

that carry and transmit the infectious agents. Huge gaps exist in the information related to these vectors, creating an essential need for campaigns to mobilise and share data. The publication of data papers is an effective tool for overcoming this challenge. These peer-reviewed articles provide scholarly credit for researchers whose vital work of assembling and publishing well-described, properly-formatted datasets often fails to receive appropriate recognition. To address this, GigaScience's sister journal *GigaByte* partnered with the Global Biodiversity Information Facility (GBIF) to publish a series of data papers, with support from the Special Programme for Research and Training in Tropical Diseases (TDR), hosted by the World Health Organisation (WHO). Here we outline the initial results of this targeted approach to sharing data and describe its importance for controlling VBDs and improving public health.

Body Text

Free and open access to biodiversity data enables research and analysis needed to confront the threats and growing burden that vector-borne diseases and their control place on human health.

The World Health Organization (WHO)'s Global Vector Control Response (GVCR) 2017–2030 calls for additional efforts on data sharing on disease vectors, and Pillar 3 of the GVCR emphasizes vector surveillance and monitoring of interventions. Universal, free and open access to data on vectors will help countries to strengthen their response against vector-borne diseases and to improve health and well-being. The achievement of this goal requires development and management of databases within accessible platforms where all stakeholders, from researchers to those implementing vector control, can find supplementary information and experience. TDR, the special program for research and training on tropical diseases, hosted by WHO, is committed to helping researchers, especially from low and lower-middle income countries (LMICs), to have access to data-sharing platforms and to enhance their capacity to publish their data and thus make them accessible to the wider community. Data papers are a cost-efficient and effective tool to increase digital availability of relevant biodiversity data, and to mainstream data openness across research communities. Support of data papers on disease vectors through sponsorship of special issues is thus fully aligned with TDR's objectives.

The Global Biodiversity Information Facility (GBIF) has in recent years identified a number of priority research and policy areas where increased availability of biodiversity data would provide a richer evidence base. These include data-intensive, cross-disciplinary research into vectors and reservoirs of human diseases, for which critical taxonomic and geographic data gaps on the occurrence of species represent a significant obstacle. In 2020, GBIF formed a task group on mobilization and use of biodiversity data for research and policy on human diseases, with a mandate to provide advice, priority directions and expert opinions. Among the measures supported by this task group is promotion of and support for the publication of data papers, as a means of encouraging data sharing by a community of data holders largely unfamiliar with the process of publishing biodiversity datasets through GBIF. As data papers ideally describe well-prepared datasets already available in a platform such as GBIF, sponsored calls for submissions with article processing charges (APCs) waived for authors, provide a direct incentive and support for publishing datasets that are correctly formatted and of high quality. With support from TDR a first sponsored call for human disease vector data papers in *GigaByte* journal was announced in November 2021, and after going through peer review the resulting first phase of papers were all published by the end of May 2022 [1]. Many biodiversity data experts within the GBIF collaborative network, as well as task group members, provided direct support for preparation of the data papers submitted to *GigaByte*, helping to bridge the domains of biodiversity and biomedical research. Following positive reaction to these activities within the GBIF network and from TDR, plans are under development to scale up this initial experiment across taxa, diseases, regions and complexity of

biological systems. And a second call for papers for the *GigaByte* series has also just been reopened.

As well as being on hand for any papers that required curation and hosting of large supplemental files, the GigaScience Press GigaDB team were also on hand to assist the data peer review process. This involved data auditing and providing a data review for each submission, ensuring the data was open and FAIR (findable, accessible, interoperable and reusable). The review included selecting a number of data points to be carefully inspected, and verifying that the total number of occurrences and the geographic range were consistent with the details in the paper. To comply with *GigaByte*'s stringent open data policies, the journal insisted upon use of CC0 public domain waivers for the datasets described in the data papers. In some of the submissions, the review also picked up inconsistencies in the metadata caused by conversion problems and non-ASCII characters. In line with Open Peer Review, the peer reviews and data review templates are available for scrutiny via the Article Review History tabs on all of the papers.

The results of this first call for data papers have been promising and the results are highlighted here. The first 11 publications present more than 500 000 occurrence records linked to 675 000 specimens across over 50 countries (Fig. 1) [1]. The data described in the papers includes occurrence records (including specimen collections) published through GBIF.org, imaging data in the EBI Bioimaging repository, DNA barcodes in NCBI, and other miscellaneous data hosted in GigaDB and Dryad repositories. Publication in *GigaScience*'s sister journal *GigaByte* has enabled use of a new state-of-the-art XML-first publishing platform, and the papers include embedded dynamics, such as interactive maps and embedded protocols. Additionally, papers are linked to preprints and there are multilingual options for many papers that allow Portuguese and Spanish speakers to better understand the implications of important work relating to the public health of their communities. In the next section we outline this first phase of submissions, what they presented, and what has been learned from the process of publishing them.

Important stories that can be told through open data

As well as mobilizing a significant volume of data, the submissions represented diverse vectors, locations and forms. And while the datasets are presented primarily for re-use by others, an analysis reveals several interesting insights from the data itself.

The online catalogue of the Coleção de Flebotomíneos (FIOCRUZ/COLFLEB) has been derived from approximately 72 000 individual specimens and 370 species deposited at the René Rachou/Fiocruz Institute in Belo Horizonte, Brazil [2]. This dataset covers over 80 years of sandfly research in Brazil and 20 countries in the Americas, making it the largest and most comprehensive collection of these insects. Approaches to capturing this data has varied, with another paper also digitising observations of all published scientific studies on sandflies in the Brazilian state of Acre [3].

Data published from the *GigaByte Vectors* of human disease series includes:

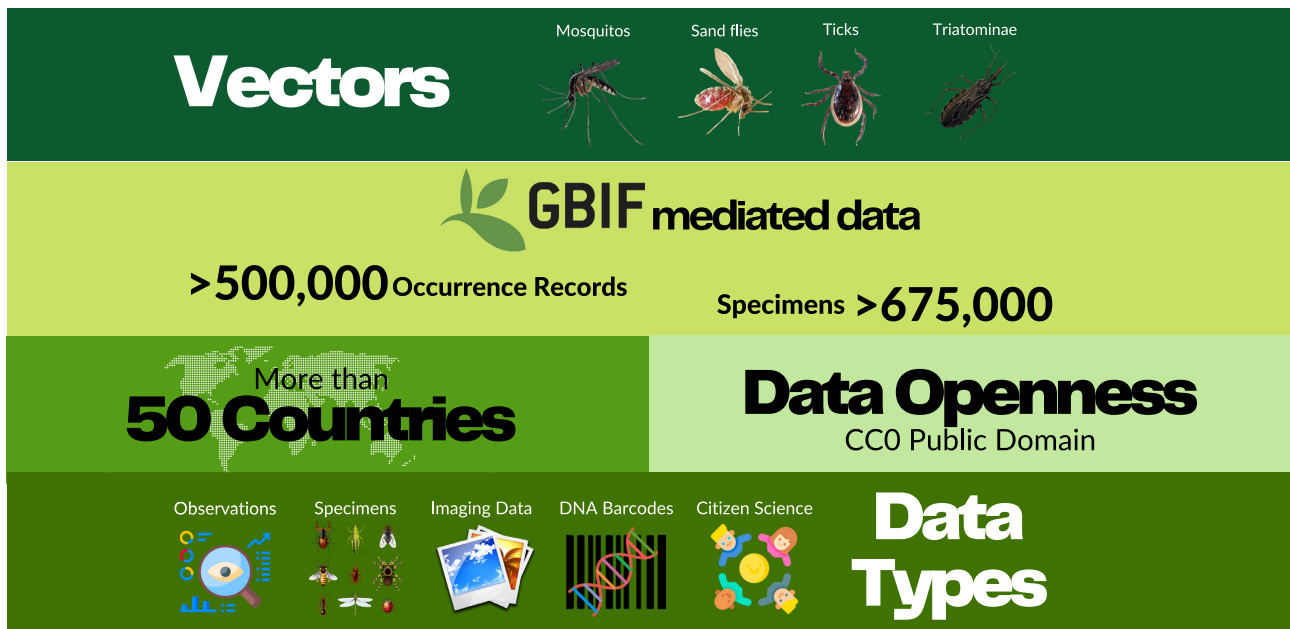


Figure 1: A summary of the disease vector data shared through the first phase of the sponsored call for data papers.

As well as national and international scale collections, papers also presented work collected in the lands of indigenous peoples in the Brazilian Amazon [4]. These sandfly vector records were obtained from areas of disease transmission where cutaneous leishmaniasis is endemic, and has grown with changes in the environment and hunting practices. The authors hope that these records will contribute to a better understanding of leishmaniasis transmission dynamics among these communities, as well as to increase data on the distribution of these insect vectors in locations that are remote and difficult to access, and that therefore are surveyed by public health systems.

The “Ana Leuch Lozovei” collection presents an incredible diversity of 100 species of Culicidae in 18 municipalities in Paraná state, Southern Brazil, collected between 1967 and 1999 [5]. It records three species for the first time in Brazil, signifying the expansion of geographical distribution of the species previously restricted to certain locations or countries.

The public health importance of this data has been very clear, such as a collection of data coming from the screening of urban households in three municipalities with a high incidence of dengue in Southwestern Colombia [6]. It presents novel data for the geographical distribution of 2383 specimens belonging to the Culicidae family, alongside the house infestation percentage per municipality and additional descriptive measures of the sampled mosquitos at each location. This type of data is not often reported due to the sampling effort involved in entomological sampling.

The MODRISK [7] and MEMO [8] projects presented the outputs of state-of-the-art monitoring of the exotic *Aedes* genus in Belgium. MODRISK uses a novel randomised approach to model mosquito biodiversity distribution at a 1-km resolution, based on longitudinal data, systematic screens and historical collections going back to 1878. MEMO (Monitoring of Exotic MOSquitoes in Belgium) looks at early detection of exotic mosquito species along

high-risk introduction routes in Belgium, where data is collected at defined points of entry. It also includes genetic sequencing data, as DNA-barcoding was used as a quality control step to validate 5% of the morphological identifications. This work showed that new exotic species could even be detected in temperate Belgium, such as *Ochlerotatus/Aedes koreicus*, which had not been reported in Europe until then.

In addition to national surveillance schemes, collaborative supranational projects are also represented in the series, with AIMSURV presenting data from the first pan-European harmonized surveillance of *Aedes* invasive mosquito species organized under the framework of the AIMCOST Action [9]. In 2020, AIMSURV was implemented by 42 teams from 24 countries. Data comprised a core file with 19 130 samples that improve knowledge of the European seasonal pattern of the *Aedes* invasive mosquito species *Ae. albopictus*, *Ae. japonicus* and *Ae. koreicus*.

Citizen scientists are accounted for, with the Mosquito Alert dataset including occurrence records of adult mosquitoes collected by citizens through the Mosquito Alert smartphone app [10]. Each record is linked to a photograph which is validated by entomological experts to assess the species. The paper shows that citizens can be part of a mature near-real-time surveillance system of targeted disease-vector mosquito species of concern in the EU. From a surveillance perspective, the system has been able to detect many appearances of *Aedes albopictus* well beyond its immediate expansion front. Another major highlight is the first detection in 2018 of *Aedes japonicus* in Spain, an isolated population located 1300 km away from its previously nearest known location in Europe.

Another project that includes citizen-collected data forms a sub-dataset of American triatomine, insect vectors involved in Chagas disease that are also known as “kissing bugs” [11]. With 90% of the US collected data obtained from the Kissing bugs and

