

# 1 Minimum Information for Reusable 2 Arthropod Abundance Data (MIReAAD)

3  
4 Myriad: a countless or extremely great number

5  
6 Samuel Rund\*, [srund@nd.edu](mailto:srund@nd.edu), VectorBase, Department of Biological Science, University of  
7 Notre Dame, IN, USA.

8  
9 Kyle Braak, [kyle.braak@gmail.com](mailto:kyle.braak@gmail.com), Global Biodiversity Information Facility (GBIF)  
10 Secretariat, Copenhagen, Denmark

11  
12 Lauren Cator, [l.cator@imperial.ac.uk](mailto:l.cator@imperial.ac.uk), Department of Life Sciences, Imperial College London,  
13 UK

14  
15 Kyle Copas, [kcopas@gbif.org](mailto:kcopas@gbif.org), Global Biodiversity Information Facility (GBIF) Secretariat,  
16 Copenhagen, Denmark

17  
18 Scott J. Emrich, [semrich@utk.edu](mailto:semrich@utk.edu), Department of Electrical Engineering and Computer Science,  
19 University of Tennessee, Knoxville, TN

20  
21 Gloria I. Giraldo-Calderón, [ggiraldo@nd.edu](mailto:ggiraldo@nd.edu), VectorBase - Bioinformatics Resource for  
22 Invertebrate Vectors of Human Pathogens, Department of Biological Science, University of  
23 Notre Dame, IN, USA.

24  
25 Michael A. Johansson, [mjohansson@cdc.gov](mailto:mjohansson@cdc.gov), Division of Vector-Borne Diseases, Centers for  
26 Disease Control and Prevention, 1324 Calle Cañada, San Juan, PR 00920; Department of  
27 Epidemiology, Harvard School of Public Health, 677 Huntington Ave, Boston, MA 02115

28  
29 Naveed Heydari, [naveedheydari@gmail.com](mailto:naveedheydari@gmail.com), Center for Global Health and Translational  
30 Science, State University of New York Upstate Medical University, Syracuse, NY  
31 544 S Vance St, Lakewood, CO, 80226, USA

32  
33 Donald Hobern, [dhobern@gbif.org](mailto:dhobern@gbif.org), Global Biodiversity Information Facility (GBIF) Secretariat,  
34 Copenhagen, Denmark

35  
36 Sarah A. Kelly, [s.kelly@imperial.ac.uk](mailto:s.kelly@imperial.ac.uk), VectorBase, Vector Immunogenomics and Infection  
37 Laboratory, Department of Life Sciences, Imperial College London, UK

38  
39 Daniel Lawson, [daniel.lawson@imperial.ac.uk](mailto:daniel.lawson@imperial.ac.uk), VectorBase and Vector Immunogenomics and  
40 Infection Laboratory, Department of Life Sciences, Imperial College London, UK

42 Cynthia Lord, [clord@ufl.edu](mailto:clord@ufl.edu), Florida Medical Entomology Lab, University of Florida-IFAS,  
43 Vero Beach, FL

44

45 Robert M MacCallum, [r.maccallum@imperial.ac.uk](mailto:r.maccallum@imperial.ac.uk), VectorBase and Vector Immunogenomics  
46 and Infection Laboratory, Department of Life Sciences, Imperial College London, UK

47

48 Dominique G. Roche, [dominique.roche@mail.mcgill.ca](mailto:dominique.roche@mail.mcgill.ca), Institute of Biology, University  
49 of Neuchâtel, 2000, Neuchâtel, Switzerland

50

51 Sadie J. Ryan, [sjryan@ufl.edu](mailto:sjryan@ufl.edu), Quantitative Disease Ecology and Conservation Lab, Department  
52 of Geography, University of Florida, Gainesville, FL 32601 USA; Emerging Pathogens Institute,  
53 University of Florida, Gainesville, FL 32610 USA; College of Life Sciences, University of Kwa-  
54 Zulu Natal, Durban, South Africa

55

56 Dmitry Schigel, [dschigel@gbif.org](mailto:dschigel@gbif.org), Global Biodiversity Information Facility (GBIF) Secretariat,  
57 Copenhagen, Denmark

58

59 Kurt Vandegrift, [kjv1@psu.edu](mailto:kjv1@psu.edu), Center for Infectious Disease Dynamics, Department of  
60 Biology, The Pennsylvania State University, 16801, PA, USA

61

62 Matthew Watts, [m.watts@imperial.ac.uk](mailto:m.watts@imperial.ac.uk), Department of Life Sciences, Imperial College  
63 London, UK

64

65 Jennifer M. Zaspel, [zaspelj@mpm.edu](mailto:zaspelj@mpm.edu), Department of Zoology, Milwaukee Public Museum, 800  
66 W Wells Street, Milwaukee, WI, 53233, USA

67

68 Samraat Pawar\*, [s.pawar@imperial.ac.uk](mailto:s.pawar@imperial.ac.uk), Department of Life Sciences, Imperial College  
69 London, Silwood Park Campus, Buckhurst Road, Ascot, Berkshire SL5 7PY, United Kingdom

70

71 \* Corresponding authors

## 72 Abstract

73 Arthropods play a dominant role in natural and human-modified terrestrial ecosystem dynamics.

74 Spatially-explicit population time-series are crucial for statistical or mathematical models of

75 these dynamics and assessment of their veterinary, medical, agricultural, and ecological impacts.

76 Arthropod data have been collected world-wide for over a century, but remain scattered and

77 largely inaccessible. With the ever-present and growing threat of arthropod vectors of infectious

78 diseases and pest species, there are enormous amounts of historical and ongoing surveillance.  
79 These data are currently reported in a wide variety of formats, typically lacking sufficient  
80 metadata to make reuse and re-analysis possible. We present the first minimum information  
81 standard for arthropod abundance. Developed with broad stakeholder collaboration, it balances  
82 sufficiency for reuse with the practicality of preparing the data for submission. It is designed to  
83 optimize data (re-)usability from the “FAIR,” (Findable, Accessible, Interoperable, and  
84 Reusable) principles of public data archiving (PDA). This standard will facilitate data unification  
85 across research initiatives and communities dedicated to surveillance for detection and control of  
86 vector-borne diseases and pests.

## 87 Introduction

88 Arthropods play a dominant role in the dynamics of practically all natural and human-modified  
89 terrestrial ecosystems<sup>1-3</sup>, and have significant economic and health effects. For example, certain  
90 insects provide significant economic benefits (*e.g.* pollination) exceeding \$57 billion a year to  
91 the United States alone<sup>4</sup>. Meanwhile, invasive insects cost an estimated \$70 billion dollars per  
92 year globally<sup>5</sup> and insect pests may reduce agricultural harvests by up to 16%, with an equal  
93 amount of further losses of harvested goods<sup>6</sup>. Particularly noteworthy is a subset of arthropods  
94 that are disease vectors, transmitting pathogens to and between animals as well as plants. Vector-  
95 borne diseases cause billions of dollars in crop and livestock losses, every year<sup>7-9</sup>. In humans,  
96 vector borne diseases account for more than 17% of all infectious diseases (*e.g.* malaria, Chagas,  
97 dengue, and leishmaniasis, Zika, West Nile, Lyme disease, and sleeping sickness), with hundreds  
98 of thousands of deaths, hundreds of millions of cases, and billions of people at risk, annually<sup>10,11</sup>.

99

100 The current economic and health burden of arthropod pests, exacerbated by invasive species, and  
101 uncertain effects of climate change<sup>12,13</sup>, has driven significant research programs and data  
102 collection efforts. These include crop pest, mosquito, and tick survey and reporting initiatives<sup>14-</sup>  
103 <sup>18</sup>, citizen science projects<sup>19-21</sup>, and digitization of museum specimen data<sup>22,23</sup>, all yielding a rich  
104 and growing trove of field-based data spanning multiple spatial and temporal scales. Monitoring  
105 arthropod abundance (*e.g.* Figure 1) in different disciplines (*e.g.*, biodiversity research, pest-  
106 control assessment, vector-borne disease monitoring, or pollination research) uses similar  
107 techniques, with similar objectives: to quantify abundance, phenology and geographical ranges  
108 of target arthropod species. Despite a growing number of data collections, they are often not  
109 reusable, or comparable to similar data, due to a lack of standardization and metadata. In  
110 contrast, the advent of the deposition of data from high-throughput technologies (*e.g.* NCBI and  
111 GenBank), data and code sharing, and other practices to improve transparency and reusability of  
112 research results are increasing rapidly across the sciences<sup>24-29</sup>. Furthering these advances through  
113 standardization and public archiving of arthropod abundance data can bring significant benefits,  
114 including (1) supporting empirical parameterization and validation of mathematical models (*e.g.*  
115 of pest or disease emergence and spread), (2) validation of model predictions, (3) reduction in  
116 the duplication of expensive empirical research, and (4) revealing new patterns and questions  
117 through meta-analyses<sup>30-33</sup>. This will also lead to substantial public benefit through improved  
118 human, animal, plant, and ecosystem health, and reduced economic costs.

119

120 A key impediment to the re-use of these data is the lack of adequate metadata or data descriptors  
121 (*i.e.* data about the data)<sup>34-37</sup>. In general, for data to be most valuable to the scientific community,

122 they should meet the FAIR Principles – they should be Findable, Accessible, Interoperable and  
123 Reusable – and delineate the key components of good data management and stewardship  
124 practices<sup>38,39</sup>. Data are Findable and Accessible when they are archived and freely downloadable  
125 from an online public data repository that is indexed and easily searchable. Interoperability and  
126 reusability describe the ease with which humans or computer programs can understand the data  
127 (*e.g.* via metadata) and explore/re-use them across a variety of non-proprietary platforms. Even  
128 when data are available, metadata for arthropod abundance data are often absent or not readily  
129 interpretable, limiting their reusability at a fundamental level.

## 130 Results

### 131 **A minimum information standard for arthropod abundance data**

132 Here, we present a Minimum Information for Reusable Arthropod Abundance Data (MIReAAD)  
133 standard for reporting primarily longitudinal (repeated, temporally explicit) field-based  
134 collections of arthropods. In the same manner as has been developed in other biological  
135 disciplines<sup>40-45</sup>, this standard is “minimum” because it defines the necessary minimal information  
136 required to understand and reuse a dataset without consulting any further text, materials, or  
137 methods<sup>46</sup>. MIReAAD is designed to facilitate data archiving efforts of publishers and field  
138 researchers. It is not a data model and therefore does not define controlled vocabularies, or  
139 specific field titles, but should be easy to understand, and interpret by the wider scientific  
140 community<sup>46</sup>.

141  
142 The minimal standards are separated into two components, metadata and data. For each  
143 component, we provide a description of the information that should be included,

144 recommendations for how to make that information as useful as possible, and examples. The  
145 metadata component (Table 1) includes information for the origin of the data set (*e.g.* study  
146 information and licensing for usage). The second component (Table 2) lists and describes  
147 specific data fields that should be included in data collection sheets. We also provide  
148 recommendations and examples to demonstrate how these recommendations can be  
149 implemented. MIREAAD was designed to match the data that are generally collected by  
150 academic researchers and surveillance initiatives, and can serve as a checklist for important  
151 information that needs to be recorded but is often unintentionally omitted (*e.g.* Figure 2A). By  
152 adhering to MIREAAD standards, omissions and ambiguity can be avoided even if the data are  
153 shared in different formats (Figure 2B and C). Finally, we identify common problems likely to  
154 be encountered across all the MIREAAD metadata and data fields, and data quality standards that  
155 can be employed to avoid confusion (Box 1).

156

### **Box 1. Data quality standards**

**No abbreviations.** Abbreviations (including in column names) are ambiguous, with the exception of measurement units (*e.g.* centigrade and meters).

**No external legend/key files.** While repetitive, all data should be explicitly given within the data table. Separate files mapping ID numbers to GPS locations, full species names, etc., should be avoided. In addition, rich metadata is essential for good data discovery and reuse.

**Unambiguous dates.** Because of country-level differences in date formats, data should be reported with 4 digit years, and months provided alphabetically and not numerically (*e.g.* 4-Jun-2017 or Nov 12, 2015).

**Machine-readable file formats.** Data should be provided in non-proprietary machine readable formats such as comma-separated text files. PDFs and multiple spreadsheets in the same document should be avoided.

**No font styling or subsection headings.** Formatting (color, bold, italics, subscripts, sheet tab names, *etc.*) should not be required for understanding the data. Subsection headings should not be required to understand data; every line of data should be interpretable in isolation from any other line of data.

**Highest precision possible.** Data should be provided at the highest temporal, spatial, numerical, and taxonomic resolution available. If location (*e.g.*, geographical coordinate) data need to be presented at a lower resolution than available for privacy reasons, this should be made clear in the submission in Study Information (Resource Metadata; Table 1).

**Language.** Once data are ready to be deposited/submitted, all fields and data are preferably written in English. This will allow researchers and data curators worldwide to understand and reuse the data. Use of other languages is better than not publishing data. Please avoid introducing data reuse barriers through incomplete translation. For example, non-English field names in an English-language submission.

157

## 158 **Examples**

159 Below we provide three examples to illustrate MIREAAD compliant data (linked to  
160 Supplemental Data Files 1-4, respectively). Researchers can use these data sheets as a basis for  
161 formatting their own data. In these examples, note that all data meet the data quality standards of  
162 Box 1; are adequately described, have columns labeled, *etc.* to eliminate ambiguity (even if the  
163 data appear repetitive; for example, the sex and life stage are repeated in every row). Examples 1  
164 and 2 should be sufficient for most data generators. Example 3 (Data Files 3-4) demonstrates a  
165 more complex data collection scenario.

166

167 1. *Long-format trapping data.* Each row captures count data for a single species' occurrence in a  
168 given sampling event. This illustrates an example of the most common mosquito collection



169 protocol. [\[Sup Datasheet 1\]](#). Also see Figure 2B.

170

171 2. *Wide format trapping data*. Each row captures count data from a given sampling event. Each  
172 identified taxonomic group is identified in a separate column. An ‘additional sample  
173 information’ field, ‘sub-location,’ has been added to describe the various locations around the  
174 village where collections were made. [\[Sup Datasheet 2\]](#). This illustrates an example of adult  
175 mosquito populations that have been tracked over time and in specific locations. Also see Figure  
176 2C.

177

178 3. *Complex trapping data scenario*. Tick surveillance performed using tick drags and flags and  
179 collections of ectoparasites on trapped mice. The tick drags/flags report three life stages  
180 independently (adult, larvae, and nymph) [\[Sup Datasheet 3\]](#) . Larvae are only identified to the  
181 genus, while adults and nymphs are identified to the species. Observations of different life  
182 stages and sexes are preferably documented in separate records. A Sample Name is used to help  
183 link these records (but would not be necessary.) The mouse survey uses an additional sample  
184 information field to record the sex of the trapped mouse from which the parasites were collected  
185 [\[Sup Datasheet 4\]](#).

## 186 Discussion

### 187 **MIReAAD as the path to FAIR data principles**

188 We designed MIReAAD to achieve a balance between standards that are too onerous for data  
189 generators and standards that are sufficient to ensure at least minimal reusability<sup>31,40</sup>. Like all

190 minimum standards, MIREAAD only aims at ensuring data ‘Reusability’. However, ultimately  
191 this will promote the implementation of data models — the explicit definition of data field  
192 names, data formats (*e.g.*, for dates and GPS locations), and controlled vocabularies (*e.g.*, the  
193 Darwin Core<sup>47</sup>). Data models enable ‘Interoperability’, and in turn facilitate structured databases,  
194 public repositories, and development of data analysis tools<sup>46,48</sup>. Deposition in open databases  
195 make data ‘Findable’ and ‘Accessible’<sup>49–51</sup>. MIREAAD compliant data contain sufficient  
196 information for established aggregators/databases such as VectorBase and SCAN (Symbiota  
197 Collections of Arthropods Network<sup>52</sup>) to process and store the data in a standardized data model  
198 [*e.g.*, Darwin Core, a widely used universal data standard that supports opportunistic observation  
199 and collection data (occurrence core) as well as presence/absence and abundance data collected  
200 using strict and documented methodology (event core)<sup>47</sup>], and ultimately facilitate data transfer  
201 to even more comprehensive biodiversity databases [*e.g.* GBIF, which contains over one billion  
202 species occurrence records, from thousands of environmental, ecological, and natural resource  
203 investigations, including research on Arthropoda in numerous ecological and monitoring  
204 projects, allowing for study of changes and trends in populations.<sup>51</sup>]. Indeed, in Supplemental  
205 File 5, we provide an example of the mapping of data fields from this minimum information  
206 standard, to DarwinCore and GBIF. In this way, MIREAAD opens the door to FAIR data and  
207 more sophisticated methods to integrate data across many scales.

208

### 209 **Benefits to field researchers**

210 It is essential that the benefits of a minimal data standard extend not just to data re-users, but also  
211 to the researchers who collect and generate data in the first place. MIREAAD provides a  
212 framework for data preparation that can help scientists achieve recognized professional merit for

213 sharing data such as increased citation rates, academic recognition, opportunities for co-  
214 authorship, and new collaborations [sensu Roche et al. 2014<sup>31</sup>]. Large, deposited data sets can  
215 now themselves be standalone, citable “data papers” (*e.g.* <sup>53–55</sup>) or even depositions without any  
216 traditional manuscript (but as an authored ‘digital product,’ with persistent identifiers, such as a  
217 DOI number), if desired. Data sets are increasingly recognized as valuable research outputs that  
218 count towards academic recognition and professional advancement (*e.g.* grants, interviews, and  
219 tenure). For example, several funders (*e.g.* United States National Science Foundation and Swiss  
220 National Science Foundation) have adopted or are in the process of adopting the Declaration on  
221 Research Assessments (DORA)<sup>56</sup>, offering further opportunities for data generators to gain  
222 recognition and publication credit for their work<sup>57</sup>. Also, an increasing number of funders are  
223 mandating public data access, and detailed data management plans are often required even at the  
224 grant proposal stage. Therefore, reporting data according to MIREAAD will provide a  
225 foundational pipeline for stipulating archival formats.

226  
227 Furthermore, many data generators are also data users. Developing analyses that rely on  
228 standardized fields can facilitate the development of generalized analytical tools that can be  
229 easily extended to datasets beyond those that were collected by a single individual or lab. In this  
230 way, they can enable extensions of work that would otherwise not happen, such as comparisons  
231 of population dynamics in different locations or assessments of interspecies interactions.  
232 Adopting MIREAAD therefore can both help data generators reap the benefits of sharing data  
233 they have collected and enable them to more readily leverage data collected by others.

234

235 **Further MIREAAD applications and extensions**

236 The creation of minimum information standards for these types of databases facilitates analyses  
237 of data at the scales that cannot be attained by a single individual or lab group. Linking records  
238 to additional information also extends the utility of these data to address population level  
239 questions. For example, a well-populated database presents opportunities to investigate  
240 interactions between populations of different species of arthropod that overlap in geography, but  
241 may be of interest individually to different realms of research. As a case in point, in the  
242 northeastern USA, *Agrilus plannipennis*, the Emerald Ash Borer (EAB), is a highly destructive  
243 invasive insect, monitored closely by both state and federal agencies for management<sup>58</sup>.  
244 Interestingly, EAB are creating lots of new habitat for carpenter bees, a species interaction that  
245 can be tracked and anticipated using large scale arthropod data.

246

247 Another example of the utility of linked data is for disease vectors. Data on insecticide resistance  
248 linked with time and place would be valuable for coordinating control strategies within and  
249 between nations and communities. Presence/absence data on infection levels would be helpful  
250 for tracking and investigating disease outbreaks, and dynamics. Standardization of these data  
251 would be particularly useful for pathogens that infect multiple vectors and hosts and would  
252 facilitate a “One Health” approach. Other important vector phenotypes that contribute to control  
253 and transmission such as pathogen susceptibility, biting preferences, and breeding behaviours  
254 could be measured over time and space.

255

256 We note that MIRreAAD is applicable not only to abundance measurements, but could be easily  
257 extended to any other kind of routinely sampled time-series field data. For example, in addition

258 to aphid abundance, plant pathogen (such as mosaic virus) infection and insecticide resistance  
259 statuses of the aphids could be reported in MIRreAAD format.

## 260 Conclusion

261 We present MIRreAAD, a minimum information standard for representing arthropod  
262 abundance data. MIRreAAD will facilitate collation and analyses of data at scales that cannot be  
263 attained by a single individual or lab, to address key questions across temporal and spatial scales,  
264 such as within and across-year phenology of abundance of target arthropod taxa over large  
265 geographical areas. This is particularly important given the pressing need to understand and  
266 predict the population dynamics of harmful (e.g., disease vectors and pests) as well as beneficial  
267 (e.g., pollinators, bio-control agents) arthropods in natural and human modified landscapes. This  
268 is the first step for achieving the broad benefits of FAIR data for arthropod abundance. We call  
269 on data generators, authors, reviewers, editors, journals, research infrastructures (e.g. data  
270 repositories) and funders to embrace MIRreAAD as a standard to facilitate FAIR data use and  
271 compliance for arthropod abundance data.

272  
273 **Table 1. The MIRreAAD Study Information (Resource metadata) fields.** The information in this  
274 table should be included with every data submission, for example by including data in the file  
275 header as demonstrated in Data Files 1-4.

276

Field	Details	Recommendations	Examples
-------	---------	-----------------	----------

Contact details	A name, person, authority, etc. that may be contacted with enquiries about the data.	Include investigator ORCID(s), email address, website (if institutional) if possible.	Kurt Vandegrift <a href="https://orcid.org/0000-0002-5690-3300">orcid.org/0000-0002-5690-3300</a> <a href="mailto:kurtvandegrift@gmail.com">kurtvandegrift@gmail.com</a>  State University Agricultural Extension John Smith ( <a href="mailto:jsmith@StateU.edu">jsmith@StateU.edu</a> ) <a href="http://www.StateU.edu/AgriculturalExtension/">www.StateU.edu/AgriculturalExtension/</a>
General description of the experiment/ collection set	A short description of the study objectives, sampling design, and hypotheses.  Used to aid in browsing multiple studies.  A short title and long form name might be helpful.	Useful things to indicate are:  Random sampling or continuous monitoring in fixed locations  General time frames and location.  General description of where data is from.	"Long term, fixed trapped, municipal surveillance of west Nile vector population in Colorado from 2000-2010" ----- "Pennsylvania <i>Ixodes scapularis</i> weekly abundance"  Continuous (weekly) monitoring of tick numbers attached to White-footed mice in fixed locations in Pennsylvania, USA (12 sites). 2003-present." ----- "Long term aphid emergence monitoring using continuous suction traps"
Citations	Reference to related publications, digital if possible (e.g. DOI(s) or PMID(s)).		"A web-based relational database for monitoring and analyzing mosquito population dynamics Sucaet Y, Van Hemert J, Tucker B, Bartholomay L."  "PMID: 18714883"  Horiuchi, Kaho, Kosei Hashimoto, and Fumio Hayashi. "Cantharidin world in air: Spatiotemporal distributions of flying canthariphilous insects in the forest interior." Entomological Science (2018).
Species Identification Method	A description of method of species identification. Particularly important for cryptic species complexes.		"Morphological"  "Genotyped, using method of Smith et al 2014, PMID: 18714883"
Not present vs zero information	Indication of what gaps, zeros, NA, etc mean.	It is imperative, especially for population surveys, to understand the difference between a species was not found when the collection method would be expected to find the given species (confirmed absence) or a species was not	"Zero indicates was looked for and not found. NA represents a trap failure etc"

		<p>looked for (e.g. a trap failure)</p> <p>Preferably, a zero indicates was looked for and not found, and a NA represents was not looked for/trap failure/ etc.</p> <p>Blank values are discouraged</p>	
GPS obfuscation information	<p>If GPS data obfuscation (e.g. GPS points are intentionally offset from their actual locations) or de-resolution occurs (e.g. GPS precision is intentionally reduced) , a statement on the manner by which this occurred.</p>	<p>The highest resolution data (e.g. trap-level, specific GPS location) are the most useful. It is hoped that no data obfuscation / de-resolution occurs</p>	<p>"GPS locations have been truncated to 3 decimals"</p> <p>"GPS locations obfuscated using N-Dispersion"</p> <p>"No GPS deresolution was performed"</p>
Data usage information	<p>The data reuse policy for your data.</p> <p>Please provide a creative commons license identification.</p> <p>See <a href="https://creativecommons.org">https://creativecommons.org</a> for more information.</p>	<p>For data to be F.A.I.R., it must be Reusable. We therefore recommend data be provided as "CC0" or "CC BY 4.0".</p> <p>"CC0", under which data are made available for any use without restriction or particular requirements on the part of users</p> <p>"CC BY 4.0", under which data are made available for any use provided that attribution is appropriately given for the sources of data used, in the manner specified by the owner (e.g. citation).</p>	<p>"CC0"</p> <p>or</p> <p>"CC BY 4.0"</p>

277

278

279 **Table 2.** The MIR<sub>e</sub>AAD data fields. Fig 1B provides an annotated example.

280

Field(s)	Details	Recommendations	Examples
<b>Start Time (for collection)</b>	<p>Start time of the data sample collection.</p> <p>e.g. The trap was set...</p>	<p>Be as specific as <i>practically</i> possible.</p> <p>Any unambiguous format is acceptable. However, do not use two-digit year abbreviations.</p> <p>If relevant, provide timezone in field or in header, a 24 hour clock is preferred, but should be made unambiguous as to which time format is being used.</p>	<p>"2012-04-27"</p> <p>"July 26, 2017"</p> <p>"2017-Jul-26"</p> <p>"2017-July-26 Morning "</p> <p>"2017-Jul-26 20:00 GMT "</p>
<b>End Time (for collection)</b>	<p>End time of the data sample collection.</p> <p>e.g. The trap was collected...</p>	<p>See above.</p> <p>If instantaneous data collection (e.g. a tick drag), End Time may be the same as Start Time.</p>	<p>See above.</p>
<b>Location</b>	<p>The geographical location of sample collection.</p>	<p>As detailed as possible. Latitude and longitude if possible with specified accuracy Providing <i>both</i> a GPS point (decimalized GPS points are preferred) field and a geographical name field is preferred.</p> <p>Note only providing location <i>names</i> is highly discouraged as they change over time and can be ambiguous. Place / Trap names and GPS fields can be provided.</p> <p>If obfuscation was used, it should be indicated in the Metadata (Table 1).</p> <p>Splitting latitude and longitude further into two columns further reduces ambiguity.</p>	<p>"Kukar Maikiya, Jigawa State, Nigeria"</p> <p>"40.697" and "-74.015"</p>



<p><b>Collection method</b></p>	<p>Sampling apparatus (e.g. trap type, observation method)</p>		<p>“CDC light trap”  “Tick drag”  “Quadrat count”  “BG Sentinel Trap”  “Pitfall trap”  “Larval dip”  “Johnson suction trap”  “Lindgren Funnel Trap”</p>
<p><b>Collection attractants</b></p>	<p>The attractant/ lures used to attract insects to a trap or collection</p>		<p>“None”  “Carbon dioxide”  “UV light”  “BG-Sweetscent Mosquito Lure”  “Human/animal bait”</p>
<p>Collection area</p>	<p>The spatial extent (area or volume) of the sample.</p>	<p>If relevant (e.g., when collection method is transect or quadrat), in units of area or volume, the spatial coverage of the sampling unit</p> <p>Note this field would not typically be used for mosquito collections.</p>	<p>“100 m<sup>2</sup>”  “1 liter”  “1 ha”  “10m<sup>3</sup>”</p>
<p><b>Taxonomy</b></p>	<p>Classification of sample collected.</p>	<p>Scientific genus and species preferred.</p> <p>Avoid abbreviation.</p>	<p>“<i>Ixodes scapularis</i>”  “<i>Aedes aegypti</i>”  “<i>Anopheles gambiae sensu stricto</i>”</p>

<p><b>Unit(s) of measurement and observation</b></p>	<p>Description of exactly what was observed, the unit for "Value" below.</p> <p>For counts, should indicate life stage, sex, etc.</p> <p>Unit measures can be encoded into value field header. Consider multiple unit fields (e.g. separate fields for sex and stage.) See Figure 2.</p>	<p>Do not abbreviate.</p> <p>Coded data key should be provided in field name (e.g. "1 = species present 0= species absent")</p>	<p>"Number of individuals per m<sup>2</sup>"</p> <p>"Female" and "Adult"</p> <p>"Male and Female" and "Nymphs"</p>
<p><b>Value</b></p>	<p>The numerical amount or result from the sample collection.</p> <p>Often this will be a quantity of observed individuals. Unit measures can be encoded into value field header. See Figure 2.</p>	<p>Units should be provided in a separate field.</p>	<p>"0"</p> <p>"23"</p> <p>"Yes"</p> <p>"Not present"</p>
<p>Additional sample information</p>	<p>This could be more than one field and should be used when more information is required to understand the experiment, for example experimental variables, sub-locations, etc.</p> <p>Some users may report wind speeds, temperatures, elevations etc.</p>	<p>Do not abbreviate.</p>	<p>"Forest" vs "Field"</p> <p>"Winter" vs "Summer"</p> <p>"Inside" vs "Outside"</p> <p>"200 meters above sea level"</p>

Sample Name	A human readable sample name.  May exist solely for the benefit of the depositor in organizing their data, use their own internal naming conventions etc.  May also be used to tie related observations together.	Naming convention is not restricted, but any encoded metadata should be revealed in the other datafields. For example, you may name a sample named 'Aphid1_StickyTrap_Jan4,' but you will still have "Sticky Trap" listed in a Collection Method field, and "Jan 4, 2017" in the date field.	"Trap1_Night1" "KissingBug_2" "00004" "Jan08_animal_4,"
-------------	---	--	--

281 Field names in bold should be considered also required. Remaining fields are optional or  
282 depend on the complexity of the experimental design  
283

## 284 Author contributions

285 The project was conceptualized by Lauren Cator and Samraat Pawar. The original draft was  
286 prepared by Michael A. Johansson, Samuel S.C. Rund, Naveed Heydari, Kurt Vandegrift,  
287 Matthew Watts, and Samraat Pawar. Visualization was prepared by Kurt Vandegrift, Samuel  
288 S.C. Rund, Samraat Pawar, and Michael A. Johansson. Review & Editing was performed by all  
289 the authors.  
290

## 291 Competing interest statement

292 The authors declare no competing interests.

## 293 Acknowledgements

294  
295 The seeds of this effort were planted in 2016 at a meeting of VectorBiTE, which is a cross-  
296 disciplinary research coordination network (RCN) for disease vectors. Samuel S.C. Rund,  
297 Matthew Watts, Kurt Vandegrift, Naveed Heydari, Cynthia Lord, Michael Johansson, Samraat  
298 Pawar, and Sadie J. Ryan, received travel funding from NIH grant 1R01AI122284-01 and  
299 BBSRC grant BB/N013573/1 as part of the joint [NIH-NSF-USDA-BBSRC] Ecology and  
300 Evolution of Infectious Diseases program.

301  
302 Samuel S.C. Rund was funded by the Royal Society (NF140517). Rund, Daniel Lawson, Robert  
303 M. MacCallum, Sarah A. Kelly, Gloria I. Giraldo-Calderon and Scott J. Emrich were supported  
304 by the National Institute of Allergy and Infectious Diseases, National Institutes of Health,  
305 Department of Health and Human Services, under Contract No. HHSN272201400029C  
306 (VectorBase Bioinformatics Resource Center).

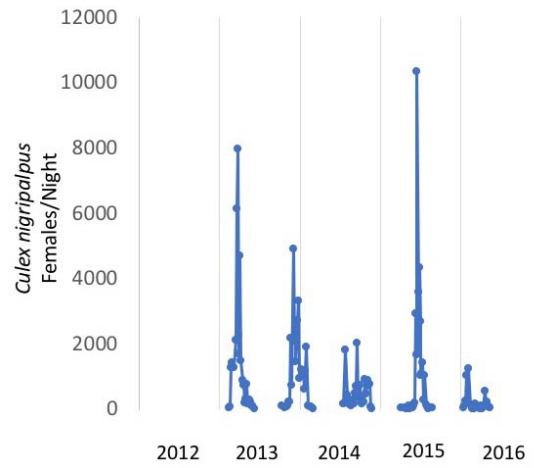
307  
308 Kurt Vandegrift was funded by the National Science Foundation Ecology and Evolution of  
309 Infectious Diseases program (1619072).

310  
311 Naveed Heydari and Sadie J. Ryan were funded by National Science Foundation (NSF DEB  
312 EEID 1518681).

313  
314 Sadie J. Ryan was additionally funded by NIH 1R01AI136035-01, and CDC grant  
315 1U01CK000510-01: Southeastern Regional Center of Excellence in Vector-Borne  
316 Diseases: the Gateway Program. This publication was supported by the Cooperative Agreement  
317 Number above from the Centers for Disease Control and Prevention. Its contents are solely the  
318 responsibility of the authors and do not necessarily represent the official views of the Centers for  
319 Disease Control and Prevention.

320  
321 Jennifer M. Zaspel was funded by the National Science Foundation Division of Biological  
322 Infrastructure (NSF 1561448, NSF 1601957).

323 **Figures**



324 **Figure 1. Example population abundance time-series.**  
325  
326

A.

Date	Trap	Count	Species
12/11/18	LT-5	6	C. inornata
12/11/18	LT-6	4	C. inornata
12/11/18	LT-7	6	C. inornata
12/11/18	LT-5	8	A. francisc.
12/11/18	LT-6	3	A. francisc.
12/11/18	LT-7	11	A. francisc.

Annotations for A:

- What is being counted? Adults? Females?
- Is this date of collection? How long was trap set?
- Ambiguous date format: Different countries use alternate day, month, and year orders
- Ambiguous genus: *Culex inornata* or *Culiseta inornata*?
- Ambiguous species: *Anopheles franciscanus* or *Anopheles Francisci*?
- Where is this trap located? What kind of trap? "LT" could be "light trap," but there are multiple kinds. Was an attractant used?

B.

Trap Set	Trap Collected	Collection Method	Attractants	Count	Life stage	Sex	Trap ID	Trap Location	Latitude	Longitude	Species
11-Dec-2018	13-Dec-2018	CDC Light Trap	Light, CO2	6	Adult	Female	LT-5	High School	22.211	84.974	Culex inornata
11-Dec-2018	13-Dec-2018	CDC Light Trap	Light, CO2	4	Adult	Female	LT-6	High School	22.209	84.894	Culex inornata
11-Dec-2018	13-Dec-2018	CDC Light Trap	Light, CO2	6	Adult	Female	LT-7	High School	22.199	85.012	Culex inornata
11-Dec-2018	13-Dec-2018	CDC Light Trap	Light, CO2	8	Adult	Female	LT-5	High School	22.211	84.974	Anopheles franciscanus
11-Dec-2018	13-Dec-2018	CDC Light Trap	Light, CO2	3	Adult	Female	LT-6	High School	22.199	85.012	Anopheles franciscanus
11-Dec-2018	11-Dec-2018	CDC Light Trap	Light, CO2	11	Adult	Female	LT-7	High School	22.199	85.012	Anopheles franciscanus
11-Dec-2018	11-Dec-2018	CDC Light Trap	Light, CO2	1	Adult	Male	LT-7	High School	22.199	85.012	Anopheles franciscanus

Annotations for B:

- The time period of collection is unambiguous
- What was counted or measured is clear
- Taxonomic ambiguity avoided by not using abbreviations.
- Unambiguous date format
- The type of trap and attractants used is clear
- Human readable (but ambiguous) place names are made unambiguous with GPS coordinates

C.

Trap Set	Trap Collected	Collection Method	Attractants	Trap ID	Latitude	Longitude	Culex inornata (Adult Female Count)	Anopheles franciscanus (Adult Female Count)
11-Dec-2018	13-Dec-2018	CDC Light Trap	Light, CO2	LT-5	22.211	84.974	3	5
13-Dec-2018	15-Dec-2018	CDC Light Trap	Light, CO2	LT-5	22.211	84.974	4	4
15-Dec-2018	17-Dec-2018	CDC Light Trap	Light, CO2	LT-5	22.211	84.974	5	6
17-Dec-2018	19-Dec-2018	CDC Light Trap	Light, CO2	LT-5	22.211	84.974	3	8
19-Dec-2018	21-Dec-2018	CDC Light Trap	Light, CO2	LT-5	22.211	84.974	4	2

Annotation for C:

- Data reported in wide format can still be unambiguous by including relevant metadata in column headings

327  
328  
329  
330  
331  
332

**Figure 2. MIREAAD reduces data ambiguity.** A. Seemingly clean data can still lack key information or have ambiguous metadata, hindering data reuse. B. MIREAAD compliant data includes the metadata necessary for data reuse and removes ambiguity. C. Note data can be formatted differently, but still be MIREAAD compliant, such as by presenting data in a wide format

## 333 References

- 334 1. Seastedt, T. R. & Crossley, D. A. The Influence of Arthropods on Ecosystems. *Bioscience*  
335 **34**, 157–161 (1984).  
336 2. Arthropod Regulation of Micro- and Mesobiota in Below-Ground Detrital Food Webs |

- 337 Annual Review of Entomology.
- 338 3. Whiles, M. R. & Charlton, R. E. The ecological significance of tallgrass prairie arthropods.
- 339 *Annu. Rev. Entomol.* **51**, 387–412 (2006).
- 340 4. Losey, J. E. & Vaughan, M. The Economic Value of Ecological Services Provided by
- 341 Insects. *Bioscience* **56**, 311–323 (2006).
- 342 5. Bradshaw, C. J. A. *et al.* Massive yet grossly underestimated global costs of invasive
- 343 insects. *Nat. Commun.* **7**, 12986 (2016).
- 344 6. Bebber, D. P., Ramotowski, M. A. T. & Gurr, S. J. Crop pests and pathogens move
- 345 polewards in a warming world. *Nat. Clim. Chang.* **3**, 985 (2013).
- 346 7. Sparling, P. F., Hamburg, M. A., Relman, D. A., Choffnes, E. R. & Mack, A. *Vector-Borne*
- 347 *Diseases : Understanding the Environmental, Human Health, and Ecological Connections,*
- 348 *Workshop Summary. Forum on Microbial Threats: Board on Global Health.* (National
- 349 Academies Press, 2008).
- 350 8. Minjauw, B. & McLeod, A. *Tick-borne diseases and poverty : the impact of ticks and tick-*
- 351 *borne diseases on the livelihoods of small-scale and marginal livestock owners in India and*
- 352 *eastern and southern Africa.* (Centre for Tropical Veterinary Medicine, 2003).
- 353 9. Van den Bossche, P., de La Rocque, S., Hendrickx, G. & Bouyer, J. A changing
- 354 environment and the epidemiology of tsetse-transmitted livestock trypanosomiasis. *Trends*
- 355 *Parasitol.* **26**, 236–243 (2010).
- 356 10. WHO | Vector-borne diseases. (2017).
- 357 11. Gubler, D. J. Resurgent vector-borne diseases as a global health problem. *Emerg. Infect.*
- 358 *Dis.* **4**, 442–450 (1998).
- 359 12. Elbers, A. R. W., Koenraadt, C. J. M. & Meiswinkel, R. Mosquitoes and Culicoides biting
- 360 midges: vector range and the influence of climate change. *Rev. Sci. Tech.* **34**, 123–137
- 361 (2015).
- 362 13. Sakai, A. K. *et al.* The Population Biology of Invasive Species. *Annu. Rev. Ecol. Syst.* **32**,
- 363 305–332 (2001).
- 364 14. Rund, S. S. C. & Martinez, M. E. Rescuing Troves of Data to Tackle Emerging Mosquito-
- 365 Borne Diseases. *bioRxiv* 096875 (2018). doi:10.1101/096875
- 366 15. Foley, D. H., Maloney, F. A., Jr, Harrison, F. J., Wilkerson, R. C. & Rueda, L. M. Online
- 367 spatial database of US Army Public Health Command Region-West mosquito surveillance
- 368 records: 1947-2009. *US Army Med. Dep. J.* 29–36 (2011).
- 369 16. Hutchinson, M. L., STROHECKER, Simmons, T. W., Kyle, A. D. & Helwig, M. W.
- 370 Prevalence Rates of *Borrelia burgdorferi* (Spirochaetales: Spirochaetaceae), *Anaplasma*
- 371 *phagocytophilum* (Rickettsiales: Anaplasmataceae), and *Babesia microti* (Piroplasmida:
- 372 *Babesiidae*) in Host-Seeking *Ixodes scapularis* (Acari: Ixodidae) from Pennsylvania.
- 373 *Journal of Medical Entomology* **52**, 693–698 (2015).
- 374 17. Magarey, R. D. *et al.* Risk maps for targeting exotic plant pest detection programs in the
- 375 United States: US risk maps for exotic plant pest detection. *EPPO Bulletin* **41**, 46–56
- 376 (2011).
- 377 18. Wilson, B. E., Beuzelin, J. M., VanWeelden, M. T., Reagan, T. E. & Way, M. O.
- 378 Monitoring Mexican Rice Borer (Lepidoptera: Crambidae) Populations in Sugarcane and
- 379 Rice With Conventional and Electronic Pheromone Traps. *J. Econ. Entomol.* **110**, 150–156
- 380 (2017).
- 381 19. Chandler, M. *et al.* Contribution of citizen science towards international biodiversity

- 382 monitoring. *Biol. Conserv.* **213**, 280–294 (2017).
- 383 20. Kampen, H. *et al.* Approaches to passive mosquito surveillance in the EU. *Parasit. Vectors*  
384 **8**, 9 (2015).
- 385 21. Suprayitno, N., Narakusumo, R. P., von Rintelen, T., Hendrich, L. & Balke, M. Taxonomy  
386 and Biogeography without frontiers - WhatsApp, Facebook and smartphone digital  
387 photography let citizen scientists in more remote localities step out of the dark. *Biodivers*  
388 *Data J* e19938 (2017).
- 389 22. Seltsmann, K. C. *et al.* LepNet: The Lepidoptera of North America Network. *Zootaxa* **4247**,  
390 73–77 (2017).
- 391 23. Short, A. E. Z., Dikow, T. & Moreau, C. S. Entomological Collections in the Age of Big  
392 Data. *Annu. Rev. Entomol.* **63**, 513–530 (2018).
- 393 24. Horton, R. (Comment) Offline: What is medicine’s 5 sigma? *The Lancet* **235**, 1380 (2015).
- 394 25. Nakagawa, S. & Parker, T. H. Replicating research in ecology and evolution: feasibility,  
395 incentives, and the cost-benefit conundrum. *BMC Biol.* **13**, 88 (2015).
- 396 26. Nosek, B. A. *et al.* Promoting an open research culture. *Science* **348**, 1422–1425 (2015).
- 397 27. Parker, T. H. *et al.* Transparency in Ecology and Evolution: Real Problems, Real Solutions.  
398 *Trends Ecol. Evol.* **31**, 711–719 (2016).
- 399 28. Smaldino, P. E. & McElreath, R. The natural selection of bad science. *R Soc Open Sci* **3**,  
400 160384 (2016).
- 401 29. Ihle, M., Winney, I. S., Krystalli, A. & Croucher, M. Striving for transparent and credible  
402 research: practical guidelines for behavioral ecologists. *Behav. Ecol.* **28**, 348–354 (2017).
- 403 30. Poisot, T. E., Mounce, R., Gravel - Ideas in Ecology and, D. & 2013. Moving toward a  
404 sustainable ecological science: don’t let data go to waste! *queens.scholarsportal.info*  
405 (2013).
- 406 31. Roche, D. G. *et al.* Troubleshooting public data archiving: suggestions to increase  
407 participation. *PLoS Biol.* **12**, e1001779 (2014).
- 408 32. Culley, T. M. The frontier of data discoverability: Why we need to share our data.  
409 *Applications in Plant Sciences* **5**, (2017).
- 410 33. Gerstner, K. *et al.* Will your paper be used in a meta-analysis? Make the reach of your  
411 research broader and longer lasting. *Wiley Online Library* (2017).
- 412 34. Ioannidis, J. P. A. *et al.* Repeatability of published microarray gene expression analyses.  
413 *Nat. Genet.* **41**, 149–155 (2009).
- 414 35. Gilbert, K. J. *et al.* Recommendations for utilizing and reporting population genetic  
415 analyses: the reproducibility of genetic clustering using the program STRUCTURE. *Mol.*  
416 *Ecol.* **21**, 4925–4930 (2012).
- 417 36. Roche, D. G., Kruuk, L. E. B., Lanfear, R. & Binning, S. A. Public Data Archiving in  
418 Ecology and Evolution: How Well Are We Doing? *PLoS Biol.* **13**, e1002295 (2015).
- 419 37. Renaut, S., Budden, A. E., Gravel, D., Poisot, T. & Peres-Neto, P. Management, Archiving,  
420 and Sharing for Biologists and the Role of Research Institutions in the Technology-Oriented  
421 Age. *Bioscience* **68**, 400–411 (2018).
- 422 38. Wilkinson, M. D. *et al.* The FAIR Guiding Principles for scientific data management and  
423 stewardship. *Sci Data* **3**, 160018 (2016).
- 424 39. Wilkinson, M. D. *et al.* A design framework and exemplar metrics for FAIRness. *Sci Data*  
425 **5**, 180118 (2018).
- 426 40. Taylor, C. F. *et al.* The minimum information about a proteomics experiment (MIAPE).



- 427 *Nat. Biotechnol.* **25**, 887–893 (2007).
- 428 41. Yilmaz, P. *et al.* Minimum information about a marker gene sequence (MIMARKS) and  
429 minimum information about any (x) sequence (MIxS) specifications. *Nat. Biotechnol.* **29**,  
430 415–420 (2011).
- 431 42. Lourenço, A. *et al.* Minimum information about a biofilm experiment (MIABiE): standards  
432 for reporting experiments and data on sessile microbial communities living at interfaces.  
433 *Pathog. Dis.* **70**, 250–256 (2014).
- 434 43. Brazma, A., Hingamp, P., Quackenbush, J., Sherlock - Nature ..., G. & 2001. Minimum  
435 information about a microarray experiment (MIAME)—toward standards for microarray  
436 data. *nature.com* (2001).
- 437 44. Bustin, S. A. *et al.* The MIQE guidelines: minimum information for publication of  
438 quantitative real-time PCR experiments. *Clin. Chem.* **55**, 611–622 (2009).
- 439 45. York, W. S. *et al.* MIRAGE: the minimum information required for a glycomics  
440 experiment. *Glycobiology* **24**, 402–406 (2014).
- 441 46. Taylor, C. F. *et al.* Promoting coherent minimum reporting guidelines for biological and  
442 biomedical investigations: the MIBBI project. *Nat. Biotechnol.* **26**, 889–896 (2008).
- 443 47. Wieczorek, J. *et al.* Darwin Core: an evolving community-developed biodiversity data  
444 standard. *PLoS One* **7**, e29715 (2012).
- 445 48. Giraldo-Calderón, G. I. *et al.* VectorBase: an updated bioinformatics resource for  
446 invertebrate vectors and other organisms related with human diseases. *Nucleic Acids Res.*  
447 **43**, D707–13 (2015).
- 448 49. Benson, D. A. *et al.* GenBank. *Nucleic Acids Res.* **41**, D36–42 (2013).
- 449 50. Berman, H., Henrick, K. & Nakamura, H. Announcing the worldwide Protein Data Bank.  
450 *Nat. Struct. Biol.* **10**, 980 (2003).
- 451 51. GBIF. Available at: <http://gbif.org>. (Accessed: 26th March 2018)
- 452 52. Heinrich, P. L., Gilbert, E., Cobb, N. S. & Franz, N. Symbiota collections of arthropods  
453 network (SCAN): A data portal built to visualize, manipulate, and export species  
454 occurrences.
- 455 53. Perryman, S. A. M. *et al.* The electronic Rothamsted Archive (e-RA), an online resource for  
456 data from the Rothamsted long-term experiments. *Sci Data* **5**, 180072 (2018).
- 457 54. Gossner, M. M. *et al.* A summary of eight traits of Coleoptera, Hemiptera, Orthoptera and  
458 Araneae, occurring in grasslands in Germany. *Sci Data* **2**, 150013 (2015).
- 459 55. Hedefalk, F., Svensson, P. & Harrie, L. Spatiotemporal historical datasets at micro-level for  
460 geocoded individuals in five Swedish parishes, 1813-1914. *Sci Data* **4**, 170046 (2017).
- 461 56. The American Society for Cell Biology. San Francisco Declaration on Research  
462 Assessment. (2012). Available at: [http://www.ascb.org/wp-](http://www.ascb.org/wp-content/uploads/2017/07/sfdora.pdf)  
463 [content/uploads/2017/07/sfdora.pdf](http://www.ascb.org/wp-content/uploads/2017/07/sfdora.pdf).
- 464 57. Chavan, V. & Penev, L. The data paper: a mechanism to incentivize data publishing in  
465 biodiversity science. *BMC Bioinformatics* **12 Suppl 15**, S2 (2011).
- 466 58. Abell, K. J., Bauer, L. S., Duan, J. J. & Van Driesche, R. Long-term monitoring of the  
467 introduced emerald ash borer (Coleoptera: Buprestidae) egg parasitoid, *Oobius agrili*  
468 (Hymenoptera: Encyrtidae), in Michigan, USA and evaluation of a newly developed  
469 monitoring technique. *Biol. Control* **79**, 36–42 (2014).

470