

1

2 DR JOEL LUKE PICK (Orcid ID : 0000-0002-6295-3742)

3 DR DANIEL NOBLE (Orcid ID : 0000-0001-9460-8743)

4

5

6 Article type : Application

7 Editor : Samantha Price

8

9

10 **Reproducible, flexible and high-throughput data extraction from primary**  
11 **literature: The *metaDigitise R* package**

12

13 Joel L. Pick<sup>1,2,\*</sup>, Shinichi Nakagawa<sup>1</sup>, Daniel W.A. Noble<sup>1</sup>

14

15 <sup>1</sup> Ecology and Evolution Research Centre, School of Biological, Earth and  
16 Environmental Sciences, University of New South Wales, Kensington, NSW  
17 2052, Sydney, Australia

18

19 <sup>2</sup> Current Address: Institute of Evolutionary Biology, School of Biological  
20 Sciences, University of Edinburgh, Edinburgh EH9 3JT, United Kingdom

21

22 \*Corresponding Author: joel.l.pick@gmail.com

23

24 **Running Head:** Data extraction from figures with *metaDigitise*

25

26

27 **Abstract**

28 1. Research synthesis, such as comparative and meta-analyses, requires the  
29 extraction of effect sizes from primary literature, which are commonly  
30 calculated from descriptive statistics. However, the exact values of such

This is the author manuscript accepted for publication and has undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the Version of Record. Please cite this article as [doi: 10.1111/2041-210X.13118](https://doi.org/10.1111/2041-210X.13118)

This article is protected by copyright. All rights reserved

31 statistics are commonly hidden in figures.

32

33 2. Extracting descriptive statistics from figures can be a slow process that is  
34 not easily reproducible. Additionally, current software lacks an ability to  
35 incorporate important meta-data (e.g., sample sizes, treatment / variable  
36 names) about experiments and is not integrated with other software to  
37 streamline analysis pipelines.

38

39 3. Here we present the *R* package *metaDigitise* which extracts descriptive  
40 statistics such as means, standard deviations and correlations from four  
41 plot types: 1) mean/error plots (e.g. bar graphs with standard errors), 2)  
42 box plots, 3) scatter plots and 4) histograms. *metaDigitise* is user-friendly  
43 and easy to learn as it interactively guides the user through the data  
44 extraction process. Notably, it enables large-scale extraction by  
45 automatically loading image files, letting the user stop processing, edit and  
46 add to the resulting data-frame at any point.

47

48 4. Digitised data can be easily re-plotted and checked, facilitating  
49 reproducible data extraction from plots with little inter-observer bias. We  
50 hope that by making the process of figure extraction more flexible and easy  
51 to conduct it will improve the transparency and quality of meta-analyses in  
52 the future.

53

54 Keywords: meta-analysis, comparative analysis, data extraction, *R*,  
55 reproducibility, figures, images, descriptive statistics

56

## 57 1 Introduction

58 In many different contexts, researchers make use of data presented in  
59 primary literature. In the fields of ecology and evolution (E&E), these data are  
60 most commonly used for comparative and meta-analyses. The use of meta-  
61 analysis in E&E in particular, is rapidly growing, not only in terms of the  
62 number of meta-analyses (in plant ecology alone the yearly number of  
63 published meta-analyses doubled from 2006 to 2012 (20-40) (Koricheva &  
64 Gurevitch, 2014)), but also in terms of their size (a recent meta-analysis, for

65 example, included 6440 effect sizes from 175 publications (Noble, Stenhouse  
66 & Schwanz, 2018)). Meta-analyses are extremely important in providing a  
67 means of quantitatively synthesizing experimental and/or observational  
68 studies to evaluate empirical support for fundamental theory in E&E  
69 (Gurevitch et al., 2018). These techniques rely heavily on descriptive statistics  
70 (e.g. means, standard deviations (SD), sample sizes, correlation coefficients)  
71 extracted from primary literature. As well as being presented in the text or  
72 tables of research papers, descriptive statistics are frequently presented in  
73 figures. For example, 42% of the papers used in a recent meta-analysis  
74 presented some or all of the required data in figures (Noble, Stenhouse &  
75 Schwanz, 2018). These data need to be manually extracted using digitising  
76 programs.

77 Although there are several tools that extract data from figures, including  
78 both standalone programs and *R* packages (reviewed in Table 1), these tools  
79 do not cater to the general needs of meta-analysts for four main reasons  
80 (here we focus on meta-analysis, although many points apply to extraction for  
81 comparative analysis). First, although meta-analysis is an important tool in  
82 consolidating the data from multiple studies, many of the processes involved  
83 in data extraction are opaque and difficult to reproduce, making extending or  
84 replicating studies problematic. Having a tool that facilitates reproducibility in  
85 meta-analyses will increase transparency and aid in resolving the  
86 reproducibility crises seen in many fields (Peng, Dominici & Zeger, 2006;  
87 Peng, 2011; Parker et al., 2016). Second, digitising programs do not allow the  
88 integration of metadata at the time of data extraction, such as experimental  
89 group or variable names, and sample sizes. This makes the downstream  
90 calculations laborious, as information has to be added later, typically using  
91 different software. Third, existing programs do not import sets of images for  
92 the user to systematically work through. Instead they require the user to  
93 manually import images and export the resulting digitised data into individual  
94 files one-by-one. These data often subsequently need to be imported and  
95 edited using different software. Finally, digitising programs typically only  
96 provide the user with calibrated x,y coordinates from imported figures, and do  
97 not differentiate between common plot types that are used to present data.  
98 Consequently, a large amount of additional data manipulation is required, that

99 is different across plots types. For example, in E&E data are commonly  
100 presented in plots with means and standard errors or confidence intervals  
101 (Figure 1A), from which the user wants a mean and SD for each group  
102 presented. From x,y coordinates, users must manually discern between mean  
103 and error coordinates and assign points to groups. The error then needs to be  
104 calculated as the deviation from the mean, and then transformed to SD,  
105 according to the type of error presented. Histograms and box plots are also  
106 frequently used in E&E to presented data, and whilst their downstream  
107 calculations are even more laborious, there are few (if any; see Table 1) tools  
108 to extract data from these plot types.

109 Data extraction from figures is therefore a time-consuming process as  
110 existing software does not provide an optimized, reproducible research  
111 pipeline to facilitate data extraction and editing. Given the ubiquity of the *R*  
112 platform in E&E, and that it hosts the most popular meta-analysis software in  
113 E&E (e.g., *metafor* (Viechtbauer, 2010) and *MCMCglmm* (Hadfield, 2010)), it  
114 is highly likely to be used for some (if not all) stages of the research synthesis  
115 process. It is therefore important to have comprehensive, robust and flexible  
116 digitisation capabilities in *R* to make the process of figure extraction more  
117 streamline, transparent and easier to reproduce. Here, we present an  
118 interactive *R* package, *metaDigitise* (available on *CRAN*), which is designed  
119 for large scale, reproducible data extraction from figures, specifically catering  
120 to the needs of meta-analysts. To this end, we provide tools to extract data  
121 from common plot types in E&E (mean/error plots, box plots, scatter plots and  
122 histograms, see Figure 1). *metaDigitise* operates within the *R* environment  
123 making data extraction, analysis and export more streamlined. The necessary  
124 calculations are carried out on calibrated data immediately after extraction so  
125 that comparable descriptive statistics can be obtained quickly. Summary data  
126 from multiple figures is returned into a single data frame which can be can  
127 easily exported or used in downstream analysis within *R*. Completed  
128 digitisations are automatically saved for each figure, meaning users can  
129 redraw their digitisations (along with metadata) on figures, make corrections  
130 and access calibration and processed (i.e., summarised) data. This makes  
131 sharing figure digitisation and reproducing the work of others simple and easy,  
132 and allows meta-analyses to be updated more efficiently.

133

## 134 **2 *metaDigitise* and Reproducibility**

135 The *metaDigitise* package has one main function, *metaDigitise()*, which  
136 interactively takes the user through the process of extracting data from figures  
137 (see Supplementary Material S1 for a full tutorial). Running *metaDigitise()*  
138 presents the user with three options; 'Process new images', 'Import existing  
139 data' or 'Edit existing data', which can be used during and after digitisation to  
140 execute a range of functions (see Figure 1 – 'Processing images' is discussed  
141 in Section 3, and 'Editing' and 'Importing' in Section 4). *metaDigitise()* works  
142 on a directory containing images of figures copied from primary literature, in  
143 .png, .jpg, .tiff, .pdf format, specified to *metaDigitise()* through the *dir*  
144 argument. *metaDigitise()* recognizes all the images in the given directory and  
145 automatically imports them one-by-one, allowing the user to extract the  
146 relevant information about a figure as they go. Figures can be organised in  
147 different ways for a project, but we would recommend having all figures for  
148 one project in a single directory with an informative and unambiguous naming  
149 scheme (e.g. paper figure trait.png). This expedites digitisation by preventing  
150 users from having to constantly change directories and / or open new images.

151 The data from each completed image is automatically saved as a  
152 *metaDigitise* object in a separate .RDS file to a *caldat* folder that is created  
153 within the parent directory when first executing *metaDigitise()*. These files  
154 enable re-plotting and editing of images at a later point (see below). When  
155 run, *metaDigitise()* also identifies the images within a directory that have been  
156 previously digitised and only imports new images to process. The data of all  
157 images is then automatically integrated into the final output. This means that  
158 all figures do not need to be extracted at one time and new figures can be  
159 added to the directory as the project develops.

160 The complete digitisation process can be reproduced at a later stage,  
161 shared with collaborators and presented as supplementary materials for a  
162 publication, regardless of the computer it is run on. To update an analysis,  
163 new figures can simply be added to the directory and *metaDigitise()* run to  
164 incorporate the new data.

165

## 166 **3 Image Processing**

167           Selecting ‘Process New Images’, after running *metaDigitise()*, starts the  
168 digitisation process on images within the directory that have not previously  
169 been digitised. For all plot types, *metaDigitise()* requires the user to calibrate  
170 the axes in the figure, by clicking on two known points on the axis in question,  
171 and entering the value of those points (Figure 1). *metaDigitise()* then  
172 calculates the value of any clicked points in terms of the figure axes. This is  
173 based on the calibration used in the digitize *R* package (Poisot, 2011). For  
174 mean/error and box plots, only the y-axis is calibrated (Figure 1), assuming  
175 the x-axis is redundant. For scatter plots and histograms both axes are  
176 calibrated (Figure 1).

177           Calibration of points in figures from older, scanned publications can be  
178 problematic, as the figures may not be perfectly orientated. *metaDigitise()*  
179 allows users to rotate the image (Figure S2A,B). Furthermore, mean/error  
180 plots, box plots and histograms may be presented with horizontal bars.  
181 *metaDigitise()* assumes that bars are vertical, but allows the user to flip the  
182 image to make the bars are vertical (Figure S2C,D). *metaDigitise* also allows  
183 back calculation of data presented on log axes.

184           *metaDigitise* recognises four main types of plot; Mean/error plots, box  
185 plots, scatter plots and histograms (Figure 1). All plot types can be extracted  
186 in a single call of *metaDigitise()* and integrated into one output. Alternatively,  
187 users can process different plot types separately, using separate directories.  
188 All four plot types are extracted slightly differently (outlined below). Upon  
189 completing all images, or quitting, either summarised or calibrated data is  
190 returned (specified by the user through the summary argument). Summarised  
191 data consists of a mean, SD and sample size, for each identified group within  
192 the plot (should multiple groups exist). In the case of scatter plots, the  
193 correlation coefficient between x and y variables within each identified group  
194 is also returned. Calibrated data consists of a list with slots for each of the four  
195 figure types, containing the calibrated points that the user has clicked. This  
196 may be particularly useful in the case of scatter plots.

197

### 198 *3.1 Mean/Error and Box Plots*

199           *metaDigitise()* handles mean/error and box plots in a very similar way.  
200 For each mean/box, the user enters group name(s) and sample size(s). If the

201 user does not enter a sample size at the time of data extraction (if, for  
202 example, the information is not readily available) a SD is not calculated.  
203 Sample sizes can, however, be entered at a later time (see next section). For  
204 mean/error plots, the user clicks on an error bar followed by the mean. Error  
205 bars above or below the mean can be clicked, as sometimes one is clearer  
206 than the other. *metaDigitise()* assumes that the error bars are symmetrical.  
207 Points are displayed where the user has clicked, with the error in a different  
208 colour to the mean (Figure 1A). The user also enters the type of error used in  
209 the figure: SD, standard error (SE) or 95% confidence intervals (CI95). For  
210 box plots, the user clicks on the maximum, upper quartile, median, lower  
211 quartile and minimum. For both plot types, the user can add, edit or remove  
212 groups while digitising for when finished. Three functions, *error to sd()*, *rqm to*  
213 *mean()* and *rqm to sd()*, that convert different error types to SD, box plot data  
214 to mean and box plot data SD, respectively, are also available in the package  
215 (see supplements for further details of these conversions).

216

### 217 *3.2 Scatter plots*

218 Users can extract points from multiple groups from scatter plots.  
219 Different groups are plotted in different colours and shapes to enable them to  
220 be distinguished, with a legend at the bottom of the figure (Figure 1D). Mean,  
221 SD and sample size are calculated from the clicked points, for each group.  
222 Data points may overlap with each other making it impossible to know  
223 whether points have been missed. This may result in the sample size of  
224 digitised groups conflicting with what is reported in the paper. However, users  
225 also have the option to input known sample sizes directly, if required.  
226 Nonetheless, it is important to recognise the impact that overlapping points  
227 can have on descriptive statistics, and in particular on sampling variance.

228

### 229 *3.3 Histograms*

230 The user clicks on the top corners of each bar, which are drawn in  
231 alternating colours (Figure 1C). Bars are numbered to allow the the user to  
232 edit them. As with scatter plots, if the sample size from the extracted data  
233 does not match a known sample size, the user can enter an alternate sample  
234 size. The formulas for calculation of mean, SD and sample size are provided

235 in the supplement.

236

## 237 **4 Importing and Editing Previously Digitised data**

238 *metaDigitise* is also able to re-import, edit and re-plot previously digitised  
239 figures. When running *metaDigitise()*, the user can choose to 'Import existing  
240 data', which returns previously digitised data, from a single figure or all  
241 figures. Alternately, the *getExtracted()* function returns the data from previous  
242 digitisations, but without user interaction, allowing easier integration into larger  
243 scripts. 'Edit existing data' allows the user to re-plot or edit information for  
244 digitisations that have previously be done. Re-plotting digitisations with all  
245 metadata is an important reproducibility feature, as it allows users to see  
246 exactly what information has been extracted, as well as making it easy to spot  
247 and data extraction errors.

248

### 249 *4.1 Adding Sample Sizes to Previous Digitisations*

250 In many cases sample sizes may not be readily available when digitising  
251 figures. This information does not need to be added at the time of digitisation.  
252 To expedite finding and adding these sample sizes at a later point,  
253 *metaDigitise()* has a specific edit option that allows users to enter previously  
254 omitted sample sizes. This first identifies missing sample sizes in the digitised  
255 output, re-plots the relevant figures and prompts the user to enter the sample  
256 sizes for the relevant groups in the figure.

257

## 258 **5 Software Validation**

259 To evaluate the consistency of digitisation with *metaDigitise* between  
260 users, fourteen people digitized sets of 14 identical images created from a  
261 simulated dataset (see supplements). We found no evidence for any inter-  
262 observer variability in digitisations for the mean (ICC = 0, 95% CI = 0 to 0.029,  
263  $p > 0.999$ ), SD (ICC = 0, 95% CI = 0 to 0.033,  $p > 0.999$ ) or correlation  
264 coefficient (ICC = 0.053, 95% CI = 0 to 0.296,  $p = 0.377$ ). There was little bias  
265 between digitised and true values, on average 1.63% (mean = 0.02%, SD =  
266 4.9%,  $R = -0.03\%$ ) and there were small absolute differences between  
267 digitised and true values, on average 2.18% (mean = 0.40%, SD = 5.81%,  $R$   
268 = 0.33%) across all three descriptive statistics. SD estimates from digitisations



269 are clearly most error prone. The mean absolute differences for each plot type  
270 clearly show that this effect is driven by extraction from box plots and  
271 histograms (% difference; box plot: 15.81, histogram: 5.21, mean/error: 1.50,  
272 scatter plot: 0.43). SD estimation from box plot descriptive statistics is known  
273 to be more error prone, especially at small sample sizes (Wan et al., 2014).

274 We also used simulated data to test the accuracy of digitisations with  
275 respect to known values (see supplements). *metaDigitise* was very accurate  
276 at matching clicked points to their true values essentially being perfectly  
277 correlated with the true simulated data for both the x-variable (Pearson's  
278 correlation;  $R > 0.999$ ,  $t = 2137.4$ ,  $df = 78$ ,  $p < 0.001$ ) and y-variable ( $r >$   
279  $0.999$ ,  $t = 1897.8$ ,  $df = 78$ ,  $p < 0.001$ ) in scatterplots.

280

## 281 **6 Limitations**

282 Although *metaDigitise* is very flexible and provides functionality not seen  
283 in any other package, there are some functions that it does not perform (see  
284 Table 1). Notably *metaDigitise* lacks automated point detection. However,  
285 from our experience, manual digitising is more reliable and often equally as  
286 fast. Given the variation in image quality, calibration for automatic point  
287 detection needs to be done for each figure individually. Additionally, auto-  
288 detection often misses points that then need to be manually added. Based on  
289 tests of *metaDigitise* (see above), figures can be extracted in around 1-2  
290 minutes, including the entry of metadata. As a result, we do not believe that  
291 current automated point detection techniques provide substantial benefits in  
292 terms of time or accuracy. Indeed, in a recent project developing automated  
293 point extraction techniques, only 15/136 (11%) of studies screened contained  
294 figures suitable for the presented method, and in only 12/27 (44%) of the  
295 resulting figures was the data correctly extracted (Hartgerink & Murray-Rust,  
296 2017).

297 *metaDigitise* also (currently) lacks the ability to zoom in on figures.  
298 Zooming may enable users to gain greater accuracy when clicking on points.  
299 However, from our own experience (see results above), with a reasonably  
300 sized screen accuracy is already high, and so relatively little gain is to be had  
301 from zooming in on points.

302 In contrast to some other packages *metaDigitise* does not extract lines

303 from figures. Although line extraction is not generally necessary in  
304 comparative and meta-analysis, outside of these fields researchers may need  
305 to extract parameters of a line from a figure. Should a user like to extract lines  
306 with metaDigitise, we would recommend extracting data as a scatter plot, and  
307 clicking along the line in question. A model can then be fitted to these points  
308 (accessed by choosing to return calibrated rather than summary data) to  
309 estimate the parameters needed.

310

## 311 **7 Conclusions**

312 Increasing the reproducibility of figure extraction for meta-analysis and  
313 making this laborious process more streamlined, flexible and integrated with  
314 existing statistical software will go a long way in facilitating the production of  
315 high quality meta-analytic studies that can be updated in the future. We  
316 believe that *metaDigitise* will improve this research synthesis pipeline, and will  
317 hopefully become an integral package that can be added to the meta-analysts  
318 toolkit.

319

## 320 **Acknowledgments**

321 We thank the I-DEEL group and colleagues at UNSW for for testing,  
322 providing feedback and digitising including: Rose O'Dea, Fonti Kar,  
323 Malgorzata Lagisz, Julia Riley, Diego Barneche, Erin Macartney, Ivan Beltran,  
324 Gihan Samarasinghe, Dax Kellie, Jonathan Noble, Yian Noble, Elena Noble  
325 and Alison Pick, as well as three anonymous reviewers for their comments on  
326 the manuscript. J.L.P. was supported by a Swiss National Science Foundation  
327 Early Mobility grant (P2ZHP3\_164962), D.W.A.N. was supported by an  
328 Australian Research Council Discovery Early Career Research Award  
329 (DE150101774) and UNSW Vice Chancellors Fellowship and S.N. an  
330 Australian Research Council Future Fellowship (FT130100268).

331

## 332 **Data Accessibility**

333 Data and code for the software validation are available at  
334 <http://doi.org/10.5281/zenodo.1311681>.

335

## 336 **Author Contributions**

337 J.L.P. and D.W.A.N. conceived the study and J.L.P., S.N. and D.W.A.N.  
338 developed the idea. J.L.P. and D.W.A.N. developed the R-package. J.L.P.  
339 and D.W.A.N. wrote the first draft of the paper and J.L.P., S.N. and D.W.A.N.  
340 contributed substantially to subsequent revisions of the manuscript and gave  
341 final approval for publication.

342

## 343 **References**

- 344 Arizona-Software (2008) *GraphClick Software, Version 3.0*.
- 345 Bormann, I. (2012) *Digitizelt Software, Version 2.0*. Braunschweig, Germany.
- 346 Gurevitch, J., Koricheva, J., Nakagawa, S. & Stewart, G. (2018) Meta-analysis  
347 and the science of research synthesis. *Nature*, **555**, 175–182.
- 348 Hadfield, J.D. (2010) MCMC methods for multi-response generalized linear  
349 mixed models: The MCMCglmm R package. *Journal of Statistical Software*,  
350 **33**, 1–22.
- 351 Hartgerink, C. & Murray-Rust, P. (2017) Extracting data from vector figures in  
352 scholarly articles. *ArXiv e-prints*.
- 353 Koricheva, J. & Gurevitch, J. (2014) Uses and misuses of meta-analysis in  
354 plant ecology. *Journal of Ecology*, **102**, 828–844.
- 355 Lajeunesse, M.J. (2016) Facilitating systematic reviews, data extraction, and  
356 meta-analysis with the metagear package for R. *Methods in Ecology and*  
357 *Evolution*, **7**, 323–330.
- 358 Noble, D.W., Stenhouse, V. & Schwanz, L.E. (2018) Developmental  
359 temperatures and phenotypic plasticity in reptiles: a systematic review and  
360 meta-analysis. *Biological Reviews*, **93**, 72–97.
- 361 Parker, T.H., Forstmeier, W., Koricheva, J., Fidler, F., Hadfield, J., En Chee,  
362 Y., Kelly, C.D., Gurevitch, J. & Nakagawa, S. (2016) Transparency in  
363 Ecology and Evolution: Real Problems, Real Solutions. *Trends in Ecology*  
364 *and Evolution*, **31**, 711–719.
- 365 Peng, R.D. (2011) Reproducible research in computational science. *Science*,  
366 **334**, 1226.
- 367 Peng, R.D., Dominici, F. & Zeger, S.L. (2006) Reproducible epidemiologic  
368 research. *American Journal of Epidemiology*, **163**, 783–789.
- 369 Poisot, T. (2011) The digitize package: extracting numerical data from  
370 scatterplots. *The R Journal*, **3**, 25–26.

371 Rohatgi, A. (2017) *WebPlotDigitizer Software, Version 4.0*. Austin, Texas,  
372 USA.

373 Tummers, B. (2006) *DataThief Software, Version 3.0*.

374 Viechtbauer, W. (2010) Conducting Meta-Analyses in *R* with the metafor  
375 Package. *Journal of Statistical Software*, **36**, 1–48.

376 Wan, X., Wang, W., Liu, J. & Tong, T. (2014) Estimating the sample mean  
377 and standard deviation from the sample size, median, range and/or  
378 interquartile range. *BMC Medical Research Methodology*, **14**, 135.

379

## 380 Figure Legends

381 **Figure 1:** Functionality of *metaDigitise*. Using the iris dataset in R, digitisation  
382 of different plot types, A) mean/error plot, B) box plot, C) histogram and D)  
383 scatter plot, is shown in *metaDigitise* (left) compared with other common  
384 softwares (right). A) and B) are plotted with the whole dataset, C) is just the  
385 data for the species *setosa* and D) a subset from all three species. Notable  
386 functions of *metaDigitise* are listed in the centre. Other software also perform  
387 points 3 and 4 (see Table 1), although these functions are more developed in  
388 *metaDigitise*. As shown on the left-hand side of the figure, *metaDigitise* clearly  
389 displays the stages of the digitisation to aid the transparency of the process,  
390 and returns concatenated summary data for all images.

391

## Tables

Table 1: Comparison of functionality between different digitisation software.

Function	metaDigitise	GraphClick <sup>1</sup>	DataThief <sup>2</sup>	DigitizeIt <sup>3</sup>	WebPlotDigitizer <sup>4</sup>	metagear <sup>5</sup>	digitize <sup>6</sup>
Scatterplots	✓	✓	✓	✓	✓	✓ <sup>7</sup>	✓
Mean/error plots	✓	✓	✓	x	x	✓ <sup>7</sup>	x
Boxplots	✓	x	x	x	x	x	x
Histograms	✓	x	x	x	✓ <sup>7</sup>	x	x
Entry of metadata	✓	x	x	x	x	x	x
Grouped Data	✓	✓	x	✓	✓	x	x
Reproducible <sup>8</sup>	✓	✓	✓	x	✓	✓	x
Summarising data	✓	x	x	x	x	x	x
Multiple image processing	✓	x	x	x	x	x	x
Automated point detection	x	✓	x	✓	✓	✓	x
Line extraction	x	✓	✓	✓	✓	x	x
Zoom	x	✓	✓	✓	✓	x	x
Graph rotation <sup>9</sup>	✓	✓	✓	✓	✓	x	x
Log axis	✓	✓	✓	✓	✓	x	x

Dates	x	x	✓	x	✓	x	x
Asymmetric error bars	x	x	✓	x	x	x	x
Freeware	✓ <sup>10</sup>	✓ <sup>11</sup>	✓ <sup>11</sup>	x <sup>11</sup>	✓ <sup>11</sup>	✓ <sup>10</sup>	✓ <sup>10</sup>

<sup>1</sup> Arizona-Software (2008) <sup>2</sup> Tummers (2006) <sup>3</sup> Bormann (2012) <sup>4</sup> Rohatgi (2017) <sup>5</sup> Lajeunesse (2016) <sup>6</sup> Poisot (2011)

<sup>7</sup> Only automated, no manual extraction.

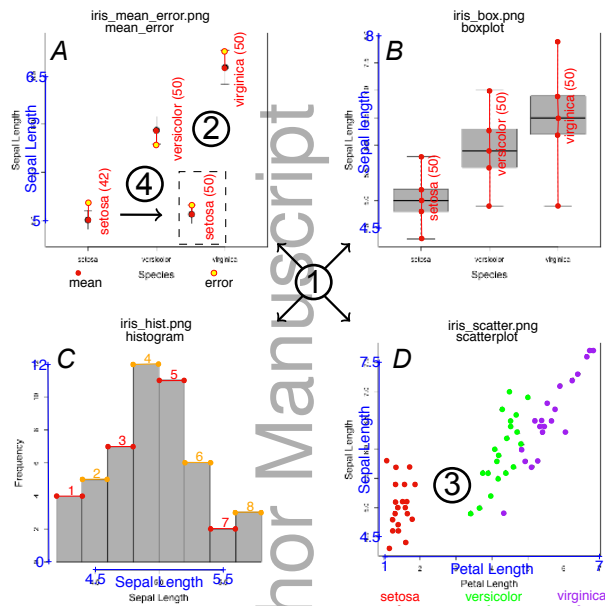
<sup>8</sup> Allows saving, re-plotting and editing of data extraction.

<sup>9</sup> Or handles rotated graphs.

<sup>10</sup> R package.

<sup>11</sup> Standalone software.

Author Manuscript



## DATA OUTPUT

File name	Variable	Method	mean	sd	n	p	plot type
iris_box.png	Sepal.Length	setosa	1.01	0.317	50	NA	boxplot
iris_box.png	Sepal.Length	versicolor	1.85	0.189	50	NA	boxplot
iris_box.png	Sepal.Length	virginica	6.19	0.603	50	NA	boxplot
iris_hist.png	Sepal.Length	setosa	4.85	0.364	50	NA	histogram
iris_hist.png	Sepal.Length	versicolor	5.01	0.550	50	NA	histogram
iris_hist.png	Sepal.Length	virginica	1.02	1.035	50	NA	histogram
iris_scatter.png	Petal.Length	setosa	1.1	0.175	25	0.159	scatterplot
iris_scatter.png	Petal.Length	versicolor	5.08	0.427	25	0.109	scatterplot
iris_scatter.png	Petal.Length	virginica	4.29	0.425	25	0.786	scatterplot
iris_mean_error.png	Sepal.Length	setosa	1.07	0.303	25	0.786	mean_error
iris_mean_error.png	Sepal.Length	versicolor	1.87	0.193	25	0.786	mean_error
iris_mean_error.png	Sepal.Length	virginica	6.86	0.688	25	0.803	mean_error

## ① Different plot types

Capable of handling A) mean error plots, B) boxplots, C) histograms and D) scatterplots

## ② Entry of Metadata

Enter sample sizes variable and group names while digitising that are displayed on plot

## ③ Grouped Data

Enter as many groups as needed to capture descriptive statistics for sub-samples of data

## ④ Digitise, edit or replot digitisations

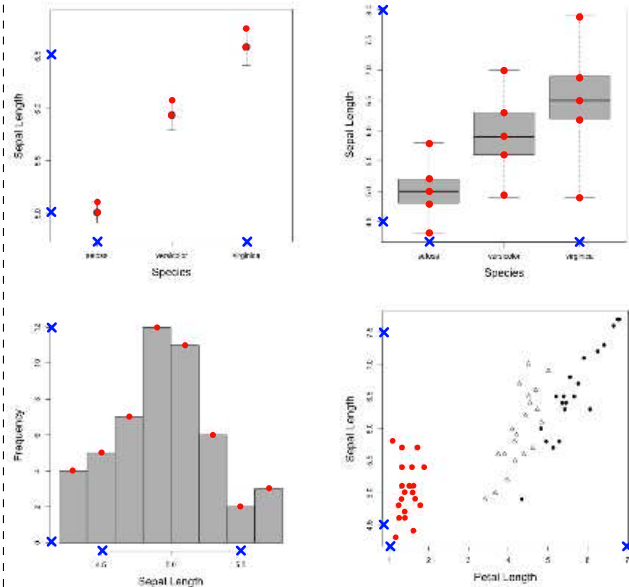
Simple user interface to guide user. Can digitise new images, edit digitisations or easily replot previous digitisations and metadata by cycling through images or choosing specific images

## ⑤ Summarising data

Get descriptive statistics automatically calculated for all plot types or use raw x,y data, if desired

## ⑥ Multiple image processing

Process as many images at once as needed and of varying types efficiently and quickly. New plots automatically plotted for digitisation



## DATA OUTPUT

Mean-Error		Box-plot		Histogram		Scatter-plot	
x	y	x	y	x	y	x	y
1.00	0.19	1.00	0.30	4.20	4.00	1.00	0.80
1.00	0.30	1.00	0.21	4.50	4.99	1.01	0.70
2.00	0.08	1.00	0.00	4.70	7.01	1.70	0.70
2.00	0.28	1.00	0.30	1.00	19.01	1.30	0.10
3.00	0.75	1.00	4.31	5.10	11.02	1.58	0.41
2.00	0.59	2.00	7.00	5.30	5.98	1.80	0.40
2.00	0.30	2.00	0.30	5.40	1.08	1.32	0.11
2.00	0.91	2.00	0.91	5.70	3.00	1.51	0.10
2.00	0.61					1.56	0.10
2.00	1.21					1.37	0.00