

1 Reproducible, flexible and high throughput data extraction from primary 2 literature: The **metaDigitise R** package

3 Joel L. Pick^{1,*}, Shinichi Nakagawa¹, Daniel W.A. Noble¹

4 ¹ Ecology and Evolution Research Centre, School of Biological, Earth and
5 Environmental Sciences, University of New South Wales, Kensington, NSW 2052,
6 Sydney, AUSTRALIA

7 *Corresponding Author: joel.l.pick@gmail.com

8 **Abstract**

9 Research synthesis, especially in the form of meta-analysis, requires data extraction
10 from primary studies. Meta-analysis synthesizes effect sizes, often calculated from
11 summary statistics of studies. However, exact values of such statistics are commonly
12 hidden in figures. The R package **metaDigitise** extracts descriptive statistics such as
13 means, standard deviations and, if applicable, correlations from the four types of plots:
14 1) mean and error plots (e.g. bar graphs with standard errors), 2) box plots, 3) scatter
15 plots and 4) histograms. The package interactively guides the user through data
16 extraction process. Notably, it enables a large-scale extraction using image files, letting
17 the user stop processing, edit and add to the resulting data frame at any point. Further,
18 it facilitates reproducible data extraction from plots with little inter-observer bias, thus,
19 allowing a group of people to participate the extraction of data collaboratively.

20 Keywords: meta-analysis, comparative analysis, data extraction, R, reproducibility,
21 figures, images, summary statistics

22 Introduction

23 In many different contexts, researchers need to make use of data presented in primary
24 literature. Most notably, this includes meta-analysis, which is becoming increasingly
25 common in many research fields. Meta-analysis uses effect size estimates and their
26 sampling variance, taken from many studies, to understand whether particular effects
27 are common across studies and to explain variation among these effects (Glass, 1976;
28 Borenstein et al., 2009; Koricheva, Gurevitch & Mengersen, 2013; Nakagawa et al.,
29 2017). Meta-analysis therefore relies foremost on data extracted from primary
30 literature, and more specifically, descriptive statistics (e.g., means, standard deviations,
31 correlation coefficients) that have been reported in the text or tables of research papers.
32 Descriptive statistics are also, however, frequently presented in figures and so need to be
33 manually extracted using digitising programs. While inferential statistics (e.g., t - and
34 F -statistics) are often presented along side descriptive statistics and can be used to
35 derive effect sizes, descriptive statistics are much more appropriate to use because
36 sources of non-independence in experimental designs can be dealt with more easily
37 (Noble et al., 2017). Although there are several existing tools to perform tasks like this
38 (e.g. **DataThief** (Tummers, 2006), **GraphClick** (Arizona-Software, 2008), **WebPlotDigitizer**
39 (Rohatgi, 2017)), these tools are not designed specifically for meta-analysis for three
40 main reasons.

41 First, they typically only provide the user with calibrated x,y coordinates from
42 imported figures, and do not differentiate between common plot types that are used to
43 present data. This means that a large amount of downstream data manipulation is
44 subsequently required, that is different across plots types. For example, data are
45 frequently presented in mean and error plots (Figure 1A), for which the user wants a
46 mean and error estimate for each group presented in the figure. With existing
47 programs, x,y coordinates of means and errors are returned, to which the user must
48 manually discern between mean and error coordinates and assign points to groups. The

error then needs to be calculated as the deviation from the mean, and then transformed to a standard deviation, depending on the type of error presented.

Second, digitising programs do not easily allow the integration of metadata at the time of data extraction, such as experimental group or variable names, and sample sizes. This makes the downstream calculations more laborious, as the information has to be added later, in most cases using different software.

Finally, existing programs do not import a set of images and allow the user to systematically work through them. Instead they require the user to manually import images one by one, and export data into individual files, that need to be imported and edited using different software. In essence, existing software does not provide an optimized research pipeline to facilitate data extraction, editing and reproducibility.

These are major issues because extracting from figures can be an incredibly time-consuming process. Furthermore, although meta-analysis is an important tool in consolidating the data from multiple studies, many of the processes involved in data extraction are opaque and difficult to reproduce, making extending studies problematic. Having a tool that facilitates reproducibility in meta-analyses will increase transparency and go a long way to resolving the reproducibility crises we are seeing in many fields (Peng, Dominici & Zeger, 2006; Peng, 2011; Sandve et al., 2013; Parker et al., 2016; Ihle et al., 2017).

Here, we present an interactive R package, **metaDigitise**, which is designed for large scale data extraction from figures, specifically catering to the the needs of meta-analysts. To this end, we provide tools specific to data extraction from common plot types (mean and error plots, box plots, scatter plots and histograms, see Figure 1). **metaDigitise** operates within the R environment making data extraction, analysis and export more streamlined. It also provides users with options to conduct the necessary calculations on processed data immediately after extraction so that comparable

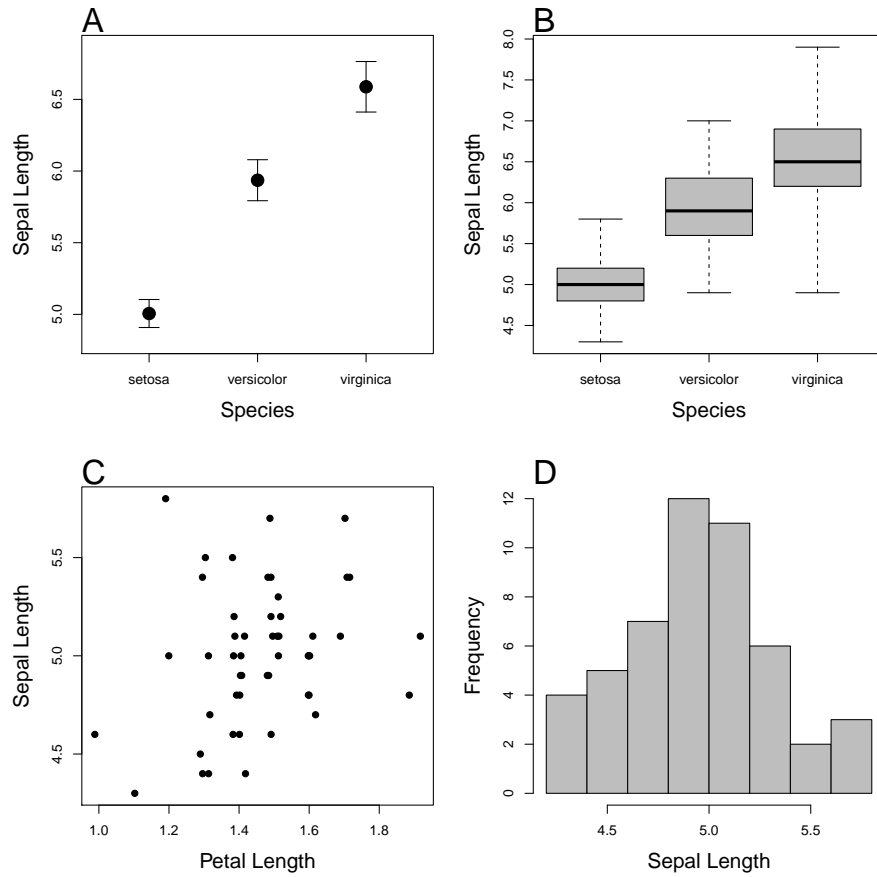


Figure 1: Four plot types that **metaDigitise** is designed to extract data from: A) mean and error plot, B) box plot, C) scatter plot and D) histogram. Data is taken from the iris dataset in R. A and B are plotted with the whole dataset, C and D are just the data for the species setosa.

76 summary statistics can be obtained quickly. **metaDigitise** condenses summary data
77 extracted from multiple figures into a single data frame which can be easily
78 exported. Processed data can also be easily extracted and analysed in any way the user
79 desires in downstream analysis within R. Conveniently, when needing to process many
80 figures at different times **metaDigitise** will only import figures not already completed
81 within a directory. This makes it easy to add new figures at any time. **metaDigitise** has
82 also been built for reproducibility in mind. It has functions that allow users to redraw
83 their digitisations on figures, make corrections and access the raw calibration data which
84 is written automatically for each figure that is digitised into a special folder within the
85 directory. This makes sharing figure digitisation and reproducing the work of others

86 simple and easy, and allows meta-analysts to update meta-analyses more easily.

87 Directory Structure, Image Processing and 88 Reproducibility

89 The **metaDigitise** package is designed to be flexible, yet simple to use. There is one
90 main function in the package, `metaDigitise()`, which interactively takes the user
91 through the process of extracting data from figures. `metaDigitise()` was created with
92 the idea that the user would likely have multiple images to extract from. It therefore
93 operates in the same way whether the user has one or multiple images.
94 `metaDigitise()` is designed to work on a directory containing images of figures copied
95 from primary literature, in .png, .jpg, .tiff, .pdf format. This directory is specified to
96 `metaDigitise()` through the `dir` argument. The user is free to set their own broad
97 directory structure (e.g. one directory for all images or one directory for each paper
98 extracted from). We would recommend having all files for one project in a single
99 directory with an informative and unambiguous naming scheme for images to make it
100 easy to identify the paper and figure the data come from. This cuts out the need to
101 change directories constantly. For example the directory structure could look like:

```
* Main project directory
  + FiguresToExtract/
    + Paper1_Figure1_trait1.png
    + Paper1_Figure2_trait2.png
    + Paper1_Figure3_trait3.png
    + Paper2_Figure1_trait1.png
    + Paper2_Figure2_trait2.png
    + Paper2_Figure3_trait3.png
```

102 It is important for the user to think about their directory structure early on in this
103 process (also more generally in the context of their entire project), especially if they
104 plan to share the extractions with collaborators or when publishing the project.

105 When `metaDigitise()` is run, it recognizes all the images in a directory and
106 automatically imports them one by one, allowing the user to click and enter relevant
107 information about a figure as they go. This expedites digitising figures by preventing
108 users from having to constantly change directories and / or open new images. The data
109 from a completed image is automatically saved as a `metaDigitise` object in an `.RDS`
110 file to a `caldat` directory that is created within the parent directory when first
111 executing the `metaDigitise()` function. These files enable re-plotting and editing of
112 images at a later point (see below).

113 A particularly powerful and flexible aspect of `metaDigitise()` is its ability to identify
114 images that have been previously digitised and only import images that have not been
115 digitised in subsequent calls of the function. This means that all figures do not need to
116 be extracted at one time and that new figures can be added as the project develops.
117 After each image is extracted, the user is asked whether they wish to continue or quit
118 the extraction process. Upon rerunning `metaDigitise()`, previously digitised figures
119 are simply ignored during processing, but their data is re-integrated within the final
120 output after new files are completed automatically.

121 After completing all images, or upon quitting, the processed data (in a form specified
122 by the user) is then returned. From all plot types, `metaDigitise()` summarises the
123 data from a figure as a mean, standard deviation and sample size, for each identified
124 group within the plot (should multiple groups exist). These are the descriptive statistics
125 needed to create many of the relevant effect sizes and sampling error for a
126 meta-analysis. In the case of scatter plots, `metaDigitise()` also returns the correlation
127 coefficient between the points within each identified group.

128 **Diverse Plot Types**

129 **metaDigitise** recognises four main types of plot; Mean and error plots, box plots,
130 scatter plots and histograms, shown in Figure 1. Each of these can be processed
131 together and integrated into a single output. Alternatively, users can keep like figures
132 together and process them separately.

133 In order to correctly extract data from figures **metaDigitise()** always requires the user
134 to calibrate the axes in the figure. To do this, the user is required to click on two known
135 points on the axis in question, and then enter the value of those points in the figure.
136 Using this information, **metaDigitise()** then calculates the value of any clicked points
137 in terms of the figure axes. In the case of mean and error plots and box plots, it
138 calibrates only the y-axis (assuming the x-axis is redundant). For scatter plots and
139 histograms both axes are calibrated.

140 **Mean and error plots**

141 **metaDigitise()** prompts the user to enter group names and allows the user to enter
142 sample sizes (n), which are used in downstream processing. The user is then prompted
143 to click on an error bar followed by the mean. Error bars above or below the mean can
144 be clicked - sometimes one is clearer than the other. **metaDigitise()** assumes that the
145 error bars are symmetrical. Where the user has clicked the error is displayed in a
146 different colour to the mean (Figure 2A). The user can subsequently add more groups,
147 edit groups or remove groups. Finally the user is asked what type of error was used in
148 the figure: standard deviation (SD, σ), standard error (SE) or 95% confidence intervals
149 (CI95). Standard deviation is calculated from standard error as

$$\sigma = SE\sqrt{n} \tag{1}$$

150 and from 95% confidence intervals as

$$\sigma = \frac{CI}{1.96} \sqrt{n} \quad (2)$$

151 If the user does not enter a sample size at the time of data extraction (if, for example,
152 the information is not readily available) the SD is not calculated. This can be entered
153 at a later time, however (see below). A function, `error_to_sd()`, that converts the
154 different error types to SD is also available in the package.

155 **Box plots**

156 As with mean and error plots, `metaDigitise()` prompts the user to enter group names
157 and allows the user to enter sample sizes (n), which are used in downstream processing.
158 The user is then prompted to click on the maximum (b), upper quartile (q_3), median
159 (m), lower quartile (q_1) and minimum (a). `metaDigitise()` will check that the
160 maximum is greater than the minimum, and return a warning if that is not the case.
161 The user can subsequently add, edit or remove groups. From the extracted data, the
162 mean (μ) and SD are calculated as

$$\mu = \frac{(n+3)(a+b) + 2(n-1)(q_1 + m + q_3)}{8n} \quad (3)$$

163 following Bland (2015) and

$$\sigma = \frac{b-a}{4\Phi^{-1}\left(\frac{n-0.375}{n+0.25}\right)} + \frac{q_3-q_1}{4\Phi^{-1}\left(\frac{0.75n-0.125}{n+0.25}\right)} \quad (4)$$

164 where $\Phi^{-1}(z)$ is the upper z th percentile of the standard normal distribution, following
165 Wan et al. (2014). As with mean and error plots, if the user does not enter a sample
166 size at the time of data extraction the SD is not calculated. Two functions,
167 `rqm_to_mean()` and `rqm_to_sd()`, that convert box plot data to mean and SD

168 respectively are also available in the package.

169 Scatter plots

170 `metaDigitise()` prompts the user to enter groups names and then to click on points.
171 Points added by mistake can be deleted. The user can subsequently add groups, edit
172 groups (add or remove points) or delete groups. Different groups are plotted in different
173 colours and shapes, with a legend at the bottom of the figure (Figure 2C). Mean, SD
174 and sample size are calculated from the clicked points, for each group. Where the
175 sample size from the clicked points does not match a known sample size (e.g. if there
176 are overlaid points), the user can enter an alternate sample size.

177 Histograms

178 `metaDigitise()` prompts the user to click on the top corners of each bar. Bars can
179 subsequently be deleted. For each bar a midpoint (m ; mean x coordinates) and a
180 frequency (f ; mean y coordinates, rounded to the nearest integer) is calculated. The
181 sample size, mean and SD are calculated as:

$$n = \sum_{i=1}^n f_i \quad (5)$$

$$\mu = \frac{\sum_{i=1}^n m_i f_i}{n} \quad (6)$$

$$\sigma = \sqrt{\frac{\sum_{i=1}^n (m_i f_i - \mu f_i)^2}{n - 1}} \quad (7)$$

182 As with the scatterplots, if the sample size from the extracted data does not match a
183 known sample size, the user can enter an alternate sample size.

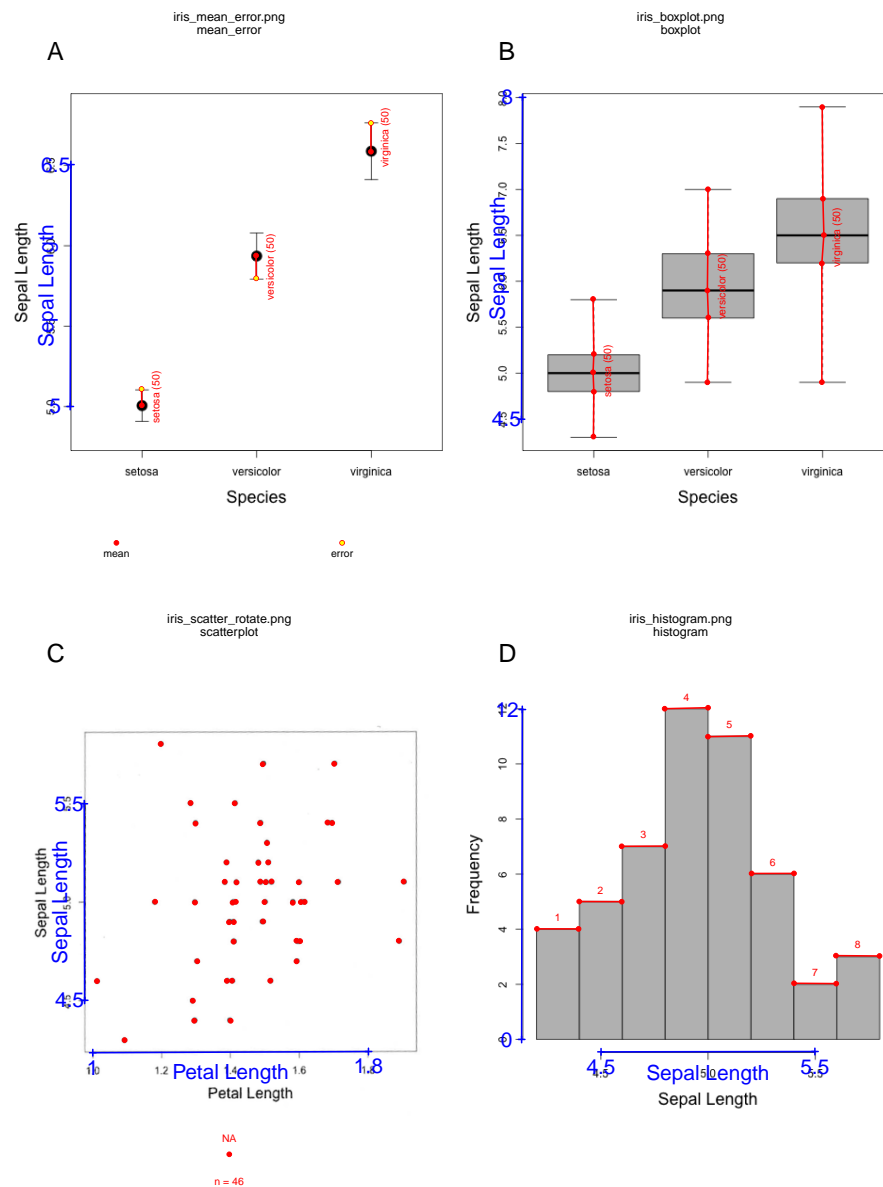


Figure 2: Demonstration of data extraction from different plot types

184 Extracting Data From Plots

185 We will now demonstrate how `metaDigitise()` works using figures generated from the
 186 well known iris data set. Users can install the **metaDigitise** package from GitHub as
 187 follows:

```
R> install.packages("devtools")
```

```
R> devtools::install_github("danielnoble/metaDigitise")
```

```
R> library(metaDigitise)
```

188 Assume that the user would like to extract descriptive statistics from studies measuring
189 sepal length or width in iris species for a fictitious project. There are a few studies that
190 only present these data in figures. As the user reads papers found from a systematic
191 search, they add figures with relevant data to a "FiguresToExtract" folder as
192 follows

```
*FiguresToExtract/  
+ 001_Anderson_1935_Fig1.png
```

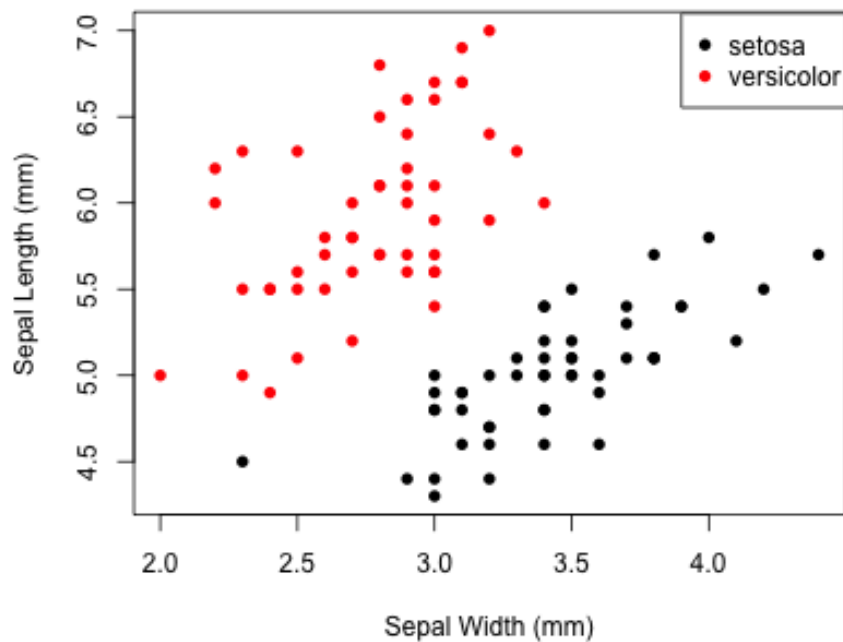


Figure 3: Example scatterplot (001_Anderson_1935_Fig1.png) of sepal length and width for two species of iris (setosa and versicolor)

193 Here, the naming of the files placed in the folder will contain the paper number, first
194 author and the figure number to keep data uniquely associated with figures. At first
195 there is one figure in the folder, shown in Figure 3. Running `metaDigitise()` brings up
196 a series of prompts for the user using a main menu that provides access to a number of

197 its features ("..." here represents the user's path to the project directory):

```
R> digitised_data <- metaDigitise("../FiguresToExtract", summary = TRUE)
```

```
Do you want to...
```

```
1: Process new images
```

```
2: Import existing data
```

```
3: Edit existing data
```

```
Selection:
```

198 The user simply enters in the numeric value that corresponds to what they would like to
199 do. In this case they want to "Process new images". The user is then asked whether
200 there are different types of plot(s) in the folder. This question is most relevant when
201 there are lots of different figures in the folder because it will then ask the user for the
202 type of figure as they are cycled through.

```
Are all plot types Different or the Same? (d/s)
```

203 metaDigitise() then asks the user whether the figure needs to be rotated or flipped.
204 This can be needed when box plots and mean and error plots are not orientated
205 correctly. In some cases, older papers can give slightly off angled images which can be
206 corrected by rotating. So, in this prompt the user has three options: **f** for "Flip", **r** for
207 "rotate" or **c** for "continue".

```
mean_error and boxplots should be vertically orientated
```

```

-
|
I.E. o    NOT  |-o-|
|
-

```

If they are not then chose flip to correct this.

If figures are wonky, chose rotate.

Otherwise chose continue

Flip, rotate or continue (f/r/c)

R> c

208 After this, `metaDigitise()` will ask the user to specify the plot type. Depending on the
209 figure, the user can specify that it is a figure containing the mean and error (**m**), a box
210 plot (**b**), a scatter plot (**s**) or a histogram (**h**). If the user has specified **d** instead of **s** in
211 response to the question about whether the plot types are the same or different, this
212 question will pop up for each plot, but will only be asked once if plots are all the
213 same.

Please specify the `plot_type` as either:

m: Mean and error

b: Box plot

s: Scatter plot

h: Histogram

R> s

214 After selecting the figure type a new set of prompts will come up that will ask the user
215 first what the y and x-axis variables are. This is useful as users can keep track of the
216 different variables across figures and papers. Here, the user can just add this
217 information in to the R console. Once complete, details on how to calibrate the x and
218 y-axis appear, so that the relevant statistics / data can be correctly calculated. When
219 working with a plot of mean and standard errors, the x-axis is rather useless in terms of

220 calibration so metaDigitise() just asks the user to calibrate the y-axis.

What is the y variable?

R> Sepal Length (mm)

What is the x variable?

R> Sepal Width (mm)

On the Figure, click IN ORDER:

y1, y2 , x1, x2

Step 1 ----> Click on known value on y axis - y1

|
|
|
|
y1
|-----
....

Step 3 ----> Click on known value on x axis - x1

|
|
|
|
|
|
|-----x1-----

....

221 The user can just follow the instructions on screen step-by-step (instructions above have
222 been truncated by ‘...’ to simplify), and in the order specified. Before moving on, the
223 user is forced to check whether or not the calibration has been set up correctly. If **n** is
224 chosen because something needs to be fixed then the user can re-calibrate.

What is the value of y1 ?

R> 4.5

What is the value of y2 ?

R> 7

What is the value of x1 ?

R> 2

What is the value of x2 ?

R> 4

Re-calibrate? (y/n)

R> n

225 Often, plots might contain multiple groups that the meta-analyst wants to extract from.
226 `metaDigitise()` handles this nicely by prompting the user to enter the group first,
227 followed by digitisation of this groups data. After digitising the first group, and having
228 exited, `metaDigitise()` will ask the user whether they would like to add another
229 group. Users can continually add groups (**a**), delete groups (**d**), edit groups (**e**) or finish
230 a plot and continue to the next one (**f** - if another plot exists). The number of groups
231 are not really limited and users can just keep adding in groups to accommodate the
232 different numbers that may be presented across figures (although it can get complicated
233 with too many).

If there are multiple groups, enter unique group identifiers (otherwise press enter)

Group identifier:

```
R> setosa
```

Click on points you want to add.

If you want to remove a point, or are finished with a group,

exit by clicking on red box in bottom left corner, then follow prompts

234 To finish selecting points, the user can exit by clicking on the red button that appears
235 when extracting points. The user is then asked if they want to add or delete points
236 from that group.

Add or Delete points to this group, or Continue? (a/d/c)

```
R> c
```

237 Once we are done digitising all the groups our plot will look something like Figure
238 4.

239 When completed `metaDigitise()` will write the digitised data as a `metaDigitise`
240 object to a RDS file in the `caldat` directory, such that our new directory structure is as
241 follows

```
*FiguresToExtract/  
  + caldat/  
    + 001_Anderson_1935_Fig1  
    + 001_Anderson_1935_Fig1.png
```

242 Users can access the `metaDigitise` object created (`001_Anderson_1935_Fig1`) at any
243 time using the `metaDigitise()` function. In the R console, the summarised data for the
244 digitised figure can be printed on screen or even written to a `.csv` file:

```
R> digitised_data
```

filename	group_id	variable	mean	error	error_type	n	r	sd	plot_type
----------	----------	----------	------	-------	------------	---	---	----	-----------

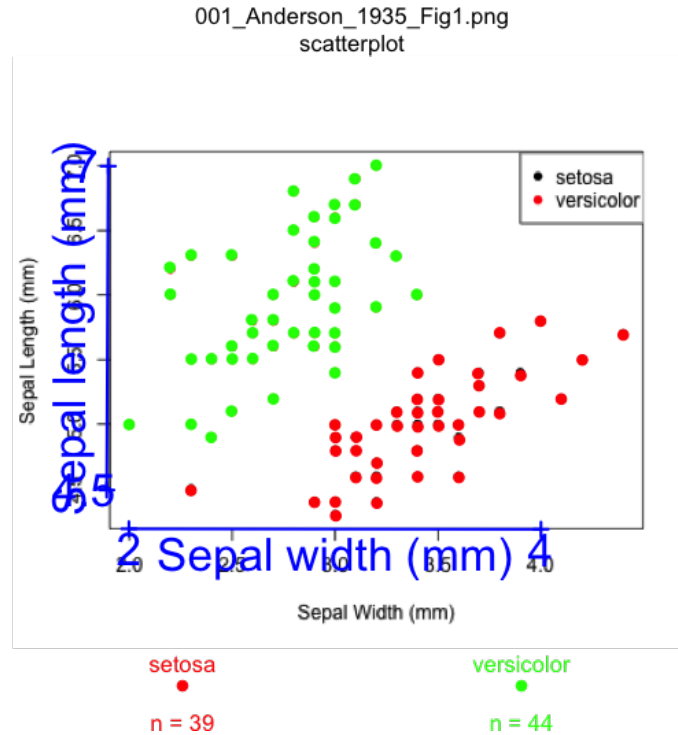


Figure 4: Digitisation of sepal length and width for two species of iris (setosa and versicolor). Names of the variables and calibration (in blue) are plotted alongside the digitised points (green = versicolor; red = setosa). The sample sizes for each group are provided on the lower part of the plot. All figures are clearly labelled at the top to remind users of the filename and plot type. This reduces errors throughout the digitisation process.

001_Anderson_1935_Fig1.png	setosa	Sepal width (mm)	3.42	0.40	sd	39	0.75	0.40	scatterplot
001_Anderson_1935_Fig1.png	setosa	Sepal length (mm)	5.00	0.38	sd	39	0.75	0.38	scatterplot
001_Anderson_1935_Fig1.png	versicolor	Sepal width (mm)	2.77	0.32	sd	44	0.52	0.32	scatterplot
001_Anderson_1935_Fig1.png	versicolor	Sepal length (mm)	5.95	0.53	sd	44	0.52	0.53	scatterplot

245 The mean for each of the two variables, along with the two species, are provided. Since
 246 this is a scatterplot, the user also gets the Person's correlation coefficient between sepal
 247 length and width for each species. These match reasonably well with the actual means
 248 of sepal length and width for each of the species in the full 'iris' dataset:

	Species	meanSL	meanSW
1	setosa	5.006	3.428
2	versicolor	5.936	2.770

249 One thing anyone with a familiarity with the iris dataset will notice is that the sample

sizes for each of these species (which are $n = 50$ each) are quite a bit lower. This is an example of some of the challenges when extracting data from scatter plots. Often data points will overlap with each other making it impossible (without having the real data) to know whether this is a problem. However, a meta-analyst will probably realise that the sample sizes here conflict with what is reported in the paper. Hence, **metaDigitise** also provides the user with options to input the sample sizes directly (see Editing section below), even for scatter plots and histograms where, strictly speaking, this should not be necessary. Nonetheless, it is important to recognise the impact that overlapping points can have on summary statistics, particularly its effects on standard deviation (SD) and standard error (SE). Here, the mean point estimates are nearly exactly the same as the true values, but the SD's are slightly over-estimated:

	Species	meanSL	meanSW
1	setosa	0.3524897	0.3790644
2	versicolor	0.5161711	0.3137983

Adding new figures

Users can add additional figures as new papers with relevant information are found. Each figure should be in its own file with unique naming, even if a single paper has multiple figures for extraction. For example, another paper on different populations (and one new species) of iris contained two additional figures where important data could be extracted. These figures can simply be named accordingly and added directly to the same extraction folder:

```
*FiguresToExtract/
+ caldat/
+ 001_Anderson_1935_Fig1
+ 001_Anderson_1935_Fig1.png
+ 002_Doe_2013_Fig1.png
```

+ 002_Doe_2013_Fig3.png

268 The user has already processed one figure (001_Anderson_1935_Fig1.png). We can tell
269 this because the caldat folder has digitised data in it (caldat/001_Anderson_1935_Fig1).
270 Now the user has two new figures that have not yet been digitised. This example will
271 nicely demonstrate how users can easily pick up from where they left off and how all
272 previous data gets re-integrated. It will also demonstrate how different plot types are
273 handled. All we have to do to begin, is again, provide the directory where all the figures
274 are located:

```
R> digitised_data <- metaDigitise("../FiguresToExtract", summary = TRUE)
```

275 The user gets the same set of prompts and simply chooses option one. This will permit
276 users to digitise new figures, and will integrate previously completed digitisations along
277 with newly digitised data together at the end of the session, or when the user decides to
278 quit. This time, 001_Anderson_1935_Fig1.png is ignored and the new plots cycle on
279 screen; first 002_Doe_2013_Fig1.png and then 002_Doe_2013_Fig3.png. Since there are a
280 few different figure types, the user answers the first question in the R console as
281 "d":

Are all plot types Different or the Same? (d/s)

```
R> d
```

**** NEW PLOT ****

mean_error and boxplots should be vertically orientated

```
-  
|  
I.E. o    NOT  |-o-|  
|  
-
```

If they are not then chose flip to correct this.

If figures are wonky, chose rotate.

Otherwise chose continue

Flip, rotate or continue (f/r/c)

R> c

Please specify the plot_type as either:

m: Mean and error

b: Box plot

s: Scatter plot

h: Histogram

R> m

282 Here, the user specifies the new plot type as *m* for 002_Doe_2013_Fig1.png because the
283 user has a plot of the mean and error of sepal length for each of the three species. The
284 user is then prompted a bit differently from our scatter plot as the x-axis is not needed
285 for calibration:

What is the y variable?

R> Sepal length

On the Figure, click IN ORDER:

y1, y2

Step 1 ----> Click on y1

```
|  
|  
|  
|  
y1  
|-----
```

Step 2 ----> Click on y2

```
|  
y2  
|  
|  
|  
|  
|-----
```

What is the value of y1 ?

R> 5

What is the value of y2 ?

R> 6.5

Re-calibrate? (y/n)

R> n

Do you know sample sizes? (y/n)

R> y

If there are multiple groups, enter unique group identifiers (otherwise press enter)

Group identifier:

```
R> setosa
```

Group sample size:

```
R> 50
```

Click on Error Bar, followed by the Mean

Add group, Edit Group, Delete group or Finish plot? (a/e/d/f)

```
R> a
```

286 Again, `metaDigitise()` will simply guide the user through digitising each of these
287 figures describing to them exactly what needs to be done. At any point if mistakes are
288 made the user can choose relevant options to edit or correct things before ending the
289 figure. This process continues for each plot so long as the user would like to continue
290 and after completing a single plot the user is always prompted as follows:

Do you want continue: 1 plots out of 2 plots remaining (y/n)

```
R> y
```

291 This continues until users have completed all non-digitised figures in the folder, at
292 which point `metaDigitise()` concatenates the new data with previously digitised data
293 in the object:

```
data
```

	filename	group_id	variable	mean	error	error_type	n	r	sd	plot_type
	001_Anderson_1935_Fig1.png	setosa	Sepal width (mm)	3.42	0.40	sd	39	0.75	0.40	scatterplot
	001_Anderson_1935_Fig1.png	setosa	Sepal length (mm)	5.00	0.38	sd	39	0.75	0.38	scatterplot
	001_Anderson_1935_Fig1.png	versicolor	Sepal width (mm)	2.77	0.32	sd	44	0.52	0.32	scatterplot
	001_Anderson_1935_Fig1.png	versicolor	Sepal length (mm)	5.95	0.53	sd	44	0.52	0.53	scatterplot
	002_Doe_2013_Fig1.png	setosa	Sepal length	5.00	0.11	se	50	NA	0.78	mean_error
	002_Doe_2013_Fig1.png	virginica	Sepal length	6.59	0.18	se	50	NA	1.26	mean_error
	002_Doe_2013_Fig1.png	versicolor	Sepal length	5.94	0.14	se	50	NA	1.01	mean_error
	003_Doe_2013_Fig3.png	catana	Sepal length	4.95	0.36	sd	50	NA	0.36	histogram

294 Re-importing, Editing and Plotting Previously

295 Digitised data

296 A particularly useful feature of **metaDigitise** is its ability to re-import, edit and re-plot
297 previously digitised figures. We can do this from the initial options from

298 `metaDigitise()`

```
R> digitised_data <- metaDigitise("../FiguresToExtract")
```

```
Do you want to...
```

```
1: Process new images
```

```
2: Import existing data
```

```
3: Edit existing data
```

```
Selection:
```

299 If the user chooses "Import existing data", they have the option of either 1) importing
300 data from all digitised images or 2) importing data from a single image that has been
301 digitised. If 2, then a list of files are provided to the user that they can select. Editing
302 existing data allows users to easily re-plot or edit information or digitisations that have
303 previously be done for any plot. This is accomplished by guiding the user through a
304 new set of options:

```
Choose how you want to edit files:
```

```
1: Cycle through images
```

```
2: Choose specific file to edit
```

```
3: Enter previously omitted sample sizes
```

```
Selection:
```

305 If the user is unsure about the name of the specific figure they need to edit or simply
306 want to just check the digitisations of figures they can choose "Cycle through images",
307 which will bring up each figure, one by one, overlaying the calibrations, group names (if

308 they exist), sample sizes (if they were entered) and the selected points. The user will
309 then be given the choice to edit individual images. Alternatively, choosing option 2, will
310 bring up a list of the completed files in the folder and the specific file can be chosen, at
311 which point it will be replotted. Either of these options will cycle through a number of
312 questions asking the user what they would like to edit:

Edit rotation? If yes, then the whole extraction will be redone (y/n)

R> n

Change plot type? If yes, then the whole extraction will be redone (y/n)

R> n

Variable entered as:

R> Sepal length

Rename Variables (y/n)

R> n

Edit calibration? (y/n)

R> n

Re-extract data (y/n)

R> y

Change group identifier? (y/n)

R> n

Add group, Delete group or Finish plot? (a/d/f)

R> d

1: setosa

2: versicolor


```
3: virginica
```

```
Selection:
```

```
R> 2
```

```
Add group, Delete group or Finish plot? (a/d/f)
```

```
R> a
```

313 A whole host of information can be edited including the rotation, plot type, the variable
314 name(s) that were provided, the calibration and even the digitisation of groups. When
315 editing the `metaDigitise` object is re-written to the caldat folder and the edits are
316 immediately integrated into the existing object once complete.

317 **Additional Features**

318 **Figure Rotation and Adjustment**

319 Figures may have been extracted from old publications, for example from scanned
320 images, and so are not perfectly orientated on the image. This will make the calibration
321 of the points in the figure from the image problematic. `metaDigitise()` allows users to
322 rotate the image. By clicking two points on the x-axis, `metaDigitise` calculates the angle
323 needed to rotate the image so the x-axis is horizontal, and rotates it. (Figure
324 5A,B)

325 Furthermore, some figures, including mean and error, boxplots or histograms, may be
326 presented with horizontal bars. `metaDigitise()` assumes that the bars are vertical, but
327 allows the user to flip the image so that the bars are vertical if provided horizontally
328 (Figure 5C,D).

329 Obtaining Processed Data

330 While `metaDigitise()` provides users with the summary statistics by default, for all
 331 plot types, in many cases the user may actually be interested in obtaining the processed
 332 digitised data from scatter plots (i.e. calibrated points). This is very easy to do by
 333 changing the default `summary` argument from `TRUE` to `FALSE` in `metaDigitise()`.

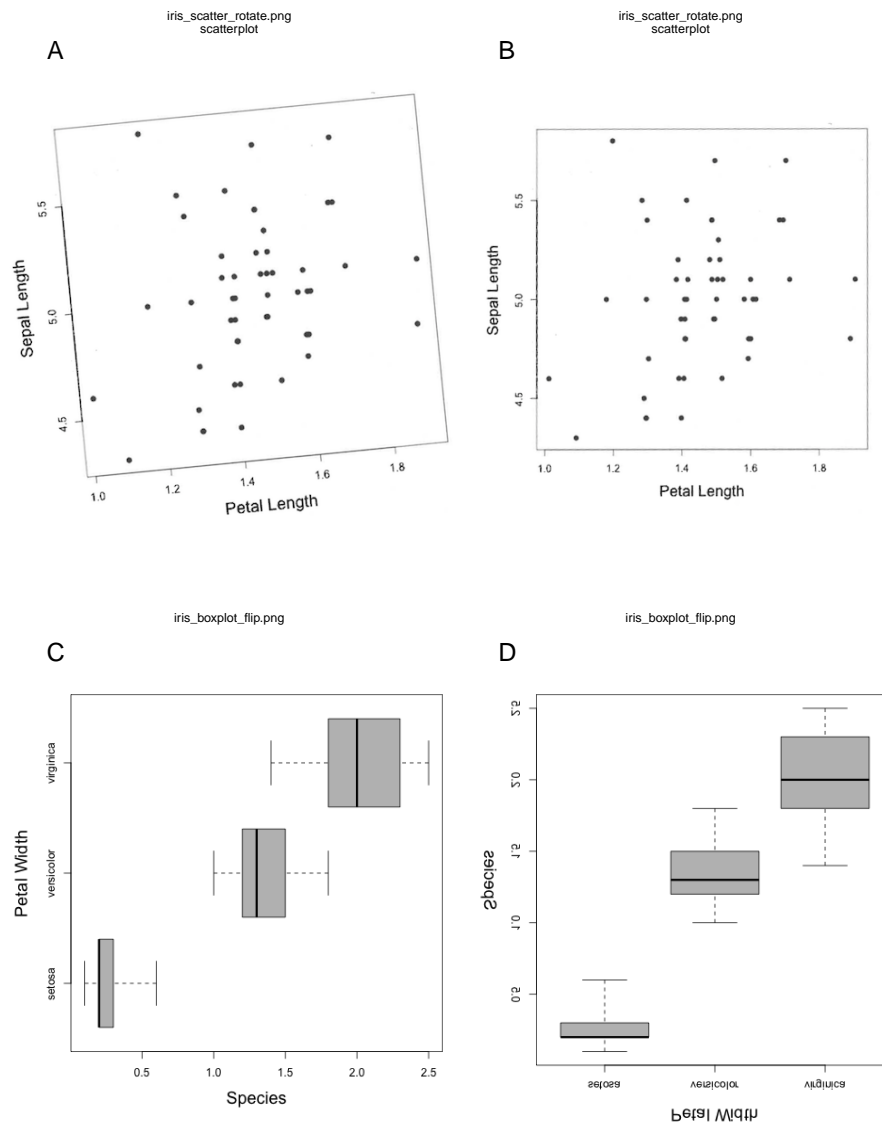


Figure 5: Figure rotation. A) and B) show how non-aligned images can be realigned through user defined rotation. C) and D) show how figures can be re-orientated so as to aid data input.

334 Instead of providing the user with summary statistics it will return a list containing
 335 four slots for each of the figure types (mean error, box plot, histogram and scatter
 336 plots). An example of a data object returned from digitising figures is as follows:

```
>R str(data)
```

```
List of 3
```

```
$ mean_error :List of 1
```

```
..$ 002_Doe_2013_Fig1.png:'data.frame': 3 obs. of 5 variables:
```

```
.. ..$ id : Factor w/ 3 levels "setosa","versicolor",...: 1 2 3
```

```
.. ..$ mean : num [1:3] 5 5.93 6.59
```

```
.. ..$ error : num [1:3] 0.111 0.148 0.178
```

```
.. ..$ n : num [1:3] 50 50 50
```

```
.. ..$ variable: chr [1:3] "Sepal length" "Sepal length" "Sepal length"
```

```
$ hist :List of 1
```

```
..$ 003_Doe_2013_Fig3.png:'data.frame': 8 obs. of 3 variables:
```

```
.. ..$ midpoints: num [1:8] 4.3 4.5 4.7 4.9 5.1 ...
```

```
.. ..$ frequency: num [1:8] 4 5 7 12 11 6 2 3
```

```
.. ..$ variable : chr [1:8] "Sepal length" "Sepal length" ...
```

```
$ scatterplot:List of 1
```

```
..$ 001_Anderson_1935_Fig1.png:'data.frame': 83 obs. of 8 variables:
```

```
.. ..$ id : Factor w/ 2 levels "setosa","versicolor": 1 1 1 1 1 ...
```

```
.. ..$ x : num [1:83] 2.3 2.9 3 3 3 ...
```

```
.. ..$ y : num [1:83] 4.5 4.4 4.41 4.3 4.8 ...
```

```
.. ..$ group : num [1:83] 1 1 1 1 1 1 1 1 1 1 ...
```

```
.. ..$ col : Factor w/ 2 levels "red","green": 1 1 1 1 1 1 1 1 1 1 ...
```

```
.. ..$ pch : num [1:83] 19 19 19 19 19 19 19 19 19 19 ...
```

```
.. ..$ y_variable: chr [1:83] "Sepal length (mm)" "Sepal length (mm)" ...
```

```
.. ..$ x_variable: chr [1:83] "Sepal width (mm)" "Sepal width (mm)" ...
```

337 Here, the user can easily access the list of processed scatter plot data by simply
338 extracting the scatter plot slot:

```
>R scatterplot <- data$scatterplot
```

339 Adding sample sizes to previous Digitisations

340 In many cases important information, such as sample sizes, may not be readily available
341 or clear when digitising figures. In these circumstances users will have answered ‘no’ to
342 the question about whether they have sample sizes or not while digitising. To expedite
343 finding and adding in these sample sizes to do the necessary calculations (if for example
344 a figure presented 95% CI’s or standard errors), `metaDigitise()` has a specific edit
345 option that allows users to enter in previously omitted sample sizes. It works by first
346 identifying the missing sample sizes in the digitised output, re-plotting the relevant
347 figure and then prompting the user to enter the sample sizes for the relevant groups in
348 the figure, one by one. As an example, assume that we were missing sample sizes for
349 two groups in 002_Doe_2013_Fig1.png:

filename	group_id	variable	mean	error	error_type	n	r	sd	plot_type
002_Doe_2013_Fig1.png	setosa	Sepal length	5.00	0.11	se	NA	NA	NA	mean_error
002_Doe_2013_Fig1.png	virginica	Sepal length	6.59	0.18	se	NA	NA	NA	mean_error

350 Here, we can see that we are missing the sample sizes for setosa and virginica, and as a
351 result, sd is not calculated because `metaDigitise()` needs this information to make the
352 calculation. If the user found this information after contacting the authors for
353 clarification then they can add these back in as follows:

```
R> digitised_data <- metaDigitise("../FiguresToExtract")
```

Do you want to...

1: Process new images

2: Import existing data

3: Edit existing data

Selection:

R> 3

Choose how you want to edit files:

1: Cycle through images

2: Choose specific file to edit

3: Enter previously omitted sample sizes

Selection:

>*R* 3

354 metaDigitise() will replot the figure after this and list, only the groups missing data,
355 for which the user can then update the data. This is then re-integrated back into the
356 data automatically and the sd calculated.

Group " setosa ": Enter sample size

R> 50

Group " viriginica ": Enter sample size

R> 50

357 Inter-observer Variability and Validation

358 Inter-observer variability in digitisations

359 In order to evaluate the consistency of digitisation using **metaDigitise** between users, we
360 simulated a dataset of two traits with two different groups. These data were then used
361 to construct plots of the four different types (scatterplot, mean and error, histogram
362 and boxplots). Each variable was plotted twice for each given plot type (figures were
363 modified slightly to give users a sense that they were digitising new data) generating a
364 total of 14 figures. 14 independent digitisers were provided with a directory with all 14
365 figures in a randomised order. Digitisers ran **metaDigitise** on their own computers,
366 across different operating systems (including Mac, Windows and Linux). Digitisers
367 varied in their level of experience, from people with experience of meta-analyses or
368 comparative work to those without any science background. We asked users to digitise
369 all 14 figures and collected the mean, standard deviation and correlation coefficient (for
370 scatterplots) generated by `metaDigitise()` for every plot digitised. We transformed
371 these data to standardized differences as

$$\frac{\theta - \hat{\theta}}{\hat{\theta}} \quad (8)$$

372 where θ is the estimate value and $\hat{\theta}$ is the true value, meaning that deviations were
373 percentage differences from the true summary statistics. The correlation coefficient
374 deviation was not divided by the true value, as it is already on a standardised scale.
375 This deviation can be seen as a measure of bias. The resulting data was used to assess
376 between- and within- user variability (i.e., the intra-class correlation coefficient) in the
377 data. This was done using linear mixed effect models with user identify as a random
378 effect. Standardised mean, standard deviation and correlation coefficients were used as
379 response variables in separate models. Sampling variance for ICC estimates was

380 generated based on 1000 parametric bootstraps of the model and the significance was
 381 tested using likelihood ratio tests. These models were run using the **lme4** (Bates et al.,
 382 2015) and **rptR** (Stoffel, Nakagawa & Schielzeth, 2017) packages in R.
 383 If digitisations were consistent across all users then we should find no significant
 384 between user variability in the data. Indeed, across plot types we found no evidence for
 385 any inter-observer variability in digitisations for the mean (ICC = 0, 95% CI = 0 to
 386 0.029, $p = 1$), standard deviation (ICC = 0, 95% CI = 0 to 0.033, $p = 0.5$) or
 387 correlation coefficient (ICC = 0.053, 95% CI = 0 to 0.296, $p = 0.377$). There were was
 388 little bias between digitised and true values, on average 1.63% (mean = 0.02%, SD =
 389 4.9%, $r = -0.03\%$) and overall there were only small absolute differences between
 390 digitised and true values, deviating, on average 2.18% (mean = 0.40%, SD = 5.81%, r
 391 = 0.33%) across all three summary statistics.
 392 SD estimates from digitisations are clearly more prone to error than means or
 393 correlation coefficients. If the mean absolute difference is calculated for each plot type,
 394 we can see that this effect is driven mainly by extraction from boxplots and histograms
 395 (% difference):

boxplot	histogram	mean_error	scatterplot
15.805	5.210	1.500	0.433

396 This is because SD estimation from the summary statistics extracted from boxplots is
 397 more error prone, especially at small sample sizes (Wan et al., 2014).

398 **Testing the accuracy of digitisations**

399 To test how accurate **metaDigitise** is at matching points to their true values, we
 400 generated four random scatterplots, each with 20 data points, and digitised these with
 401 **metaDigitise()**. This was done by one digitiser, as there is no detectable between user
 402 variation. Data digitised using **metaDigitise** was essentially perfectly correlated with

the true simulated data for both the x -variable (Pearson's correlation; $r = 0.9999915$, $t = 2137.4$, $df = 78$, $p < 0.001$) and y -variable ($r = 0.9999892$, $t = 1897.8$, $df = 78$, $p < 0.001$).

Discussion and Conclusions

Although **metaDigitise** is already very flexible, and provides functionality not seen in any other package (Table 1) it is clear that there are some functions that it does not perform. A notable feature that **metaDigitise** lacks is automated point detection. Point detection is available in several packages (Table 1). However, from our experience of using these functions, manual digitising is more reliable and often equally as fast. Particularly given that calibration (for point detection) needs to be done for each plot individually in any case. Additionally, auto-detection often misses many points which then subsequently need to be manually added. Based on tests of **metaDigitise** (see above), figures can be extracted in around 1-2 minutes, including the entry of metadata. As a result, we do not believe that current automated point detection provides substantial benefits in terms of time or accuracy.

Another feature that **metaDigitise** (currently) lacks, is an ability to zoom in on plots. Zooming may enable users to gain greater accuracy when clicking on points. However, from our own experience (and indeed from the results above), if you are using a reasonably sized screen then the accuracy is already high from these programs, and there is not much gain to be had from zooming in on points in many circumstances.

In contrast to some other packages, **metaDigitise** currently also does not extract lines from figures. In our own experience, line extraction is not particularly useful for meta-analysis, although we recognise that it may be useful in other fields. Should a user like to extract lines with **metaDigitise**, we would recommend extracting data as a

Function	metaDigitise	GraphClick ¹	DataThief ²	DigitizeIt ³	WebPlotDigitizer ⁴	metagear ⁵	digitize ⁶
Scatterplots	✓	✓	✓	✓	✓	✓ ⁷	✓
Mean and error plots	✓	✓	✓	×	×	✓ ⁷	×
Boxplots	✓	×	×	×	×	×	×
Histograms	✓	×	×	×	✓ ⁷	×	×
Graph rotation ⁸	✓	✓	✓	✓	✓	×	×
Groups	✓	✓	×	✓	✓	×	×
Entry of metadata	✓	×	×	×	×	×	×
Summarising data	✓	×	×	×	×	×	×
Multiple image processing	✓	×	×	×	×	×	×
Reproducible ⁹	✓	✓	✓	×	✓	×	×
Automated point detection	×	✓	×	✓	✓	✓	×
Line extraction	×	✓	✓	✓	✓	×	×
Zoom	×	✓	✓	✓	✓	×	×
Log axis	×	✓	✓	✓	✓	×	×
Dates	×	×	✓	×	✓	×	×
Asymmetric error bars	×	×	✓	×	×	×	×
Freeware	✓ ¹⁰	✓ ¹¹	✓ ¹¹	×	✓ ¹¹	✓ ¹⁰	✓ ¹⁰

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

⁷ Only automated, no manual extraction.

⁸ Or handles rotated graphs.

⁹ Allows saving, re-plotting and editing of data extraction.

¹⁰ R package.

¹¹ Standalone software.

Table 1: Comparison of functionality between different digitisation softwares.

428 scatter plot, and clicking along the line in question. A model can then be fitted to these
429 points (setting the argument "summary = FALSE" in **metaDigitise** - will provide access
430 to the processed data) to estimate the parameters needed.

431 Finally, **metaDigitise** currently does not allow for asymmetric error bars. At present
432 this is a deliberate omission, as it is not clear how best to derive SD from such data,
433 given also that such asymmetric error bars may represent different things in different
434 figures.

435 Descriptive statistics are usually the most robust sources of information for calculating
436 effect size statistics (Noble et al., 2017). These are most often presented in figures.
437 Users may therefore also want to compare effect size estimates from inferential statistics
438 with those derived from descriptive statistics (obtained for example using **metaDigitise**)
439 from a paper. Comparing these different effects sizes can be useful in identifying
440 uncertainties and problems within a paper. In the future, we hope to provide functions
441 to easily convert inferential statistics to standardised effect size estimates, which can
442 seamlessly be integrated with summary statistics from **metaDigitise**, to calculate
443 equivalent standardised effect size estimates and their sampling variance.

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

450 Acknowledgments

451 We thank the I-DEEL group at UNSW for incredibly useful feedback, and a host of
452 colleagues for testing, providing feedback and digitising including: Rose O'Dea, Fonti

453 Kar, Malgorzata Lagisz, Julia Riley, Diego Barneche, Erin Macartney, Ivan Beltran,
 454 Gihan Samarasinghe, Dax Kellie, Jonathan Noble, Yian Noble and Alison Pick. JLP
 455 was supported by a Swiss National Science Foundation Early Mobility grant
 456 (P2ZHP3_164962), DWAN was supported by an Australian Research Council Discovery
 457 Early Career Research Award (DE150101774) and UNSW Vice Chancellors Fellowship
 458 and SN an Australian Research Council Future Fellowship (FT130100268).

459 References

- 460 Arizona-Software (2008) *GraphClick Software, Version 3.0*.
- 461 Bates, D., Maechler, M., Bolker, B. & Walker, S. (2015) Fitting Linear Mixed-Effects
 462 Models Using lme4. *Journal of Statistical Software*, **67**, 1–48.
- 463 Bland, M. (2015) Estimating Mean and Standard Deviation from the Sample Size,
 464 Three Quartiles, Minimum, and Maximum. *International Journal of Statistics in*
 465 *Medical Research*, **4**, 57–64.
- 466 Borenstein, M., Hedges, L., Higgins, J. & Rothstein, H. (2009) Introduction to
 467 meta-analysis. *John Wiley Sons. Ltd. West Sussex, UK*.
- 468 Bormann, I. (2012) *Digitizelt Software, Version 2.0*. Braunschweig, Germany.
- 469 Glass, G. (1976) Primary, secondary, and meta-analysis research. *Educational*
 470 *Researcher*, **5**, 3–8.
- 471 Ihle, M., Winney, I.S., Krystalli, A. & Croucher, M. (2017) Striving for transparent and
 472 credible research: practical guidelines for behavioral ecologists. *Behavioral Ecology*,
 473 **28**, 348–354.
- 474 Koricheva, J., Gurevitch, J. & Mengersen, K. (2013) Handbook of Meta-Analysis in
 475 Ecology and Evolution. *Princeton University Press, Princeton, New Jersey*.

476 Lajeunesse, M.J. (2016) Facilitating systematic reviews, data extraction, and
 477 meta-analysis with the metagear package for R. *Methods in Ecology and Evolution*, **7**,
 478 323–330.

479 Nakagawa, S., Noble, D.W., Senior, A.M. & Lagisz, M. (2017) Meta-evaluation of
 480 meta-analysis: ten appraisal questions for biologists. *BMC Biology*, **15**, 18; DOI
 481 10.1186/s12915-017-0357-7.

482 Noble, D.W., Lagisz, M., O’Dea, R.E. & Nakagawa, S. (2017) Nonindependence and
 483 sensitivity analyses in ecological and evolutionary meta-analyses. *Molecular Ecology*,
 484 **26**, 2410–2425.

485 Parker, T.H., Forstmeier, W., Koricheva, J., Fidler, F., Hadfield, J., En Chee, Y., Kelly,
 486 C.D., Gurevitch, J. & Nakagawa, S. (2016) Transparency in Ecology and Evolution:
 487 Real Problems, Real Solutions. *Trends in Ecology and Evolution*, **31**, 711–719.

488 Peng, R.D. (2011) Reproducible research in computational science. *Science*, **334**, 1226.

489 Peng, R.D., Dominici, F. & Zeger, S.L. (2006) Reproducible epidemiologic research.
 490 *American Journal of Epidemiology*, **163**, 783–789.

491 Poisot, T. (2011) The digitize package: extracting numerical data from scatterplots.
 492 *The R Journal*, **3**, 25–26.

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

494 Sandve, G.K., Nekrutenko, A., Taylor, J. & Hovig, E. (2013) Ten simple rules for
 495 reproducible computational research. *PLoS Computational Biology*, **9**, e1003285.

496 Stoffel, M.A., Nakagawa, S. & Schielzeth, H. (2017) rptR: repeatability estimation and
 497 variance decomposition by generalized linear mixed-effects models. *Methods in*
 498 *Ecology and Evolution*, **8**, 1639–1644.

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

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