Taxonomic Facilities | Distributed Systems of Scientific Collections – DiSSCo [Internet]. 2020 [cited 2020 Aug 31]. Available from: <https://cetaf.org/covid19-taf-communities-taking-action>
4. ViralMuse. iDigBio Wiki- ViralMuse Task Force [Internet]. iDigBio Wiki. 2020 [cited 2020 Oct 30]. Available from: https://www.idigbio.org/wiki/index.php/ViralMuse_Task_Force
5. Research Data Alliance. RDA-COVID19 [Internet]. RDA. 2020 [cited 2020 Oct 30]. Available from: <https://www.rd-alliance.org/groups/rda-covid19>
- 250 6. Faust CL, McCallum HI, Bloomfield LSP, Gottdenker NL, Gillespie TR, Torney CJ, et al. Pathogen spillover during land conversion. *Ecol Lett* [Internet]. 2018 [cited 2020 Sep 18];21(4):471–83. Available from: <https://onlinelibrary.wiley.com/doi/abs/10.1111/ele.12904>
- 255 7. Agosti D, Catapano T, Sautter G, Egloff W. The Plazi Workflow: The PDF prison break for biodiversity data. *Biodivers Inf Sci Stand* [Internet]. 2019 Jun 13 [cited 2021 May 19];3:e37046. Available from: <https://biss.pensoft.net/article/37046/>
8. Penev L, Dimitrova M, Senderov V, Zhelezov G, Georgiev T, Stoev P, et al. OpenBiodiv: A Knowledge Graph for Literature-Extracted Linked Open Data in Biodiversity Science. *Publications* [Internet]. 2019 Jun [cited 2020 Oct 29];7(2):38. Available from: <https://www.mdpi.com/2304-6775/7/2/38>
- 260 9. Page R. Towards a biodiversity knowledge graph. *Res Ideas Outcomes* [Internet]. 2016 Apr 7 [cited 2020 Aug 18];2:e8767. Available from: <http://rio.pensoft.net/articles.php?id=8767>
10. Berners-Lee T. Linked Data - Design Issues [Internet]. 2006 [cited 2020 Sep 14]. Available from: <https://www.w3.org/DesignIssues/LinkedData.html>
- 265 11. Burgdorf A, Pomp A, Meisen T. Towards NLP-supported Semantic Data Management. *ArXiv200506916 Cs* [Internet]. 2020 May 14 [cited 2020 Nov 4]; Available from: <http://arxiv.org/abs/2005.06916>
- 270 12. Bellomarini L, Sallinger E, Vahdati S. Chapter 6 Reasoning in Knowledge Graphs: An Embeddings Spotlight. In: Janev V, Graux D, Jabeen H, Sallinger E, editors. *Knowledge Graphs and Big Data Processing* [Internet]. Cham: Springer International Publishing; 2020 [cited 2021 May 21]. p. 87–101. (Lecture Notes in Computer Science). Available from: https://doi.org/10.1007/978-3-030-53199-7_6

13. Bisby FA. The Quiet Revolution: Biodiversity Informatics and the Internet. Science [Internet]. 2000 Sep 29 [cited 2021 May 19];289(5488):2309–12. Available from: <https://science.sciencemag.org/content/289/5488/2309>
14. Ge X-Y, Li J-L, Yang X-L, Chmura AA, Zhu G, Epstein JH, et al. Isolation and characterization of a bat SARS-like coronavirus that uses the ACE2 receptor. Nature [Internet]. 2013 Nov [cited 2020 Nov 12];503(7477):535–8. Available from: <https://www.nature.com/articles/nature12711>
15. Burgin CJ. *Rhinolophus sinicus* K. Andersen 1905. Fam Rhinolophidae Horseshoe Bats Pp 280-332 Handb Mamm World Vol 9 Lynx Edicions Pp 325-326 [Internet]. 2019 Oct 31 [cited 2020 Nov 12]; Available from: <https://zenodo.org/record/3808964>
16. Burgin CJ, Colella JP, Kahn PL, Upham NS. How many species of mammals are there? J Mammal [Internet]. 2018 Feb 1 [cited 2018 Feb 6];99(1):1–14. Available from: <https://academic.oup.com/jmammal/article/99/1/1/4834091>
17. International Committee on Taxonomy of Viruses. ICTV Historical taxonomy releases [Internet]. 2020 [cited 2020 May 3]. Available from: https://talk.ictvonline.org/taxonomy/p/taxonomy_releases
18. Agosti D, Egloff W. Taxonomic information exchange and copyright: the Plazi approach. BMC Res Notes [Internet]. 2009 [cited 2020 Aug 11];2(1):53. Available from: <http://bmresnotes.biomedcentral.com/articles/10.1186/1756-0500-2-53>
19. Poelen JH, Simons JD, Mungall CJ. Global biotic interactions: An open infrastructure to share and analyze species-interaction datasets. Ecol Inform [Internet]. 2014 Nov 1 [cited 2020 Apr 17];24:148–59. Available from: <http://www.sciencedirect.com/science/article/pii/S1574954114001125>
20. Agosti D. Time for an interim review of Plazi's Covid-19 related activities [Internet]. 2020. Available from: <http://plazi.org/news/beitrag/time-for-an-interim-review-of-plazis-covid-19-related-activities/3e26b3bc95a4b39f0a2a9d7fcce8b19/>
21. Agosti D, Catapano T, Sautter G, Kishor P, Nielsen L, Ioannidis-Pantopikos A, et al. Biodiversity Literature Repository (BLR), a repository for FAIR data and publications. Biodivers Inf Sci Stand [Internet]. 2019 Jun 19 [cited 2020 May 1];3:e37197. Available from: <https://zenodo.org/record/3257816#.XqzzfJopC7M>
22. Penev L, Catapano T, Agosti D, Georgiev T, Sautter G, Stoev P. Implementation of TaxPub, an NLM DTD extension for domain-specific markup in taxonomy, from the experience of a biodiversity publisher [Internet]. Journal Article Tag Suite Conference (JATS-Con) Proceedings 2012 [Internet]. National Center for Biotechnology Information (US); 2012 [cited 2020 Sep 27]. Available from: <https://www.ncbi.nlm.nih.gov/books/NBK100351/>
23. Senderov V, Simov K, Franz N, Stoev P, Catapano T, Agosti D, et al. OpenBiodiv-O: ontology of the OpenBiodiv knowledge management system. J Biomed Semant [Internet]. 2018 Jan 18 [cited 2020 Mar 10];9(1):5. Available from: <https://doi.org/10.1186/s13326-017-0174-5>

24. Dimitrova M, Poelen J, Zhelezov G, Georgiev T, Agosti D, Penev L. Semantic Publishing Enables Text Mining of Biotic Interactions. *Biodivers Inf Sci Stand* [Internet]. 2020 Sep 28 [cited 2020 Sep 30];4:e59036. Available from: <https://biss.pensoft.net/article/59036/>
25. Poelen J, Upham N, Agosti D, Ruschel T, Guidoti M, Reeder D, et al. CETAF-DiSSCo/COVID19-TAF biodiversity-related knowledge hub working group: indexed biotic interactions and review summary [Internet]. Zenodo; 2020 [cited 2021 May 21]. Available from: <https://doi.org/10.5281/zenodo.3838240>
26. Olival KJ, Hosseini PR, Zambrana-Torrel C, Ross N, Bogich TL, Daszak P. Host and viral traits predict zoonotic spillover from mammals. *Nature* [Internet]. 2017 Jun [cited 2019 Sep 12];546(7660):646–50. Available from: <https://www.nature.com/articles/nature22975>
27. Nagaraj A, Shears E, Vaan M de. Improving data access democratizes and diversifies science. *Proc Natl Acad Sci* [Internet]. 2020 Sep 22 [cited 2020 Oct 4];117(38):23490–8. Available from: <https://www.pnas.org/content/117/38/23490>

Acknowledgments: We thank A. Casino, D. Koureas, and W. Addink for organizing the CETAF-DiSSCo COVID-19 Taskforce that resulted in this research. Illustrations were created in Inkscape (<http://www.inkscape.org/>) with images from BioRender.com and CC-BY licenses; Fig 2 reuses bat images from <http://zenodo.org/record/3756730> and Michigan Science Art with permission. E. Florsheim provided valuable conversations. **Funding:** Our efforts were supported by CETAF (Consortium of European Taxonomic Facilities), DiSSCo (Distributed System of Scientific Collections), the Biodiversity Knowledge Integration Center at Arizona State University (N.S.U.; President's Special Initiative Funds to N.M.F. and B.S.), the SYNTHESYS+ Research and Innovation action (Q.J.G., grant no. H2020-EU.1.4.1.2823827), the Arcadia charitable fund of Lisbet Rausing and Peter Baldwin (D.A.), the National Science Foundation Advancing Digitization of Biodiversity Collections Program DBI-1547229 (D.P.), the Special Research Fund of Hasselt University (M.P.M.V., BOF20TT06), and National Science Foundation award "Collaborative Research: Digitization TCN: Digitizing collections to trace parasite-host associations and predict the spread of vector-borne disease," Award numbers DBI:1901932 and DBI:1901926 (J.H.P.).

Author contributions: N.S.U, D.P., Q.J.G., N.B.S., J.H.P., and D.A. designed the conceptual arguments of this research. J.H.P., M.G., D.A., L.P., and N.S.U. worked on methods and created

software for data liberation. Curation of the resulting data was performed by J.H.P., D.A., and L.P., while N.S.U. and J.H.P. performed validations. N.S.U. and D.A. wrote the initial draft, N.S.U. created the figures with help from C.B.S., and all authors reviewed and edited the manuscript.

Competing interests: The authors declare no competing interests. **Data and materials**

availability: All data liberated as a result of these efforts are available at <https://doi.org/10.5281/zenodo.4068958>. **Search strategy and selection criteria:** Included

datasets were identified from 14 April to 6 October 2020 through CETAF-DiSSCo Taskforce activities and subsequently indexed by Global Biotic Interactions (GloBI, <https://globalbioticinteractions.org>). A full list of sources indexed through GloBI is provided with

the archived data at <https://doi.org/10.5281/zenodo.4068958>.