

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, VectorBase, Department of Biological Science, University of
7 Notre Dame, IN, USA.

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

11
12 Lauren Cator, l.cator@imperial.ac.uk, Department of Life Sciences, Imperial College London,
13 UK

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

17
18 Scott J. Emrich, 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, 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, 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, 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, Global Biodiversity Information Facility (GBIF) Secretariat,
34 Copenhagen, Denmark

35
36 Sarah A. Kelly, 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, VectorBase and Vector Immunogenomics and
40 Infection Laboratory, Department of Life Sciences, Imperial College London, UK

42 Cynthia Lord, 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, 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, Institute of Biology, University
49 of Neuchâtel, 2000, Neuchâtel, Switzerland

50
51 Sadie J. Ryan, 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, Global Biodiversity Information Facility (GBIF) Secretariat,
57 Copenhagen, Denmark

58
59 Kurt Vandegrift, kjv1@psu.edu, Center for Infectious Disease Dynamics, Department of
60 Biology, The Pennsylvania State University, PA, USA

61
62 Matthew Watts, m.watts@imperial.ac.uk, Department of Life Sciences, Imperial College
63 London, UK

64
65 Jennifer M. Zaspel, 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, 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

74 Introduction

75 Arthropods play a dominant role in the dynamics of practically all natural and human-modified
76 terrestrial ecosystems¹⁻³, and have significant economic and health effects. For example, certain

77 insects provide significant economic benefits (*e.g.* pollination) exceeding \$57 billion a year to
78 the United States alone⁴. Meanwhile, invasive insects cost an estimated \$70 billion dollars per
79 year globally⁵ and insect pests may reduce agricultural harvests by up to 16%, with an equal
80 amount of further losses of harvested goods⁶. Particularly noteworthy is a subset of arthropods
81 that are disease vectors, transmitting pathogens to and between animals as well as plants. Vector-
82 borne diseases cause billions of dollars in crop and livestock losses, every year⁷⁻⁹. In humans,
83 vector borne diseases account for more than 17% of all infectious diseases (*e.g.* malaria, Chagas,
84 dengue, and leishmaniasis, Zika, West Nile, Lyme disease, and sleeping sickness), with hundreds
85 of thousands of deaths, hundreds of millions of cases, and billions of people at risk, annually^{10,11}.

86

87 The current economic and health burden of arthropod pests, exacerbated by invasive species, and
88 uncertain effects of climate change^{12,13}, has driven significant research programs and data
89 collection efforts. These include crop pest, mosquito, and tick survey and reporting initiatives¹⁴⁻
90 ¹⁸, citizen science projects¹⁹⁻²¹, and digitization of museum specimen data^{22,23}, all yielding a rich
91 and growing trove of field-based data spanning multiple spatial and temporal scales. Monitoring
92 arthropod abundance (*e.g.* Figure 1) in different disciplines (*e.g.*, biodiversity research, pest-
93 control assessment, vector-borne disease monitoring, or pollination research) uses similar
94 techniques, with similar objectives: to quantify abundance, phenology and geographical ranges
95 of target arthropod species. Despite a growing number of data collections, they are often not
96 reusable, or comparable to similar data, due to a lack of standardization and metadata. In
97 contrast, the advent of the deposition of data from high-throughput technologies (*e.g.* NCBI and
98 GenBank), data and code sharing, and other practices to improve transparency and reusability of
99 research results are increasing rapidly across the sciences²⁴⁻²⁹. Furthering these advances through

100 standardization and public archiving of arthropod abundance data can bring significant benefits,
101 including (1) supporting empirical parameterization and validation of mathematical models (*e.g.*
102 of pest or disease emergence and spread), (2) validation of model predictions, (3) reduction in
103 the duplication of expensive empirical research, and (4) revealing new patterns and questions
104 through meta-analyses^{30–33}. This will also lead to substantial public benefit through improved
105 human, animal, plant, and ecosystem health, and reduced economic costs.

106

107 A key impediment to the re-use of these data is the lack of adequate metadata or data descriptors
108 (*i.e.* data about the data)^{34–37}. In general, for data to be most valuable to the scientific community,
109 they should meet the FAIR Principles – they should be Findable, Accessible, Interoperable and
110 Reusable – and delineate the key components of good data management and stewardship
111 practices^{38,39}. Data are Findable and Accessible when they are archived and freely downloadable
112 from an online public data repository that is indexed and easily searchable. Interoperability and
113 reusability describe the ease with which humans or computer programs can understand the data
114 (*e.g.* via metadata) and explore/re-use them across a variety of non-proprietary platforms. Even
115 when data are available, metadata for arthropod abundance data are often absent or not readily
116 interpretable, limiting their reusability at a fundamental level.

117

118 **A minimum information standard for arthropod abundance data**

119 Here, we present a Minimum Information for Reusable Arthropod Abundance Data (MIReAAD)
120 standard for reporting primarily longitudinal (repeated, temporally explicit) field-based
121 collections of arthropods. In the same manner as has been developed in other biological
122 disciplines^{40–45}, this standard is “minimum” because it defines the necessary minimal information

123 required to understand and reuse a dataset without consulting any further text, materials, or
124 methods⁴⁶. MIREAAD is designed to facilitate data archiving efforts of publishers and field
125 researchers. It is not a data model and therefore does not define controlled vocabularies, or
126 specific field titles, but should be easy to understand, and interpret by the wider scientific
127 community⁴⁶.

128

129 The minimal standards are separated into two components, metadata and data. For each
130 component, we provide a description of the information that should be included,
131 recommendations for how to make that information as useful as possible, and examples. The
132 metadata component (Table 1) includes information for the origin of the data set (*e.g.* study
133 information and licensing for usage). The second component (Table 2) lists and describes
134 specific data fields that should be included in data collection sheets. We also provide
135 recommendations and examples to demonstrate how these recommendations can be
136 implemented. MIREAAD was designed to match the data that are generally collected by
137 academic researchers and surveillance initiatives, and can serve as a checklist for important
138 information that needs to be recorded but is often unintentionally omitted (*e.g.* Figure 2A). By
139 adhering to MIREAAD standards, omissions and ambiguity can be avoided even if the data are
140 shared in different formats (Figure 2B and C). Finally, we identify common problems likely to
141 be encountered across all the MIREAAD metadata and data fields, and data quality standards that
142 can be employed to avoid confusion (Box 1).

143

Box 1. Data quality standards

Language. Once data are ready to be deposited/submitted, all fields and data need to be written in English. This will allow researchers and data curators worldwide to understand and reuse the data.

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.

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).

144

145 **Examples**

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

153

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

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

157

158 2. *Wide format trapping data*. Each row captures count data from a given sampling event. Each
159 identified taxonomic group is identified in a separate column. An ‘additional sample
160 information’ field, ‘sub-location,’ has been added to describe the various locations around the
161 village where collections were made. [[Sup Datasheet 2](#)]. This illustrates an example of adult
162 mosquito populations that have been tracked over time and in specific locations. Also see Figure
163 2C.

164

165 3. *Complex trapping data scenario*. Tick surveillance performed using tick drags and flags and
166 collections of ectoparasites on trapped mice. The tick drags/flags report three life stages
167 independently (adult, larvae, and nymph) [[Sup Datasheet 3](#)]. Larvae are only identified to the
168 genus, while adults and nymphs are identified to the species. A Sample Name is used to help link
169 these records (but would not be necessary.) The mouse survey uses an additional sample
170 information field to record the sex of the trapped mouse from which the parasites were collected
171 [[Sup Datasheet 4](#)].

172 Discussion

173 **MIReAAD as the path to FAIR data principles**

174 We designed MIReAAD to achieve a balance between standards that are too onerous for data
175 generators and standards that are sufficient to ensure at least minimal reusability^{31,40}. Like all
176 minimum standards, MIReAAD only aims at ensuring data ‘Reusability’. However, ultimately

177 this will promote the implementation of data models — the explicit definition of data field
178 names, data formats (*e.g.*, for dates and GPS locations), and controlled vocabularies (*e.g.*, the
179 Darwin Core⁴⁷). Data models enable ‘Interoperability’, and in turn facilitate structured databases,
180 public repositories, and development of data analysis tools^{46,48}. Deposition in open databases
181 make data ‘Findable’ and ‘Accessible’^{49–51}. MIREAAD compliant data contain sufficient
182 information for established aggregators/databases such as VectorBase and SCAN (Symbiota
183 Collections of Arthropods Network⁵²) to process and store the data in a standardized data model
184 (*e.g.*, Darwin Core⁴⁷), and ultimately facilitate data transfer to even more comprehensive
185 biodiversity databases (*e.g.* GBIF, which contains over one billion species occurrence records⁵¹).
186 In this way, MIREAAD opens the door to FAIR data and more sophisticated methods to integrate
187 data across many scales.

188

189 **Benefits to field researchers**

190 It is essential that the benefits of a minimal data standard extend not just to data re-users, but also
191 to the researchers who collect and generate data in the first place. MIREAAD provides a
192 framework for data preparation that can help scientists achieve recognized professional merit for
193 sharing data such as increased citation rates, academic recognition, opportunities for co-
194 authorship, and new collaborations [sensu Roche et al. 2014³¹]. Large, deposited data sets can
195 now themselves be standalone, citable “data papers” (*e.g.* ^{53–55}) or even depositions without any
196 traditional manuscript (but with persistent identifiers, such as a DOI number), if desired. Data
197 sets are increasingly recognized as valuable research outputs that count towards academic
198 recognition and professional advancement (*e.g.* grants, interviews, and tenure). For example,
199 several funders (*e.g.* United States National Science Foundation and Swiss National Science

200 Foundation) have adopted or are in the process of adopting the Declaration on Research
201 Assessments (DORA)⁵⁶, offering further opportunities for data generators to gain recognition and
202 publication credit for their work⁵⁷. Also, an increasing number of funders are mandating public
203 data access, and detailed data management plans are often required even at the grant proposal
204 stage. Therefore, reporting data according to MIREAAD will provide a foundational pipeline for
205 stipulating archival formats.

206

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

214

215 **Further MIREAAD applications and extensions**

216 The creation of minimum information standards for these types of databases facilitates analyses
217 of data at the scales that cannot be attained by a single individual or lab group. Linking records
218 to additional information also extends the utility of these data to address population level
219 questions. For example, a well-populated database presents opportunities to investigate
220 interactions between populations of different species of arthropod that overlap in geography, but
221 may be of interest individually to different realms of research. As a case in point, in the
222 northeastern USA, *Agilus plannipennis*, the Emerald Ash Borer (EAB), is a highly destructive

223 invasive insect, monitored closely by both state and federal agencies for management⁵⁸.

224 Interestingly, EAB are creating lots of new habitat for carpenter bees, a species interaction that
225 can be tracked and anticipated using large scale arthropod data.

226

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

235

236 We note that MIRreAAD is applicable not only to abundance measurements, but could be easily
237 extended to any other kind of routinely sampled time-series field data. For example, in addition
238 to aphid abundance, plant pathogen (such as mosaic virus) infection and insecticide resistance
239 statuses of the aphids could be reported in MIRreAAD format.

240 Conclusion

241 We present MIRreAAD, a minimum information standard for representing arthropod
242 abundance data. MIRreAAD will facilitate collation and analyses of data at scales that cannot be
243 attained by a single individual or lab, to address key questions across temporal and spatial scales,

244 such as within and across-year phenology of abundance of target arthropod taxa over large
 245 geographical areas. This is particularly important given the pressing need to understand and
 246 predict the population dynamics of harmful (e.g., disease vectors and pests) as well as beneficial
 247 (e.g., pollinators, bio-control agents) arthropods in natural and human modified landscapes. This
 248 is the first step for achieving the broad benefits of FAIR data for arthropod abundance. We call
 249 on data generators, authors, reviewers, editors, journals, research infrastructures (e.g. data
 250 repositories) and funders to embrace MIREAAD as a standard to facilitate FAIR data use and
 251 compliance for arthropod abundance data.

252
 253 **Table 1. The MIREAAD Study Information (Resource metadata) fields.** The information in this
 254 table should be included with every data submission, for example by including data in the file
 255 header as demonstrated in Data Files 1-4.
 256

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 orcid.org/0000-0002-5690-3300 kurtvandegrift@gmail.com State University Agricultural Extension John Smith (jsmith@StateU.edu) www.StateU.edu/AgriculturalExtension/
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)).		<p>“A web-based relational database for monitoring and analyzing mosquito population dynamics Sucaet Y, Van Hemert J, Tucker B, Bartholomay L.”</p> <p>“PMID: 18714883”</p> <p>Horiuchi, Kaho, Kosei Hashimoto, and Fumio Hayashi. "Cantharidin world in air: Spatiotemporal distributions of flying canthariphilous insects in the forest interior." Entomological Science (2018).</p>
Species Identification Method	A description of method of species identification. Particularly important for cryptic species complexes.		<p>“Morphological”</p> <p>“Genotyped, using method of Smith et al 2014, PMID: 18714883”</p>
Not present vs zero information	Indication of what gaps, zeros, NA, etc mean.	<p>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 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>	<p>“Zero indicates was looked for and not found. NA represents a trap failure etc”</p>
GPS obfuscation information	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.	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>“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 https://creativecommons.org 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”.</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”, 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”</p>
------------------------	---	--	---------------------------------------

257

258

259 **Table 2.** The MIR_eAAD data fields. Fig 1B provides an annotated example.

260

Field(s)	Details	Recommendations	Examples
Start Time (for collection)	<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>Month should be written out so that day and month are not confused.</p> <p>Provide timezone in field or in header if hourly data, a 24 hour clock is preferred, but should be made unambiguous as to which time format is being used.</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>
End Time (for collection)	<p>End time of the data sample</p>	<p>See above.</p>	<p>See above.</p>

	<p>collection.</p> <p>e.g. The trap was collected...</p>	<p>If instantaneous data collection (e.g. a tick drag), End Time may be the same as Start Time.</p>	
Location	<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>

Collection method	Sampling apparatus (e.g. trap type, observation method)		<p>"CDC light trap"</p> <p>"Tick drag"</p> <p>"Quadrat count"</p> <p>"BG Sentinel Trap"</p> <p>"Pitfall trap"</p> <p>"Larval dip"</p> <p>"Johnson suction trap"</p> <p>"Lindgren Funnel Trap"</p>
Collection attractants	The attractant/ lures used to attract insects to a trap or collection		<p>"None"</p> <p>"Carbon dioxide"</p> <p>"UV light"</p> <p>"BG-Sweetscent Mosquito Lure"</p> <p>"Human/animal bait"</p>
Collection area	The spatial extent (area or volume) of the sample.	<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 m²"</p> <p>"1 liter"</p> <p>"1 ha"</p> <p>"10m³"</p>
Taxonomy	Classification of sample collected.	<p>Scientific genus and species preferred.</p> <p>Avoid abbreviation.</p>	<p><i>"Ixodes scapularis"</i></p> <p><i>"Aedes aegypti"</i></p> <p><i>'Anopheles gambiae sensu stricto'</i></p>

<p>Unit(s) of measurement and observation</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 m²"</p> <p>"Female" and "Adult"</p> <p>"Male and Female" and "Nymphs"</p>
<p>Value</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,"
-------------	---	--	--

261 Field names in bold should be considered also required. Remaining fields are optional or
262 depend on the complexity of the experimental design
263

264 Acknowledgements

265
266 The seeds of this effort were planted in 2016 at a meeting of VectorBiTE, which is a cross-
267 disciplinary research coordination network (RCN) for disease vectors. Samuel S.C. Rund,
268 Matthew Watts, Kurt Vandegrift, Naveed Heydari, Cynthia Lord, Michael Johansson, Samraat
269 Pawar, and Sadie J. Ryan, received travel funding from NIH grant 1R01AI122284-01 and
270 BBSRC grant BB/N013573/1 as part of the joint [NIH-NSF-USDA-BBSRC] Ecology and
271 Evolution of Infectious Diseases program.

272
273 Samuel S.C. Rund was funded by the Royal Society (NF140517). Rund, Daniel Lawson, Robert
274 M. MacCallum, Sarah A. Kelly, Gloria I. Giraldo-Calderon and Scott J. Emrich were supported
275 by the National Institute of Allergy and Infectious Diseases, National Institutes of Health,
276 Department of Health and Human Services, under Contract No. HHSN272201400029C
277 (VectorBase Bioinformatics Resource Center).

278

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

281

282 Naveed Heydari and Sadie J. Ryan were funded by National Science Foundation (NSF DEB
283 EEID 1518681).

284

285 Sadie J. Ryan was additionally funded by NIH 1R01AI136035-01, and CDC grant
286 1U01CK000510-01: Southeastern Regional Center of Excellence in Vector-Borne

287 Diseases: the Gateway Program. This publication was supported by the Cooperative Agreement

288 Number above from the Centers for Disease Control and Prevention. Its contents are solely the

289 responsibility of the authors and do not necessarily represent the official views of the Centers for

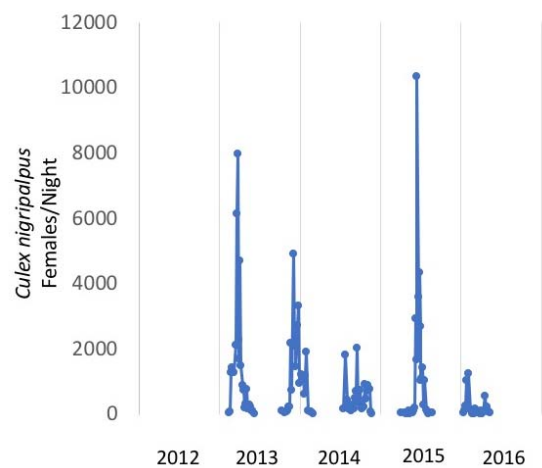
290 Disease Control and Prevention.

291

292 Jennifer M. Zaspel was funded by the National Science Foundation Division of Biological

293 Infrastructure (NSF 1561448, NSF 1601957).

294 **Figures**



295 **Figure 1. Example population abundance time-series.**
296
297

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

298

299

300

301

302

303

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

304

References

305

306

307

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

21

- 308 Annual Review of Entomology.
- 309 3. Whiles, M. R. & Charlton, R. E. The ecological significance of tallgrass prairie arthropods.
310 *Annu. Rev. Entomol.* **51**, 387–412 (2006).
- 311 4. Losey, J. E. & Vaughan, M. The Economic Value of Ecological Services Provided by
312 Insects. *Bioscience* **56**, 311–323 (2006).
- 313 5. Bradshaw, C. J. A. *et al.* Massive yet grossly underestimated global costs of invasive
314 insects. *Nat. Commun.* **7**, 12986 (2016).
- 315 6. Bebber, D. P., Ramotowski, M. A. T. & Gurr, S. J. Crop pests and pathogens move
316 polewards in a warming world. *Nat. Clim. Chang.* **3**, 985 (2013).
- 317 7. Sparling, P. F., Hamburg, M. A., Relman, D. A., Choffnes, E. R. & Mack, A. *Vector-Borne*
318 *Diseases : Understanding the Environmental, Human Health, and Ecological Connections,*
319 *Workshop Summary. Forum on Microbial Threats: Board on Global Health.* (National
320 Academies Press, 2008).
- 321 8. Minjauw, B. & McLeod, A. *Tick-borne diseases and poverty : the impact of ticks and tick-*
322 *borne diseases on the livelihoods of small-scale and marginal livestock owners in India and*
323 *eastern and southern Africa.* (Centre for Tropical Veterinary Medicine, 2003).
- 324 9. Van den Bossche, P., de La Rocque, S., Hendrickx, G. & Bouyer, J. A changing
325 environment and the epidemiology of tsetse-transmitted livestock trypanosomiasis. *Trends*
326 *Parasitol.* **26**, 236–243 (2010).
- 327 10. WHO | Vector-borne diseases. (2017).
- 328 11. Gubler, D. J. Resurgent vector-borne diseases as a global health problem. *Emerg. Infect.*
329 *Dis.* **4**, 442–450 (1998).
- 330 12. Elbers, A. R. W., Koenraadt, C. J. M. & Meiswinkel, R. Mosquitoes and Culicoides biting
331 midges: vector range and the influence of climate change. *Rev. Sci. Tech.* **34**, 123–137
332 (2015).
- 333 13. Sakai, A. K. *et al.* The Population Biology of Invasive Species. *Annu. Rev. Ecol. Syst.* **32**,
334 305–332 (2001).
- 335 14. Rund, S. S. C. & Martinez, M. E. Rescuing Troves of Data to Tackle Emerging Mosquito-
336 Borne Diseases. *bioRxiv* 096875 (2018). doi:10.1101/096875
- 337 15. Foley, D. H., Maloney, F. A., Jr, Harrison, F. J., Wilkerson, R. C. & Rueda, L. M. Online
338 spatial database of US Army Public Health Command Region-West mosquito surveillance
339 records: 1947-2009. *US Army Med. Dep. J.* 29–36 (2011).
- 340 16. Hutchinson, M. L., STROHECKER, Simmons, T. W., Kyle, A. D. & Helwig, M. W.
341 Prevalence Rates of *Borrelia burgdorferi* (Spirochaetales: Spirochaetaceae), *Anaplasma*
342 *phagocytophilum* (Rickettsiales: Anaplasmataceae), and *Babesia microti* (Piroplasmida:
343 Babesiidae) in Host-Seeking *Ixodes scapularis* (Acari: Ixodidae) from Pennsylvania.
344 *Journal of Medical Entomology* **52**, 693–698 (2015).
- 345 17. Magarey, R. D. *et al.* Risk maps for targeting exotic plant pest detection programs in the
346 United States: US risk maps for exotic plant pest detection. *EPPO Bulletin* **41**, 46–56
347 (2011).
- 348 18. Wilson, B. E., Beuzelin, J. M., VanWeelden, M. T., Reagan, T. E. & Way, M. O.
349 Monitoring Mexican Rice Borer (Lepidoptera: Crambidae) Populations in Sugarcane and
350 Rice With Conventional and Electronic Pheromone Traps. *J. Econ. Entomol.* **110**, 150–156
351 (2017).
- 352 19. Chandler, M. *et al.* Contribution of citizen science towards international biodiversity

- 353 monitoring. *Biol. Conserv.* **213**, 280–294 (2017).
- 354 20. Kampen, H. *et al.* Approaches to passive mosquito surveillance in the EU. *Parasit. Vectors*
355 **8**, 9 (2015).
- 356 21. Suprayitno, N., Narakusumo, R. P., von Rintelen, T., Hendrich, L. & Balke, M. Taxonomy
357 and Biogeography without frontiers - WhatsApp, Facebook and smartphone digital
358 photography let citizen scientists in more remote localities step out of the dark. *Biodivers*
359 *Data J* e19938 (2017).
- 360 22. Seltsmann, K. C. *et al.* LepNet: The Lepidoptera of North America Network. *Zootaxa* **4247**,
361 73–77 (2017).
- 362 23. Short, A. E. Z., Dikow, T. & Moreau, C. S. Entomological Collections in the Age of Big
363 Data. *Annu. Rev. Entomol.* **63**, 513–530 (2018).
- 364 24. Horton, R. (Comment) Offline: What is medicine’s 5 sigma? *The Lancet* **235**, 1380 (2015).
- 365 25. Nakagawa, S. & Parker, T. H. Replicating research in ecology and evolution: feasibility,
366 incentives, and the cost-benefit conundrum. *BMC Biol.* **13**, 88 (2015).
- 367 26. Nosek, B. A. *et al.* Promoting an open research culture. *Science* **348**, 1422–1425 (2015).
- 368 27. Parker, T. H. *et al.* Transparency in Ecology and Evolution: Real Problems, Real Solutions.
369 *Trends Ecol. Evol.* **31**, 711–719 (2016).
- 370 28. Smaldino, P. E. & McElreath, R. The natural selection of bad science. *R Soc Open Sci* **3**,
371 160384 (2016).
- 372 29. Ihle, M., Winney, I. S., Krystalli, A. & Croucher, M. Striving for transparent and credible
373 research: practical guidelines for behavioral ecologists. *Behav. Ecol.* **28**, 348–354 (2017).
- 374 30. Poisot, T. E., Mounce, R., Gravel - Ideas in Ecology and, D. & 2013. Moving toward a
375 sustainable ecological science: don’t let data go to waste! *queens.scholarsportal.info*
376 (2013).
- 377 31. Roche, D. G. *et al.* Troubleshooting public data archiving: suggestions to increase
378 participation. *PLoS Biol.* **12**, e1001779 (2014).
- 379 32. Culley, T. M. The frontier of data discoverability: Why we need to share our data.
380 *Applications in Plant Sciences* **5**, (2017).
- 381 33. Gerstner, K. *et al.* Will your paper be used in a meta-analysis? Make the reach of your
382 research broader and longer lasting. *Wiley Online Library* (2017).
- 383 34. Ioannidis, J. P. A. *et al.* Repeatability of published microarray gene expression analyses.
384 *Nat. Genet.* **41**, 149–155 (2009).
- 385 35. Gilbert, K. J. *et al.* Recommendations for utilizing and reporting population genetic
386 analyses: the reproducibility of genetic clustering using the program STRUCTURE. *Mol.*
387 *Ecol.* **21**, 4925–4930 (2012).
- 388 36. Roche, D. G., Kruuk, L. E. B., Lanfear, R. & Binning, S. A. Public Data Archiving in
389 Ecology and Evolution: How Well Are We Doing? *PLoS Biol.* **13**, e1002295 (2015).
- 390 37. Renaut, S., Budden, A. E., Gravel, D., Poisot, T. & Peres-Neto, P. Management, Archiving,
391 and Sharing for Biologists and the Role of Research Institutions in the Technology-Oriented
392 Age. *Bioscience* **68**, 400–411 (2018).
- 393 38. Wilkinson, M. D. *et al.* The FAIR Guiding Principles for scientific data management and
394 stewardship. *Sci Data* **3**, 160018 (2016).
- 395 39. Wilkinson, M. D. *et al.* A design framework and exemplar metrics for FAIRness. *Sci Data*
396 **5**, 180118 (2018).
- 397 40. Taylor, C. F. *et al.* The minimum information about a proteomics experiment (MIAPE).

- 398 *Nat. Biotechnol.* **25**, 887–893 (2007).
- 399 41. Yilmaz, P. *et al.* Minimum information about a marker gene sequence (MIMARKS) and
400 minimum information about any (x) sequence (MIxS) specifications. *Nat. Biotechnol.* **29**,
401 415–420 (2011).
- 402 42. Lourenço, A. *et al.* Minimum information about a biofilm experiment (MIABiE): standards
403 for reporting experiments and data on sessile microbial communities living at interfaces.
404 *Pathog. Dis.* **70**, 250–256 (2014).
- 405 43. Brazma, A., Hingamp, P., Quackenbush, J., Sherlock - Nature ..., G. & 2001. Minimum
406 information about a microarray experiment (MIAME)—toward standards for microarray
407 data. *nature.com* (2001).
- 408 44. Bustin, S. A. *et al.* The MIQE guidelines: minimum information for publication of
409 quantitative real-time PCR experiments. *Clin. Chem.* **55**, 611–622 (2009).
- 410 45. York, W. S. *et al.* MIRAGE: the minimum information required for a glycomics
411 experiment. *Glycobiology* **24**, 402–406 (2014).
- 412 46. Taylor, C. F. *et al.* Promoting coherent minimum reporting guidelines for biological and
413 biomedical investigations: the MIBBI project. *Nat. Biotechnol.* **26**, 889–896 (2008).
- 414 47. Wieczorek, J. *et al.* Darwin Core: an evolving community-developed biodiversity data
415 standard. *PLoS One* **7**, e29715 (2012).
- 416 48. Giraldo-Calderón, G. I. *et al.* VectorBase: an updated bioinformatics resource for
417 invertebrate vectors and other organisms related with human diseases. *Nucleic Acids Res.*
418 **43**, D707–13 (2015).
- 419 49. Benson, D. A. *et al.* GenBank. *Nucleic Acids Res.* **41**, D36–42 (2013).
- 420 50. Berman, H., Henrick, K. & Nakamura, H. Announcing the worldwide Protein Data Bank.
421 *Nat. Struct. Biol.* **10**, 980 (2003).
- 422 51. GBIF. Available at: <http://gbif.org>. (Accessed: 26th March 2018)
- 423 52. Heinrich, P. L., Gilbert, E., Cobb, N. S. & Franz, N. Symbiota collections of arthropods
424 network (SCAN): A data portal built to visualize, manipulate, and export species
425 occurrences.
- 426 53. Perryman, S. A. M. *et al.* The electronic Rothamsted Archive (e-RA), an online resource for
427 data from the Rothamsted long-term experiments. *Sci Data* **5**, 180072 (2018).
- 428 54. Gossner, M. M. *et al.* A summary of eight traits of Coleoptera, Hemiptera, Orthoptera and
429 Araneae, occurring in grasslands in Germany. *Sci Data* **2**, 150013 (2015).
- 430 55. Hedefalk, F., Svensson, P. & Harrie, L. Spatiotemporal historical datasets at micro-level for
431 geocoded individuals in five Swedish parishes, 1813-1914. *Sci Data* **4**, 170046 (2017).
- 432 56. The American Society for Cell Biology. San Francisco Declaration on Research
433 Assessment. (2012). Available at: [http://www.ascb.org/wp-](http://www.ascb.org/wp-content/uploads/2017/07/sfdora.pdf)
434 [content/uploads/2017/07/sfdora.pdf](http://www.ascb.org/wp-content/uploads/2017/07/sfdora.pdf).
- 435 57. Chavan, V. & Penev, L. The data paper: a mechanism to incentivize data publishing in
436 biodiversity science. *BMC Bioinformatics* **12 Suppl 15**, S2 (2011).
- 437 58. Abell, K. J., Bauer, L. S., Duan, J. J. & Van Driesche, R. Long-term monitoring of the
438 introduced emerald ash borer (Coleoptera: Buprestidae) egg parasitoid, *Oobius agrili*
439 (Hymenoptera: Encyrtidae), in Michigan, USA and evaluation of a newly developed
440 monitoring technique. *Biol. Control* **79**, 36–42 (2014).

441