

1 A database of egg size and shape from more than 6,700 insect 2 species

3 Samuel H. Church^{*,1,†}, Seth Donoughe^{*,1,2}, Bruno A. S. de Medeiros¹, Cassandra G. Extavour^{1,3}

4 Contents

5	1 Abstract	2
6	2 Background & summary	2
7	3 Methods	3
8	3.1 Gathering primary literature with egg descriptions	3
9	3.2 Defining egg traits	5
10	3.3 Extracting egg descriptions from text sources	6
11	3.4 Measuring published images of eggs	6
12	3.5 Assessing the accuracy of image measuring software	7
13	3.6 Calculating final and transformed values	8
14	3.7 Cross-referencing entries with taxonomic and genetic databases	8
15	3.8 Assessing intraspecific variation	9
16	3.9 Assessing the precision of entries	9
17	3.10 Assessing the phylogenetic sampling	10
18	4 Code availability	10
19	5 Data records	10
20	6 Technical validation	10
21	7 Acknowledgements	13
22	8 Author contributions	13
23	9 Competing interests	13
24	References	17

* Samuel H. Church and Seth Donoughe contributed equally to this work.

1 Department of Organismic and Evolutionary Biology, Harvard University, Cambridge, MA 02138, United States

2 *Current address:* Department of Cell and Molecular Biology, University of Chicago, Chicago, IL 60637, United States

3 Department of Molecular and Cellular Biology, Harvard University, Cambridge, MA 02138, United States

† Correspondence to church@g.harvard.edu

25 **1 Abstract**

26 Offspring size is a fundamental trait in disparate biological fields of study. This trait can be measured as the size
27 of plant seeds, animal eggs, or live young, and it influences ecological interactions, organism fitness, maternal
28 investment, and embryonic development. Although multiple evolutionary processes have been predicted to drive the
29 evolution of offspring size, the phylogenetic distribution of this trait remains poorly understood, due to the difficulty
30 of reliably collecting and comparing offspring size data from many species. Here we present a database of 10,449
31 morphological descriptions of insect eggs, with records for 6,706 unique insect species and representatives from
32 every extant hexapod order. The dataset includes eggs whose volumes span more than eight orders of magnitude. We
33 created this database by partially automating the extraction of egg traits from the primary literature. In the process,
34 we overcame challenges associated with large-scale phenotyping by designing and employing custom bioinformatic
35 solutions to common problems. We matched the taxa in this database to the currently accepted scientific names in
36 taxonomic and genetic databases, which will facilitate the use of this data for testing pressing evolutionary hypotheses
37 in offspring size evolution.

38 **2 Background & summary**

39 The size of a reproductive propagule, for example an animal egg or a plant seed, has crucial implications for the
40 biology of both the parent and the offspring¹⁻³. From the perspective of the parent organism, propagule size is a
41 component of the maternal investment in each offspring², and propagule size is predicted to be positively correlated
42 with adult body size and negatively correlated with propagule number³⁻⁵. From the perspective of the offspring,
43 the size of the propagule is relevant to the starting material for embryonic development, and it can impact both life
44 history and ecological interactions^{2,6}. Evolutionary hypotheses have been proposed to explain patterns in the diversity
45 of propagule size, yet the robustness or generality of the patterns themselves have rarely been tested across species³.
46 To understand the evolutionary forces driving propagule size evolution, we need large-scale, reliable descriptions of
47 the distribution of propagule size across the evolutionary tree.

48 Insect eggs come in an incredible diversity of shapes and sizes^{7,8}. The thousands of egg descriptions in the ento-
49 mological literature, however, have never to our knowledge been systematically compiled across insects. Without a
50 comparison of egg sizes across insects, we cannot ascertain basic information such as the extant range of insect egg
51 sizes, or the relationship between size and ecology or development. To address this problem, we created a database of
52 quantitative parameters describing egg morphology from the entomological literature. All data were collected from
53 published records, including both measurements reported in text descriptions of insect eggs, as well as our own new
54 measurements of published images. We developed custom software that allowed us to collect data from thousands
55 of publications efficiently and reproducibly (Figure 1). We provide this software as a set of tools that can assist other
56 scientists in collecting phenotypic data from the literature (see Methods).

57 Using this software we extracted egg descriptions from 1,756 publications from the past 250 years (Table 1). The
58 database has 10,449 entries representing every extant order of insects, and 6,706 unique insect species. The insect
59 egg database includes descriptions of egg size and shape (Table 2), and the scientific name of each entry has been

60 matched to current taxonomic and genetic databases. The egg database is made publicly available for download (see
61 Methods).

62 Insect egg sizes vary between species, within species, and within a single individual⁷, and the database described here
63 contains variation from all of these sources. We calculated the degree of intraspecific variation in egg length for all
64 taxa where these data were available in the literature. We additionally assessed the variation in the precision used to
65 record data for all database entries. This provides the necessary information to account for sources of variation in a
66 comparative study of insect egg morphology.

67 The insect egg database includes representatives of all insect orders (Table 1), but these orders are not equivalent to
68 each other either in terms of number of extant species or in the historical degree of entomological study^{9,10}. We
69 therefore assessed the phylogenetic coverage of the insect egg database relative to the number of species estimated for
70 each clade. This enables evaluation of the potential bias present in the database, and highlights undersampled clades
71 as potential priorities for future study.

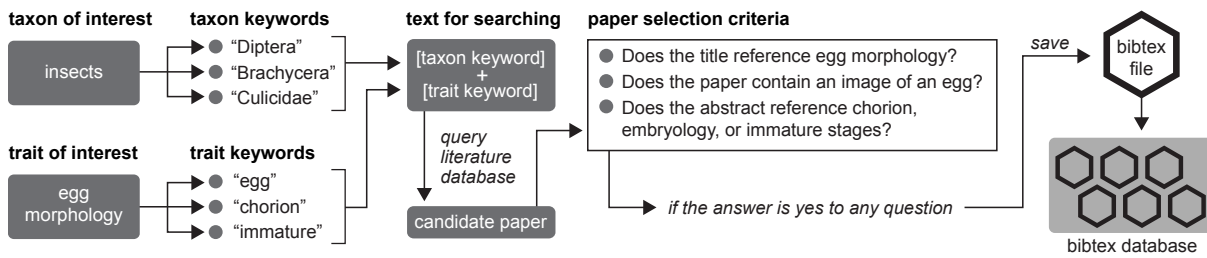
72 The methods used to create the insect egg database include solutions to challenges in assembling phenotypic data
73 from large groups of organisms. Phenotypic descriptions can require great resources and expertise to reliably collect,
74 identify, and describe morphological features across thousands of species¹¹. This expense can limit macroevolutionary
75 studies of morphological evolution. One way to overcome this barrier is to rely on the thousands of data points
76 already reported by experts in the scientific literature. However, this method brings its own challenges, such as
77 assigning concordance between taxonomic names and extracting data from published text or images¹¹. To address
78 these needs, we include bioinformatic approaches that can be used by future researchers. Both the egg database and
79 the software solutions used to generate it will have broad value for researchers interested in studying questions of
80 morphological evolution across large evolutionary scales.

81 **3 Methods**

82 **3.1 Gathering primary literature with egg descriptions**

83 The workflow used to assembling the database is shown in Figure 1. Publications were identified for potential inclu-
84 sion in the egg database using the following online literature databases: Google Scholar (scholar.google.com),
85 Web of Knowledge (webofknowledge.com), and Harvard's HOLLIS library system (hollis.harvard.edu).
86 We searched these databases continuously during the period of from October 2015 – August 2017 with a predeter-
87 mined set of word pairs that included an insect common or taxonomic name (e.g. 'fly', 'Diptera', 'Nematocera')
88 and one of the following egg related terms: 'egg', 'chorion', 'immature', or 'embryo'. Insect clade names included
89 all insect order names and all insect families from the five largest insect orders (Coleoptera, Diptera, Lepidoptera,
90 Hymenoptera, and Hemiptera). Following a search, all publications returned by the search were manually eval-
91 uated for inclusion in the database. The criteria for this evaluation were as follows: [1] Does the title or abstract
92 of the paper suggest that the paper contains insect egg information? [2] If the publication could be immedi-
93 ately previewed on the Harvard library system, does it contain an egg measurement in the text or an egg image
94 with a scale bar? [3] If the publication could not be immediately previewed, does the title or abstract refer to

A Assembling a database of published sources



B Extracting data from published sources

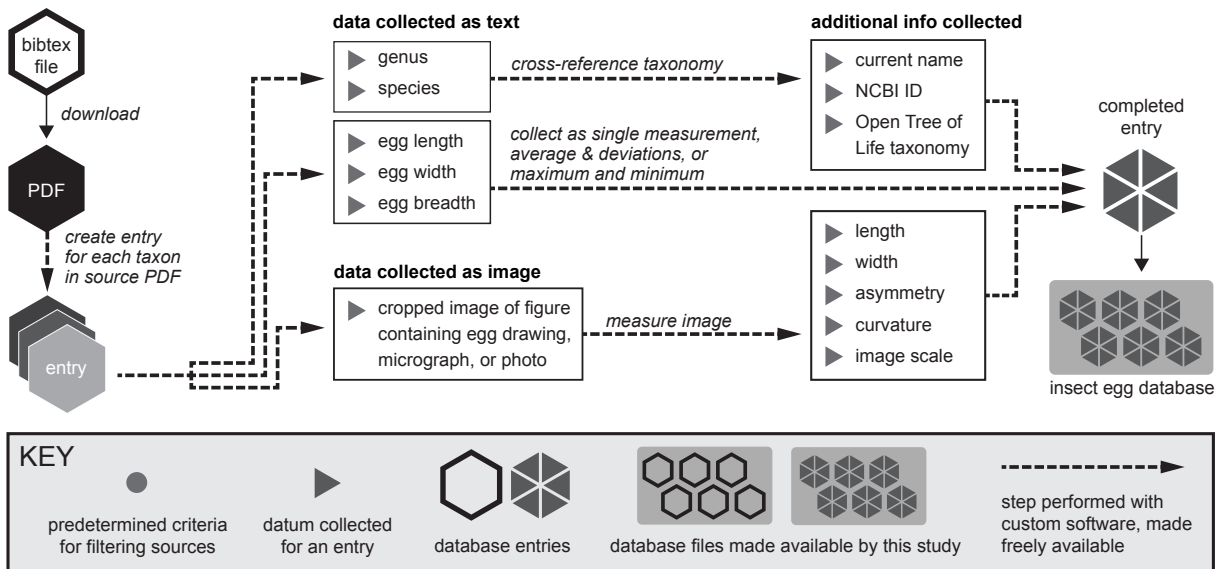


Figure 1: **The workflow used to create the insect egg database.** The database was compiled from the insect literature following the discrete steps shown here, using custom bioinformatic software to maximize reproducibility, consistency, and efficiency. **A**, The workflow used to evaluate candidate publications for inclusion in the database. **B**, The workflow used to extract egg descriptions from the text of published sources and to remeasure published images of eggs. Steps performed with custom software are shown in dashed lines.

95 descriptions of the chorion, immature stages, or embryology? If a publication met at least one of these criteria,
 96 complete bibliographic information for the reference was stored in a master BibTeX reference file (available at Dryad
 97 <https://datadryad.org/review?doi=doi:10.5061/dryad.pv40d2r>). Publications were continually
 98 added to the database throughout the study, and the final count of publications that met these criteria were 2,900,
 99 of which 1,756 contained egg morphological data. The language of the publication was not a criterion for inclusion
 100 in the database. However, due to the nature of the online search engines that we used, the database is enriched for
 101 papers published with at least an abstract in English. A formatted list of the references cited in the egg database is
 102 available in the supplemental file 'bibliography_egg_database'.

103 3.2 Defining egg traits

104 The egg traits in the database are listed in Table 2. For each trait listed below we used the descriptions of egg length
105 and width as presented in the original publications. Given that conventions vary across entomologists and insect
106 taxonomic groups, we present the following definitions to resolve ambiguous cases and to serve as a suggestion for
107 future egg descriptions.

108 *Egg*: The term *egg* is used in the literature to describe several successive developmental stages, including the mature
109 oocyte, the zygote cell, and the developing embryo in its eggshell. For consistency we selected measurements that were
110 recorded closest to the time of fertilization, when multiple descriptions were available within a single publication,
111 given that in some insects it has been documented that the dimensions of the egg change over time (typically <20%
112 change in length due to water exchange during embryonic development)^{7,12–15}. In most insects the egg is oviposited
113 outside the adult body; however in viviparous insects, eggs proceed through some or all of embryonic development
114 within the body of the mother. The egg is often enveloped in a secreted eggshell called the chorion¹⁵, which may
115 have elaborations (e.g. dorsal appendages or opercula)¹⁶. We selected egg measurements that excluded chorionic
116 elaborations over those that included them, as our goal was to measure the comparable cellular material across
117 species.

118 *Length*: To resolve ambiguous cases, and when measuring egg features from published images, we defined egg length
119 as the distance in millimeters (mm) of the axis of rotational symmetry. This definition maximizes consistency with
120 published descriptions of egg length. Under this definition, length is not always longer than width (as defined
121 below). For some insect groups (e.g. Lepidoptera) the axis of rotational symmetry is sometimes referred to in the
122 literature as *height*^{17–19}. For published images with a scale bar, we measured both the straight and curved length of
123 the egg (for those eggs that are curved), but for all analyses and figures, we used the straight length of the egg to
124 maximize consistency with published records.

125 *Width* and *breadth*: To resolve ambiguous cases, and when measuring egg features from images, we defined width
126 as the widest diameter (mm), measured perpendicular to the axis of rotational symmetry of the egg. For some insect
127 groups this axis is referred to in the literature as *diameter*¹⁷ or *breadth*²⁰. For eggs described in published records
128 as having a length, width, and breadth or depth (i.e., the egg is a flattened ellipsoid²¹), we considered *width* as the
129 wider of the two diameters, and *breadth* as the diameter perpendicular to both width and length. For published
130 images with a scale bar, we measured width as the widest of the three egg diameters at the first quartile, midpoint,
131 and third quartile of the length axis. We did not measure breadth from published images.

132 *Volume*: Volume (mm³) was calculated using the equation for the volume of an ellipsoid, following previous
133 studies^{22,23}. The formula is $\frac{1}{6}\pi lwb$, with *l*, *w*, and *b* as length, width, and breadth, respectively. This simplifies to
134 $\frac{1}{6}\pi lw^2$ when the egg is rotationally symmetric. For records in which the volume was reported but egg length and
135 width were not, we used the reported volume. For all other entries, we recalculated volume from the measurements
136 in the text and from measurements of images published with a scale bar.

137 *Aspect ratio*: We calculated aspect ratio as the ratio of length to width. An aspect ratio of one corresponds to a
138 spherical egg. An aspect ratio less than one corresponds to an egg that is wider than long (oblate ellipsoid). An aspect

139 ratio greater than one corresponds to an egg that is longer than it is wide (prolate ellipsoid). Analyses testing the
140 sensitivity of our measurement software (see “Assessing the accuracy of image measuring software” below) for egg
141 images indicated that the variance in measured aspect ratio increases sharply when aspect ratio is much higher than
142 typical (Table 3). Therefore we excluded the eggs in the top 0.1 percentile of aspect ratio from the final database. We
143 recorded the aspect ratio from images published with or without a scale bar, as aspect ratio is a scale-free attribute.

144 *Asymmetry:* We defined asymmetry as $\frac{\max(q_1, q_3)}{\min(q_1, q_3)} - 1$, where q_1 and q_3 are the egg diameters at the first and third
145 quartile of the curved length axis. Therefore an egg with an asymmetry of zero has quartile diameters with equal
146 length. Baker’s λ value, used to measure asymmetry in bird eggs²⁴, can be converted to the asymmetry parameter
147 used in the present study. Analyses testing the sensitivity of our image measuring software (see “Assessing the
148 accuracy of image measuring software” below) indicated that the variance increases sharply near the extreme high
149 values of asymmetry (Table 3). We therefore excluded the eggs in the top 0.1 percentile of asymmetry from the final
150 database. Asymmetry was only recorded from published egg images.

151 *Angle of curvature:* We defined the angle of egg curvature as the angle of the arc (measured in degrees) created by the
152 endpoints and midpoint of the length axis. Analyses testing the sensitivity of our image measuring software (see
153 “Assessing the accuracy of image measuring software” below) indicated that the variance in curvature increases when
154 the curvature and aspect ratio are low (Table 3). We therefore did not calculate curvature for eggs with an aspect
155 ratio of one or less. Angle of curvature was only recorded from published egg images.

156 3.3 Extracting egg descriptions from text sources

157 Information was extracted from publications using a custom text parsing tool that automatically opened and
158 searched the text of a PDF of the publication (https://github.com/shchurch/Insect_Egg_Evolution,
159 file ‘parsing_eggs.py’, commit bd765c8). The tool, written in Python2, uses a text scoring formula to identify
160 candidate blocks of text that contain egg descriptions and corresponding names. Each database entry was manually
161 verified and stored in tab delimited format.

162 All entries included, at a minimum, a genus name and an egg measurement in one dimension or egg volume.
163 Measurements were recorded as either an average and deviation, a range of measurements, or a single value, with
164 precedence for inclusion given in that order. A text description of the volume of the egg was included only in cases in
165 which there were no available data on the linear dimensions of the egg. The majority of the descriptions are reported
166 as single values (Table 1).

167 3.4 Measuring published images of eggs

168 Published images of eggs were measured using a custom tool ([https://github.com/sdonoughe/Insect_](https://github.com/sdonoughe/Insect_Egg_Image_Parser)
169 [Egg_Image_Parser](https://github.com/sdonoughe/Insect_Egg_Image_Parser), commit faee2e8) that enabled the user to calculate aspect ratio, curvature, and asymmetry of
170 the egg by dropping guided landmarks on the published egg image (Figure 2). If the published image included a
171 scale bar, the program also measured the absolute length and width of the egg. The final output of this tool was
172 combined with the corresponding text description of the egg of that species. Images were included regardless of

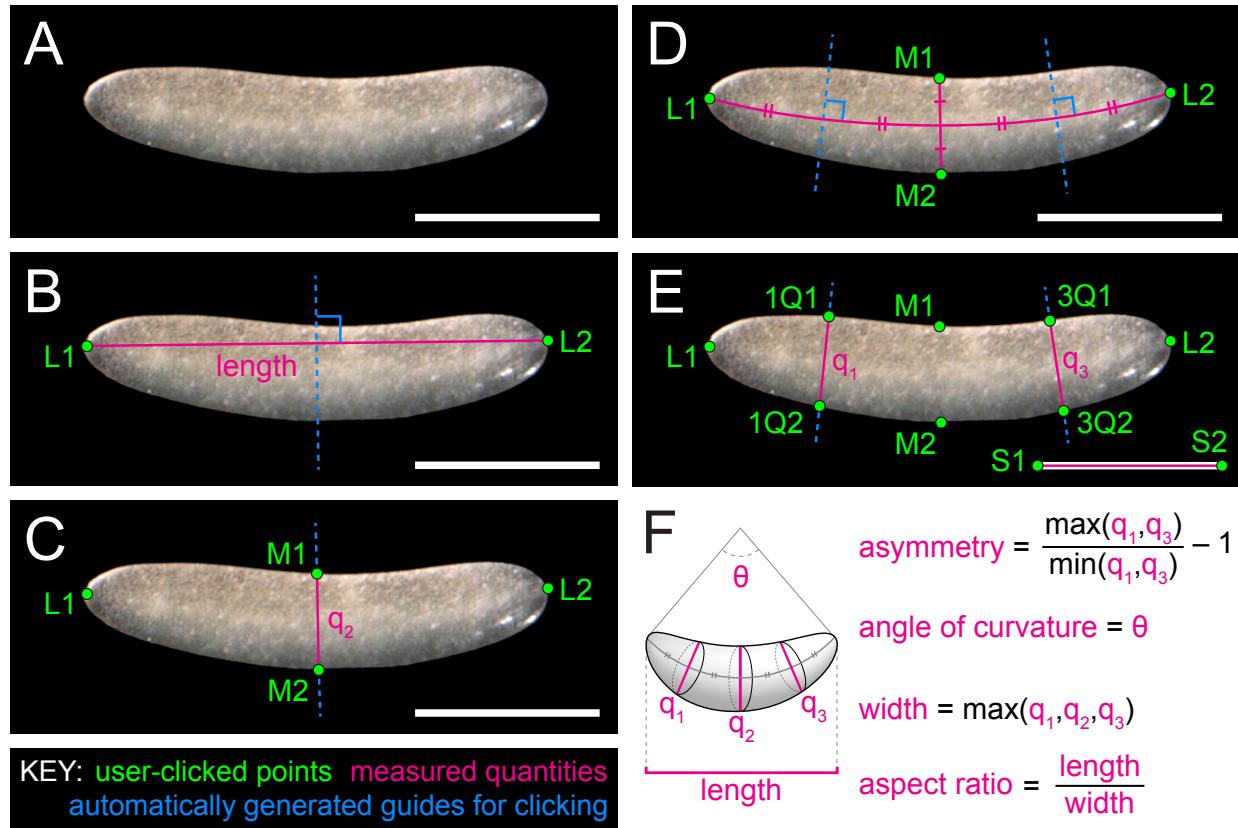


Figure 2: **Demonstration of guided landmark-based measurement of egg shape traits.** **A**, An example micrograph of an egg, in this case from the cricket *Gryllus bimaculatus*. **B**, The user places points L1 and L2 at the poles of the egg. We define egg ‘poles’ as the points on opposite sides of the egg where the curvature of the egg margin is steepest. The tool draws a line segment connecting L1 and L2 (length) and then draws its perpendicular bisector (dashed blue line). **C**, The user uses the blue line as a guide to place points M1 and M2 where the line meets the egg margin. The tool draws a line segment connecting M1 and M2 (q_2). **D**, The tool draws a curved segment connecting the midpoint of q_1 with L1 and L2, and then draws two perpendicular bisectors of the curved segment (dashed blue lines). **E**, The user uses the blue lines as a guide to place points 1Q1, 1Q2, 3Q1, and 3Q2 where the lines meet the egg margin. The tool draws two lines connecting these points (q_1 and q_3). The user places points S1 and S2 at the ends of the scale bar. **F** Collected measurements from this image are as follows: Length is the distance from L1 to L2. Asymmetry is the ratio of the larger distance among q_1 and q_2 to the smaller. Angle of curvature is calculated as the angle formed by points L1, L2 and the midpoint of q_2 . Width is the longest distance between q_1 , q_2 , and q_3 . Aspect ratio is the ratio of length to width. See Table 2 for additional details.

173 type (e.g. light micrograph, scanning electron micrograph, drawing). However, images of low quality were excluded
 174 by manually evaluating cases where landmarks could not be placed unambiguously.

175 3.5 Assessing the accuracy of image measuring software

176 To examine the possible interactions between shape parameters and the accuracy of the image measuring software,
 177 an array of 24 egg silhouettes were simulated with combinations of known parameter values (Figure 3). Each of
 178 these eggs was measured five times with the custom image measurement tool to calculate aspect ratio, asymmetry,
 179 and the angle of curvature (Table 3).

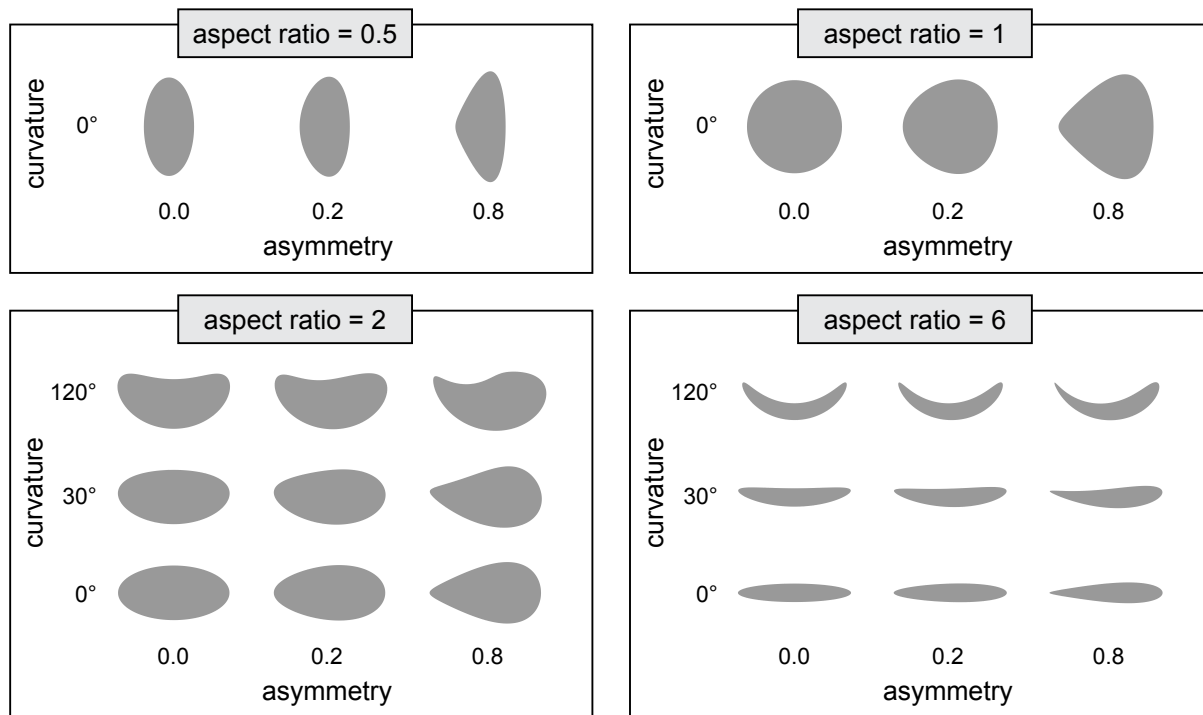


Figure 3: **Assessing the accuracy of the egg image measuring software.** Simulated egg silhouettes with known combinations of shape parameter values used to assess accuracy of image measurement software. Each egg was remeasured five times using the image measurement software and the results are reported in Table 3.

180 3.6 Calculating final and transformed values

181 Following data extraction from text and image sources, final values (e.g. volume, aspect ratio) were calculated. Egg
182 length, width, breadth, volume, aspect ratio were log₁₀ transformed, and egg curvature and asymmetry were square
183 root transformed. For entries that had both a text description of egg size as well as an image with a scale bar, the
184 text description was used in the final calculations. Both the raw and processed final database are freely available for
185 download (Dryad <https://datadryad.org/review?doi=doi:10.5061/dryad.pv40d2r>).

186 3.7 Cross-referencing entries with taxonomic and genetic databases

187 Taxonomic names parsed from the literature occasionally contained errors, including published typographical
188 errors and optical character recognition errors. These errors needed to be corrected and the taxonomic names
189 also had to be reconciled with currently accepted taxonomy in order to link egg morphology data with other data
190 sources (e.g. published phylogenies). To address these issues, we developed a tool called TaxReformer (<https://github.com/brunoasm/TaxReformer>, commit 1831a11) that searches the Global Names Architecture
191 (GN)^{25,26}, Open Tree Taxonomy (OTT)^{27,28}, and Global Biodiversity Information Facility (GBIF)²⁹ databases,
192 taking advantage of the strengths of each database. For the taxa included in the insect egg database, GN had the
193 most effective fuzzy matching algorithm and broadest database. OTT provided a better control of the context of
194

195 each taxonomic query, enabling one to search names only among insects and avoiding homonyms in kingdoms
196 regulated by different codes of nomenclature. OTT's fuzzy matching algorithm, however, often returned matches
197 to the correct species name but wrong genus name with a high confidence score. OTT and GBIF both contain
198 information about higher taxonomy, which is not standardized in records obtained from GN.

199 Names obtained from the literature were first parsed with Global Names Parser v. 0.3.1³⁰ to obtain genus and
200 species name in canonical forms. The full species name was then used to search in GN with fuzzy matching to
201 allow for correction of optical character recognition errors. If a match to a species or genus was found, the matched
202 name was recorded and then searched in OTT to obtain higher taxonomy and identifier numbers from OTT and
203 the National Center for Biotechnology Information. If the name was not found in OTT, higher taxonomy was
204 alternatively obtained from GBIF. In all cases, if databases contained information about synonyms, the currently
205 accepted name for each taxon was retrieved.

206 **3.8 Assessing intraspecific variation**

207 We assessed intraspecific variation in egg size descriptions using four methods:

208 First, for database entries that reported egg size variation (e.g. egg descriptions that included a range of egg length or
209 an average egg length with deviation), the percent difference in egg size was calculated as follows: for egg descriptions
210 recorded as ranges, percent difference was calculated as $100 * \frac{\max l - \min l}{\text{median} l}$; for egg descriptions recorded as average
211 and deviations, percent difference was calculated as $100 * \frac{(2 * \text{deviation})}{\text{mean} l}$.

212 Second, independent observations of a single species were identified as two entries for the same species that differed
213 in the calculated volume by more than $1.0 * 10^{-5} \text{ mm}^3$. This excluded entries that were repeated publications of
214 the same description, such as an observation repeated in a subsequent review (Table 1). The percent difference in
215 egg length was calculated as $100 * \frac{\max l - \min l}{\text{median} l}$.

216 Third, for entries that had both a text description of egg length as well as a published image with a scale bar, the
217 difference in the reported egg length and our re-measurement of the image was assessed. The percent difference
218 between these two measurements was calculated as $100 * \frac{\max l - \min l}{\text{median} l}$.

219 Fourth, for eggs that were measured as triaxial ellipsoids (length, width, and breadth measured all separately), the
220 percent difference was calculated from the change in egg volume if the egg had been assumed to be a rotationally
221 symmetric ellipsoid (volume = $\frac{1}{6} \pi l w b$ vs volume = $\frac{1}{6} \pi l w^2$). Given that more eggs are likely triaxial ellipsoids than
222 are reported in the egg database, this metric gives insight into the variation in egg volume that might be masked
223 when only two dimensions are reported.

224 **3.9 Assessing the precision of entries**

225 The distribution of precision in the insect egg database was assessed using two metrics. First, the number of decimal
226 places used in the length measurement was calculated for each database entry from a base of millimeters (e.g. '1 mm'
227 has 0 decimal places, while '1.00 mm' has 2 decimal places). Second, the relative precision of each measurement was

228 calculated by dividing the total length of the egg by the smallest unit used to measure it, and multiplying this value
229 by 100. This gives the percent of egg length captured by the unit of measurement (i.e. an egg measured as 1.00 mm
230 was measured within 1% of egg length).

231 **3.10 Assessing the phylogenetic sampling**

232 The phylogenetic coverage of the insect egg database was assessed by comparing the number of egg entries for a
233 taxonomic rank to the number of species in that rank, estimated by the number of tips in the Open Tree of Life²⁸.
234 This assay was performed for all extant hexapod orders and for all insect families in the insect egg database.

235 **4 Code availability**

236 All code used to generate the insect egg database as well as reproduce the tables and plots shown here is made freely
237 available. Python code used to compile the database and extract text information from text sources, as well as the
238 R code used to convert the raw database to the final database and to generate the tables and figures shown here
239 is available at https://github.com/shchurch/Insect_Egg_Evolution. Python code used to measure
240 published images of eggs is available at https://github.com/sdonoughe/Insect_Egg_Image_Parser,
241 and python code to cross-reference the egg database with taxonomic tools is available at <https://github.com/brunoasm/TaxReformer>. Statistical analyses were performed using R version 3.4.2³¹.

243 **5 Data records**

244 The final data files include the raw database in tab delimited format, which includes all values extracted from the text
245 and images, as well as the final database in tab delimited format. The code to convert the raw database to the final
246 database is located in https://github.com/shchurch/Insect_Egg_Evolution, directory ‘analyze_data’.
247 Additionally, all data files have been uploaded to Dryad <https://datadryad.org/review?doi=doi:10.5061/dryad.pv40d2r>.

249 **6 Technical validation**

250 The accuracy of the image measuring software was assessed using an array of 24 simulated egg silhouettes with
251 known combinations of parameter values (Figure 4). We found that as the actual angle of curvature increases,
252 the difference between the actual and measured values increases (that is, the software underestimates the angle of
253 curvature), and this difference is larger in eggs with lower aspect ratio and higher asymmetry (Table 3). As the actual
254 asymmetry increases the variance in measured asymmetry increases, and in eggs with low aspect ratio this results in
255 an overestimation of asymmetry. As the actual aspect ratio increases, the software overestimates the total aspect

256 ratio by up to 0.75 (12.5% of the total aspect ratio). Given these results we removed eggs in the top 0.1 percentile of
257 values for asymmetry and aspect ratio when creating the final database.

258 Intraspecific variation in insect egg size was assessed using four metrics (see Methods section “Assessing intraspecific
259 variation”). The first two describe the percent difference in egg size reported in the literature, either as variation
260 recorded in an egg description (Figure 4A), or as variation recorded across multiple independent observations of
261 eggs from the same species (Figure 4B). In both cases the percent difference in egg length averaged 10% and ranged
262 from 1% to 100% (i.e., for an insect species with an average egg length of 1 mm, it was common to observe eggs from
263 0.9 to 1.1 mm and occasional outliers at 0.5 and 2 mm).

264 Additionally we re-measured published images of eggs and calculated the percent difference between our measure-
265 ments and the text description (Figure 4C). The variation between observations of the same species was consistent
266 with the reported intraspecific variation (average around 10%).

267 Although the majority of eggs in the database are described as rotationally symmetric ellipsoids (Table 1), for a
268 few clades of insects it is common to measure eggs as triaxial ellipsoids, with length, width, and breadth measured
269 separately (Table 2). Calculating the egg volume using two different methods — one taking into account breadth,
270 and the other assuming rotational symmetry — showed that the percent difference in calculated volume ranges
271 between 10% and 100% (Figure 4D). Eggs from additional clades might be more accurately modeled as triaxial
272 ellipsoids than currently reported in the literature, but this percent difference likely represents the upper range of
273 the error in volume, because the clades typically measured as triaxial ellipsoids are those that are most obviously
274 flattened along one axis.

275 The text descriptions in the insect egg database were extracted from a diverse set of sources published over hundreds
276 of years, and the precision used to measure eggs varies across these sources (Figure 4). Most entomologists measured
277 eggs in tenths or hundredths of a millimeter (Figure 4E). In terms of the total length of the egg, most measurements
278 in the database are precise to within 1% to 10% (Figure 4F). Given that intraspecific variation is also around 10% of
279 total egg length, it is likely that some of this variation is due to measurement error.

280 The egg database contains descriptions of eggs from every insect order and from hundreds of insect families (Table
281 1). Given that the number of species varies greatly across taxonomic ranks we assessed the phylogenetic coverage of
282 the egg database (Figure 4G, H). We found that families and orders with the highest number of estimated species are
283 represented by the greatest number of entries in the egg database. Additionally, most families in the egg database
284 have more than 1 entry per 100 species.

285 There are several orders represented in the database by fewer than ten entries (Figure 4H). We suggest that this is
286 likely due in part to idiosyncracies of the entomological research for certain clades. For example, although many
287 descriptions of mantis and cockroach oothecae exist, measurements or images of individual eggs within the oothecae
288 are rare in the published literature, which leaves these groups undersampled for propagule size in the literature. The
289 orders with the lowest representation—Trichoptera, Psocoptera, and Zygentoma—are potentially rich new datasets
290 to target for future study.

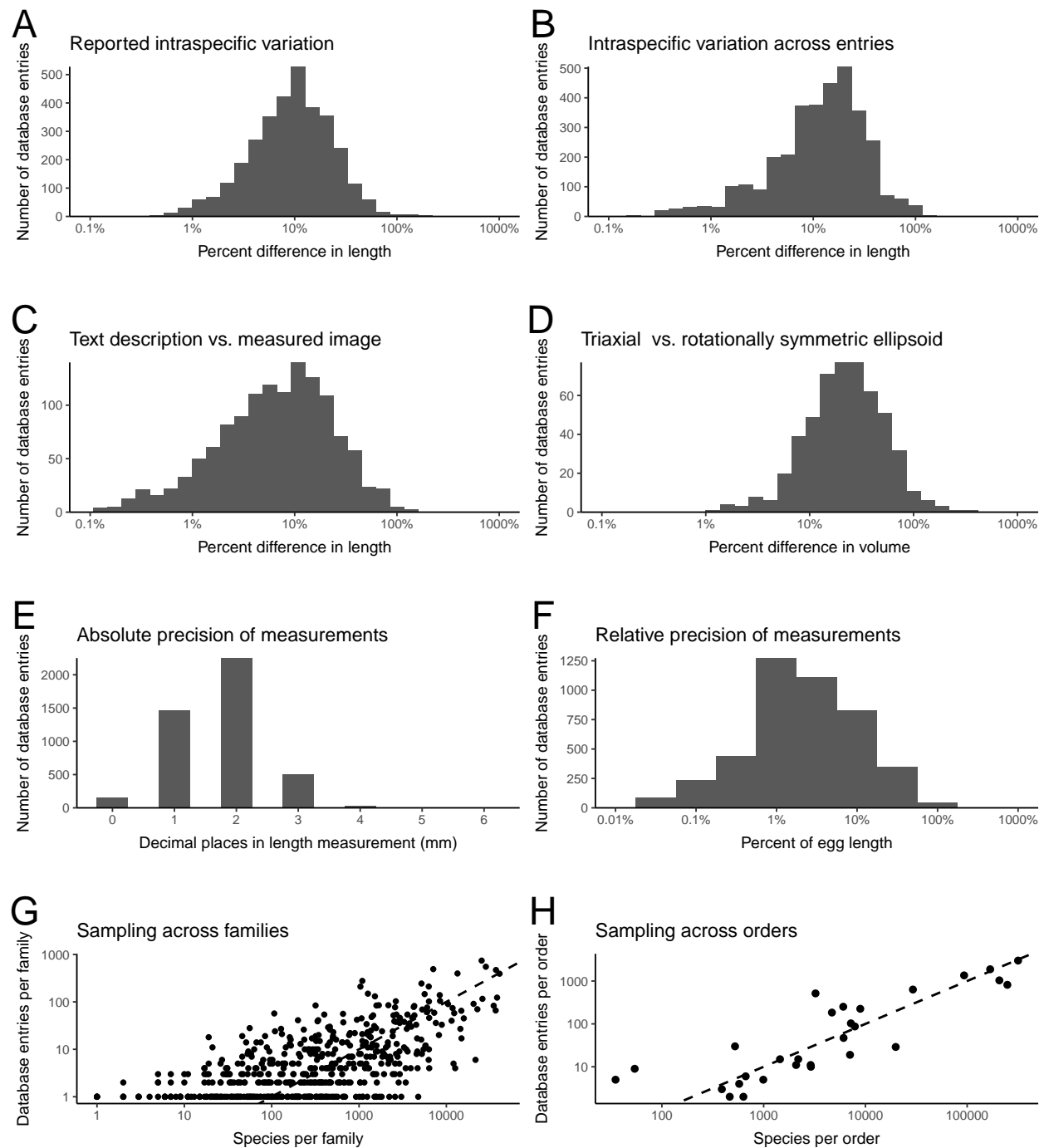


Figure 4: Assessing intraspecific variation, precision, and sampling within the insect egg database. **A**, The distribution of the percent difference between the largest and smallest egg length reported for a species within a publication. **B**, The distribution of the percent difference between the largest and smallest egg length reported for a species across different publications. **C**, The distribution of the percent difference between the largest and smallest egg length, comparing the reported length and the remeasured image from the same publication. **D**, The distribution of the percent difference between the largest and smallest egg volume, measured as triaxial ellipsoids (length, width, and breadth) vs. rotationally symmetric ellipsoids (length and width). **E**, The distribution of the absolute precision of each measurement (decimal places in the egg length measurement in millimeters). **F**, The distribution of the relative precision of each measurement (percent of egg length of the smallest unit used to measure insect egg length). **G**, A comparison of the number of database entries to the number of species estimated in every family present in the insect egg database. **H**, A comparison of the number of database entries to the number of species estimated in every extant insect order. In **H-I** the dotted line shows an arbitrary standard of 1 entry per 100 estimated species.

291 **7 Acknowledgements**

292 This work was supported by the National Science Foundation (NSF) Grant No. IOS-1257217 to CGE, NSF
293 Graduate Research Fellowship No. DGE1745303 to SHC, and by a Jorge Paulo Lemann Fellowship to BdM
294 from Harvard University. We acknowledge Jordan Hoffman and Casey W. Dunn for initial code advice and
295 troubleshooting. We thank the Extavour lab and Brian Farrell for discussion, and Arpita Kulkarni, Angela de Pace,
296 Benjamin Goulet, and Tarun Kumar for suggestions on initial versions of this manuscript. We acknowledge the
297 Ernst Mayr Library at the Museum of Comparative Zoology at Harvard, and specifically Mary Sears, for countless
298 hours of support in gathering the references used in this study.

299 **8 Author contributions**

300 SHC and SD wrote all code to parse egg descriptions from the literature, and contributed equally to database
301 creation, study design, writing, and figure preparation. SHC wrote code to manipulate the database and perform
302 statistical analyses. SD wrote code to measure published images. BdM wrote code to correct taxonomic information.
303 BdM and CGE contributed to study design, interpretation, and writing.

304 **9 Competing interests**

305 The authors declare no competing interests.

Bibliographic statistics	
references examined	2900
references with egg information	1756
unique authors	1498
unique journals / books	491

Data type statistics	
Total entries in egg database	10449
Entries with text description of length and width	7672
Length reported as average and deviation	1065
Length reported as range	2188
Single length value reported	4419
Only volume reported	1368
Entries with an image	4774
Images re-measured	2004
Entries with both text and image measurements	1205

Taxonomic statistics	
unique hexapod species	6706
unique hexapod genera	4077
unique hexapod families	526
unique hexapod orders	32

Table 1: **Bibliographic, data type, and taxonomic statistics of the insect egg database**

Text measurements		Standardized text measurements		
Name	Units	Name	Units	Method
length or height	mm	length, l	mm	as recorded
width or diameter	mm	width, w	mm	$\max(w, b)$
breadth or depth	mm	breadth, b	mm	$\min(w, b)$
volume*	mm ³	volume, v	mm ³	$\frac{1}{6}\pi lwb$ OR $\frac{1}{6}\pi lw^2$ OR v
		aspect ratio	ratio, no units	$\frac{l}{w}$

Image measurements		Standardized image measurements		
Name	Units	Name	Units	Method
curved length	px	length**, l	mm	straight length
straight length	px	width**, w	mm	$\max(q_1, q_2, q_3)$
1st quartile width, q_1	px	volume**	mm ³	$\frac{1}{6}\pi lw^2$
2nd quartile width, q_2	px	aspect ratio	ratio, no units	$\frac{l}{w}$
3rd quartile width, q_3	px	asymmetry	ratio, no units	$\frac{\max(q_1, q_3)}{\min(q_1, q_3)} - 1$
angle of curvature	degrees, radians	angle of curvature	radians	as recorded

Final database measurements			
Name	Units	Transformation	Method
length	mm	\log_{10}	used text measurement, when both text and image were available
width	mm	\log_{10}	used text measurement, when both text and image were available
breadth	mm	\log_{10}	used text measurement, when both text and image were available
volume	mm ³	\log_{10}	used text measurement, when both text and image were available
aspect ratio	ratio, no units	\log_{10}	used text measurement, when both text and image were available, removed egg images in the top 0.1%
asymmetry	ratio, no units	sq. root	removed egg images in the top 0.1%
angle of curvature	radians	sq. root	did not record for eggs with an aspect ratio ≤ 1

Table 2: **Trait definitions and standardizations** * volume was included only when length and width measurements were not available from text. ** measurements included only when a scale bar was published with the image.

Aspect ratio	Actual value		Aspect ratio	Mean discrepancy	
	Asymmetry	Angle of curvature		Asymmetry	Angle of curvature
0.5	0	0	-0.01	-0.05	
0.5	0.2	0	-0.01	-0.08	
0.5	0.8	0	-0.02	0.02	
1	0	0	-0.02	-0.05	
1	0.2	0	-0.03	-0.07	
1	0.8	0	-0.03	-0.13	
2	0	0	-0.03	-0.04	-2.68
2	0	30	-0.06	-0.04	8.74
2	0	120	-0.18	-0.05	15.49
2	0.2	0	-0.06	-0.05	-2.99
2	0.2	30	-0.05	-0.07	6.66
2	0.2	120	-0.17	-0.02	16.75
2	0.8	0	-0.09	-0.08	-0.65
2	0.8	30	-0.10	-0.14	15.02
2	0.8	120	-0.18	-0.06	23.84
6	0	0	-0.36	-0.06	-1.63
6	0	30	-0.15	-0.04	-1.47
6	0	120	-0.32	-0.05	2.52
6	0.2	0	-0.24	-0.06	-0.66
6	0.2	30	-0.50	-0.19	-0.80
6	0.2	120	-0.45	-0.06	3.32
6	0.8	0	-0.36	-0.25	-2.61
6	0.8	30	-0.56	-0.13	-0.16
6	0.8	120	-0.40	-0.14	2.28

Table 3: Results of image measurement software accuracy assessment. Mean discrepancy calculated as the average difference between the actual and measured values, n = 5.

References

- 307 1. Smith, C. C. & Fretwell, S. D. The optimal balance between size and number of offspring. *The American*
308 *Naturalist* **108**, 499–506 (1974).
- 309 2. Bernardo, J. The particular maternal effect of propagule size, especially egg size: patterns, models, quality of
310 evidence and interpretations. *American Zoologist* **36**, 216–236 (1996).
- 311 3. Fox, C. W. & Czesak, M. E. Evolutionary ecology of progeny size in arthropods. *Annual Review of Entomology*
312 **45**, 341–369 (2000).
- 313 4. Berrigan, D. The allometry of egg size and number in insects. *Oikos* **60**, 313 (1991).
- 314 5. García-Barros, E. Body size, egg size, and their interspecific relationships with ecological and life history traits in
315 butterflies (Lepidoptera: Papilionoidea, Hesperioidea). *Biological Journal of the Linnean Society* **70**, 251–284
316 (2000).
- 317 6. Blackburn, T. M. *Comparative and experimental studies of animal life history variation* PhD thesis (University
318 of Oxford, 1990).
- 319 7. Hinton, H. E. *Biology of insect eggs* (Pergammon Press, Oxford, 1981).
- 320 8. Legay, J. M. Allometry and systematics of insect egg form. *Journal of Natural History* **11**, 493–499 (1977).
- 321 9. Misof, B. *et al.* Phylogenomics resolves the timing and pattern of insect evolution. *Science* **346**, 763–767
322 (2014).
- 323 10. Rainford, J. L., Hofreiter, M., Nicholson, D. B. & Mayhew, P. J. Phylogenetic distribution of extant richness
324 suggests metamorphosis is a key innovation driving diversification in insects. *PLoS One* **9**, 1–7 (2014).
- 325 11. Dahdul, W. M. *et al.* Evolutionary characters, phenotypes and ontologies: curating data from the systematic
326 biology literature. *PLoS One* **5**, e10708 (2010).
- 327 12. Kobayashi, Y. Embryogenesis of the fairy moth, *Nemophora albiantennella* Issiki (Lepidoptera, Adelidae), with
328 special emphasis on its phylogenetic implications. *International Journal of Insect Morphology and Embryology*
329 **27**, 157–166 (1998).
- 330 13. Chaves, L. F., Ramoni-Perazzi, P., Lizano, E. & Añez, N. Morphometrical changes in eggs of *Rhodnius prolixus*
331 (Heteroptera: Reduviidae) during development. *Entomotropica* **18**, 83–88 (2003).
- 332 14. Donoughe, S. & Extavour, C. G. Embryonic development of the cricket *Gryllus bimaculatus*. *Developmental*
333 *Biology* **411**, 140–156 (2016).
- 334 15. Rezende, G. L., Vargas, H. C. M., Moussian, B. & Cohen, E. in *Extracellular composite matrices in arthropods*
335 325–366 (Springer, Cham, 2016).
- 336 16. Hinton, H. Respiratory systems of insect egg shells. *Annual Review of Entomology* **14**, 343–368 (1969).
- 337 17. Dolinskaya, I. V. Comparative morphology on the egg chorion characters of some Noctuidae (Lepidoptera).
338 *Zootaxa* **4085**, 374–392 (2016).
- 339 18. Dahlan, A. & Gordh, G. Development of *Trichogramma australicum* Girault (Hymenoptera: Trichogram-
340 matidae) in eggs of *Helicoverpa armigera* Hiibner (Lepidoptera: Noctuidae) and in artificial diet. *Austral*
341 *Entomology* **37**, 254–264 (1998).
- 342 19. Zompro, O., Adis, J. & Weitschat, W. A review of the order Mantophasmatodea (Insecta). *Zoologischer*
343 *Anzeiger-A Journal of Comparative Zoology* **241**, 269–279 (2002).

- 344 20. Duffy, E. A. J. *A monograph of the immature stages of oriental timber beetles (Cerambycidae)* (The British
345 Museum (Natural History), London, 1968).
- 346 21. Clark, J. T. The eggs of stick insects (Phasmida): a review with descriptions of the eggs of eleven species.
347 *Systematic Entomology* **1**, 95–105 (1976).
- 348 22. Markow, T. A., Beall, S. & Matzkin, L. M. Egg size, embryonic development time and ovoviviparity in
349 *Drosophila* species. *Journal of Evolutionary Biology* **22**, 430–434 (2009).
- 350 23. García-Barros, E. Egg size in butterflies (Lepidoptera: Papilionoidea and Hesperioidea): a summary of data.
351 *Journal of Research on the Lepidoptera* **35**, 90–136 (2000).
- 352 24. Stoddard, M. C., Yong, E. H., Akkaynak, D., Sheard, C., Tobias, J. A. & Mahadevan, L. Avian egg shape:
353 Form, function, and evolution. *Science* **356**, 1249–1254 (2017).
- 354 25. Patterson, D., Mozzherin, D., Shorthouse, D. P. & Thessen, A. Challenges with using names to link digital
355 biodiversity information. *Biodiversity Data Journal* (2016).
- 356 26. Pyle, R. L. Towards a global names architecture: The future of indexing scientific names. *ZooKeys* **2016**,
357 261–281 (2016).
- 358 27. Rees, J. & Cranston, K. Automated assembly of a reference taxonomy for phylogenetic data synthesis. *Biodi-
359 versity Data Journal* **5**, e12581 (2017).
- 360 28. Hinchliff, C. E. *et al.* Synthesis of phylogeny and taxonomy into a comprehensive tree of life. *Proceedings of
361 the National Academy of Sciences of the United States of America* **112**, 12764–12769 (2015).
- 362 29. GBIF. *GBIF: The Global Biodiversity Information Facility* 2018.
- 363 30. Mozzherin, D. Y., Myltsev, A. A. & Patterson, D. J. “gnparser”: A powerful parser for scientific names based
364 on Parsing Expression Grammar. *BMC Bioinformatics* **18**, 1–14 (2017).
- 365 31. R Core Team. *R: A Language and Environment for Statistical Computing* Vienna, Austria: R Foundation for
366 Statistical Computing, 2017. <https://www.R-project.org/>.