# Electronically Cataloging Butterflies

*Part IB Group Project - Group D*

## Tutorial

James Alner, Yulong Huang, Francesca Iovu, Jack Parkinson, Suzie Welby and Abigail Wilkinson
2020

# Prerequisites

## Scanning Instructions:

1) We assume that the user will scan the logbooks with a phone app since they are fragile documents.

2) The user can use a desk phone holder clip to hold the phone in the same position above the pages and avoid unnecessary shadows.

3) The user should try their best to avoid glares in the scans.

4) They need to make sure that the file is of type PDF of JPEG, and has no cover pages or pages other than simple spreadsheets. Our program only works with consistently formatted spreadsheets.

5) The user should make sure the PDF or JPEG is oriented correctly, with the writing horizontal. The program does not care if the file is in landscape or portrait mode.

6) The user should make sure to scan one page at a time, even if a logbook page is made out of two physical pages. This way a higher resolution image will be produced and word recognition will be more accurate.

For example:

The scan below comes from a logbook that has 2 physical pages per logbook page (you can see the black dots partially on the first page and fully on the second one) This gives a bad image quality and makes the words harder to get out. Plus, the corners are hard to find on a white background.
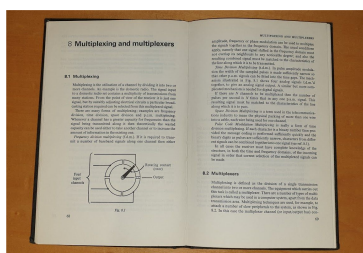
Instead, you can scan each page separately and on a black page background so we can identify the corners easier. Also, make sure that the corners are flat (can use blu tack for that). This is how the logbook page above looks when properly scanned:



For an actual book, in order to make sure one page is scanned at a time,
a black page (larger than the book) can be put in between the top page currently being scanned and the previous one. Another black page can be used to cover up the second page on top. For example, for a book looking like this:



Place black paper like so to scan the 1st page:

And this to scan the 2nd page;



After opening the app, you will see this:

# Using the interface

There are three buttons available from the homepage:
- Upload PDF
- Convert CSV to Standard Format
- View Help Guide

## Upload pdf

1. Click file icon to upload
2. Select pdf file
3. Click open
4. Click Read Page
5. You will be asked the following:



6. Input the number of scanned pages that make up a logbook page:
   For example, for the logbook below, the number is 2:



7. The PDF will be displayed like so:

This is an example for one scanned page per logbook page, for more than one, they would be merged together. Like this:



8. Now you can add or delete columns and move them around. Make sure the columns line up with the actual ones. Please also remove the header and the footer. You can also get rid of margin columns that are just white and have no header name.
9. Rename all the columns to have the correct headers, these are the headers that will appear in the finished CSV.
10. Before confirming, you can add a txt dictionary that will be used for spelling correction



The dictionary should look like this:

Window title: NHM_butterfly_dict - Notepad

```
mesotype
sclerocona
macariini
cribrella
eupithecia
alnifoliae
cnephasiini
primaria
yellow-ringed
empetrella
reticulata
erminea
pallidactyla
ruficinctata
fuscocuprea
biren
valerianata
chimabachidae
luculella
lancealana
molliculana
pigmy
adelaidae
limnaecia
munitata
fen
pirithous
atricapitana
centonalis
porrittia
shining
subcognata
roseana
punctata
uncula
douglasii
pallida
picaepennis
ruddy
ambigua
thoracella
```

11. After you've made sure the columns are correctly placed and the header and footer removed, you can click confirm.
12. The program will work its magic now!

## Convert CSV to Standard format

This feature can be used to convert an existing CSV to a standard format (i.e. to be uploaded to an existing database in a specific format)

1. First select a CSV using the file browser by clicking the folder icon and click next to continue. The only accepted format is a non-empty .CSV file.



Next

2.

Once on the Rules Page, choose a column to split using the drop-down selection. Type the name of a new column and include the advanced parameters. This parameter can take multiple forms:

 - A single index (i.e. "0" to select the first word in each cell), with indices starting from 0.

 -A set of indices (i.e. "[0, 1, 3]" selects the 1st, 2nd and 4th words and joins them in that order with the joiner)

 - A range of indices (i.e. "1:5" selects the 2nd to 6th words and joins them with the joiner).

 -A wildcard "*" which will take the words which are not assigned to any other column for each cell.

Note that negative indices can be used to index from the end (i.e. "-2" refers to the 2nd to last word).

Multiple new columns can be added using the + button, and can be removed using the x button. To remove the current rule press the x button when there is only one new column.

Select the resolution from the drop down list, the default is that there is no clash. You can also provide which characters to split or join words on, the default is a space.

Once all rules have been filled in and checked, click the Confirm Rules and Continue button to move to the Mappings page.

3. Having created new columns using the rules page, you can then specify what to put in each column of the standard format.

You can either define a mapping - takes the values from 1 or more columns and combines them - or a constant value for every record in that column.

Once done with these mappings, you can save the CSV by clicking 'Confirm all Mappings and Save'

## View Help Guide:

1. Will open this tutorial pdf in the operating system's default application.