



Deposited via The University of Sheffield.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/id/eprint/227048/>

Version: Accepted Version

Article:

He, Y., Cooney, C.R., Maddock, S. et al. (2025) PhenoLearn: a user-friendly toolkit for image annotation and deep learning-based phenotyping for biological datasets. *Journal of Evolutionary Biology*, 38 (8). pp. 1152-1162. ISSN: 1010-061X

<https://doi.org/10.1093/jeb/voaf058>

© The Author(s) 2025. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

Reuse

This article is distributed under the terms of the Creative Commons Attribution (CC BY) licence. This licence allows you to distribute, remix, tweak, and build upon the work, even commercially, as long as you credit the authors for the original work. More information and the full terms of the licence here:

<https://creativecommons.org/licenses/>

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.

PhenoLearn: A user-friendly Toolkit for Image Annotation and Deep Learning-Based Phenotyping for Biological Datasets

Authors: Yichen He^{1,2*}, Christopher R. Cooney¹, Steve Maddock³, Gavin H. Thomas^{1,4}

¹ Ecology and Evolutionary Biology, School of Biosciences, University of Sheffield; Alfred Denny Building, Western Bank, Sheffield S10 2TN, UK.

² Department of Life Sciences, Natural History Museum; Cromwell Rd, South Kensington, London SW7 5BD, UK

³ Department of Computer Science, University of Sheffield; Regent Court, 211 Portobello, Sheffield S1 4DP, UK.

⁴ Bird Group, Department of Life Sciences, Natural History Museum; Akeman Street, Tring, HP23 6AP, UK.

*Corresponding authors. Email: csyichenhe@gmail.com

Acknowledgements. We thank Thomas Guillerme, Kathryn Harris and Eleftherios Ioannou for testing this toolkit.

Funding. This work was funded by a Leverhulme Early Career Fellowship (ECF-2018-101) and Natural Environment Research Council Independent Research Fellowship (NE/T01105X/1) to C.R.C, and a European Research Council grant (615709, Project 'ToLERates') and Royal Society University Research Fellowship (UF120016, URF\R\180006) to G.H.T.

Authors contributions. All authors designed the tool; Y.H. developed the tool; C.R.C. and G.H.T. tested the tool. Y.H. wrote the manuscript with input from all authors.

Competing interests. All authors have no competing interests.

Data and materials availability. All code, datasets, and binaries used in this study are publicly archived and available:

- **Source code:** The PhenoLearn source code is available on GitHub for continued development: <https://github.com/echanhe/phenolearn>. A snapshot of the code corresponding to the version used in this paper has been archived on Zenodo and assigned a DOI: <https://zenodo.org/records/15350513>.
- **Example datasets:** The bird and Littorina test datasets used for evaluation are available on Zenodo: <https://zenodo.org/records/8152784>.
- **Binary executable:** The compiled Windows binary of the PhenoLabel annotation tool is also archived on Zenodo: <https://zenodo.org/records/10909841>.

37

38 **Abstract**

39 The digitisation of natural history specimens has unlocked opportunities for large-scale
40 phenotypic trait analysis. In recent years, deep learning has shown significant results in
41 accurately predicting annotations on 2D specimen photographs. However, it can be
42 challenging for biologists without extensive related expertise to easily use deep learning.
43 Here, we introduce PhenoLearn, a toolkit developed for biologists to generate
44 annotations on 2D specimen images using deep learning. PhenoLearn integrates
45 graphical user interfaces (GUIs) within its two main modules, PhenoLabel for image
46 annotation and PhenoTrain for model training and prediction. GUIs increase
47 accessibility and reduce the need for computational expertise, allowing biologists to
48 intuitively go through a workflow of labelling training sets, using deep learning, and
49 reviewing predictions in the same tool. We demonstrate PhenoLearn's capabilities
50 through a case study involving the segmentation of plumage areas on bird images,
51 showcasing prediction accuracy and the running time with and without GPU,
52 highlighting its potential to generate annotations with minimal computational cost and
53 time. The toolkit's modular design and flexibility ensure adaptability, allowing for
54 integration with other tools amidst rapidly evolving deep learning approaches.
55 PhenoLearn bridges the gap between specimen digitisation and downstream analysis,
56 providing biologists with broader access to deep learning. The source code, installation
57 guides, tutorials with screenshots, and a small demo dataset for PhenoLearn can be
58 found at <https://github.com/echanhe/phenolearn>.

59 **Keywords:** Deep Learning, Phenotyping, Image Annotation, Phenotypic Trait, Toolkit
60 with User Interface.

ORIGINAL UNEDITED MANUSCRIPT

61 Introduction

62 The process of measuring phenotypic traits on 2D digitised specimen images is
63 increasingly used to phenotype specimens for a range of tasks. Through the use of
64 annotations such as points (Chang & Alfaro, 2016; Zelditch et al., 2004) and
65 segmentations (Cooney et al., 2022; Y. He et al., 2022), researchers can extract and
66 analyse a variety of morphological measurements from specimens to provide insights
67 into evolutionary and ecological questions. Digitisation allows rapid and non-invasive
68 measurements of natural history collections and mobilises specimens for further
69 analyses, helping to unlock their full potential. Techniques such as tray scanning
70 (Blagoderov et al., 2012) have significantly accelerated the digitisation of entomological
71 collections, by leveraging robotic automation to automatically capture 2D images of
72 specimens directly from museum trays. In addition, many computational tools for
73 analysing phenotypes like shape (Adams & Otárola-Castillo, 2013) and colouration
74 (Maia et al., 2019) have been developed, expanding the breadth of tools available to
75 analyse phenotypic traits. However, manually pre-processing images (e.g., placing
76 annotations) is time-consuming, especially with large datasets such as hundreds of
77 thousands of observations (Cooney et al., 2022). To prevent manual annotation from
78 becoming a bottleneck for mobilising large digital datasets, efficient high-throughput
79 data extraction tools are essential.

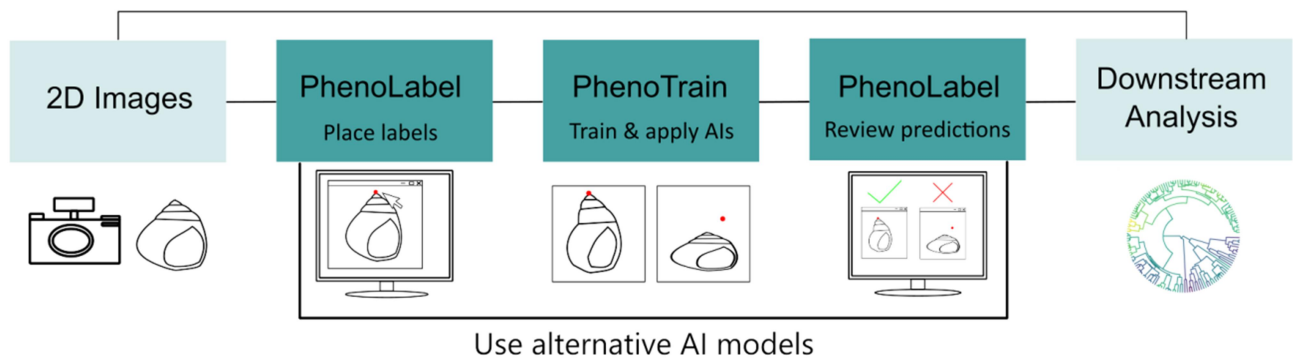
80 Classic computer vision algorithms like thresholding, connected components and region
81 growing have been used for extracting phenotypic information from images,
82 representing a significant increase in measurement speed compared to manual
83 methods (Lürig, 2022; Pennekamp & Schtickzelle, 2013). Deep learning-based methods
84 have recently become state-of-the-art for various computer vision tasks, including object
85 segmentation in images with complex backgrounds. In particular, deep learning
86 applications for measuring digitised specimens have demonstrated success with
87 different types of annotations, including points (Mathis et al., 2018; Porto & Voje, 2020),
88 bounding boxes (John et al., 2024; Shedrawi et al., 2024) and segmentations (Y. He et
89 al., 2022; Schwartz & Alfaro, 2021). These methods yield high-throughput pipelines and
90 accurate results, illustrating the potential for expanding deep learning to other
91 biological datasets. However, several barriers remain preventing the widespread
92 application of deep learning in ecology and evolutionary biology.

93 A significant barrier to the wider adoption of deep learning is the generally high level of
94 technical skill required for implementation. This issue is often compounded by the lack
95 of intuitive platforms that allow non-specialists to use deep learning for phenotyping.
96 Recent phenotyping toolkits, such as DeepLabCut (Mathis et al., 2018) and Argos (Ray
97 & Stopfer, 2022), have focused on improving accessibility through graphical user
98 interfaces (GUIs). The integration of deep learning models with GUIs can greatly
99 increase accessibility, allowing researchers with limited technical knowledge to utilise
100 these advanced techniques. Furthermore, the development of fully-integrated toolkits for
101 performing a complete workflow, including labelling training sets, training models and
102 reviewing predictions, can significantly improve the accessibility and efficiency for
103 biologists to apply deep learning in biological research. Such a toolkit, tailored for
104 extracting traits from 2D specimen photographs for ecological and evolutionary studies,

105 would serve as a much-needed bridge between digitisation and downstream biological
106 analysis.

107 Here, we introduce PhenoLearn, a user-friendly toolkit for generating annotations using
108 deep learning. PhenoLearn comprises two main modules, PhenoLabel and PhenoTrain,
109 covering three main functions (**Figure 1**). PhenoLabel implements both image labelling
110 and reviewing, whereas PhenoTrain implements the functions for deep learning. As an
111 open-source tool with GUIs, PhenoLearn aims to minimise the computational expertise
112 required to generate point or segmentation predictions using deep learning for 2D
113 biological image datasets. While PhenoLearn is designed to facilitate the entire
114 annotation generation workflow, its modular design allows users to use individual
115 functions for desired tasks. For instance, PhenoLabel can be used to review predictions
116 from other methods. Likewise, labels generated elsewhere can be used to train models
117 implemented in PhenoTrain. The PhenoLearn pipeline has already been successfully
118 used to generate annotations in several large-scale research projects (Cooney et al.,
119 2022; Y. He et al., 2022, 2023). In this paper, we provide a detailed explanation, a user
120 guideline, and an example of using PhenoLearn.

121



122

123 **Figure 1. Workflow overview of using PhenoLearn to generate annotations for**
124 **biological datasets.** Green boxes indicate steps involving PhenoLearn modules (PhenoLabel
125 and PhenoTrain), which offer a connection between digitisation (2D imaging) and
126 downstream biological analysis.

127 Installation

128 PhenoLearn was developed using Python 3, with the following libraries and their
129 versions tested during its development:

- 130 • Python: 3.10
- 131 • PyQt: 5.15.9
- 132 • NumPy: 1.25.1
- 133 • pandas: 2.0.3
- 134 • opencv-python: 4.8.0.74
- 135 • PyTorch: 2.0.1
- 136 • TensorBoard: 2.13.0

137 For deep learning, PhenoLearn is optimised to utilise NVIDIA graphics processing units
138 (GPUs) through CUDA (<https://developer.nvidia.com/cuda-toolkit>). While it is possible to
139 train models using the CPU on systems without CUDA-supported GPUs, this will
140 generally lead to slower running time. We recommend using a GPU with at least 8GB of
141 video memory for faster running time.

142 PhenoLearn's two main modules, PhenoLabel and PhenoTrain have their own GUIs.
143 PhenoLabel implements the labelling and reviewing functions and can be accessed by
144 running **phenolabel.py**. PhenoTrain handles deep learning training and prediction and
145 is accessed by running **phenotrain.py**. It was tested on Windows 10, macOS 13.6, and
146 Ubuntu 22.04.3 LTS.

147 The source code, installation guides, tutorials with screenshots, and a small demo
148 dataset for PhenoLearn can be found at <https://github.com/echanhe/phenolearn>.
149 Datasets used in the example section are available at
150 <https://zenodo.org/records/8152784>. For Windows users, a binary version of
151 PhenoLabel (e.g., a .exe file) is available at <https://zenodo.org/records/10909841>.
152 Detailed file introductions can be found in the supplementary material.

153 Design and Implementation

154 Labelling

155 This section outlines using PhenoLabel for labelling, including creating a project, placing
156 points/segmentations, and managing progress. To start, select 'Open Dir' in the File
157 menu (**Figure 2a**) to open a folder of images for labelling. PhenoLabel uses the imread
158 function from opencv-python (Bradski, 2000), which supports common formats including
159 jpg, png and tiff. PhenoLabel lists all images in the File panel (**Figure 2d**) and displays
160 the selected image in the Main panel (**Figure 2e**). Users can zoom the image and view
161 the cursor coordinate and RGB values in the status bar (**Figure 2g**).

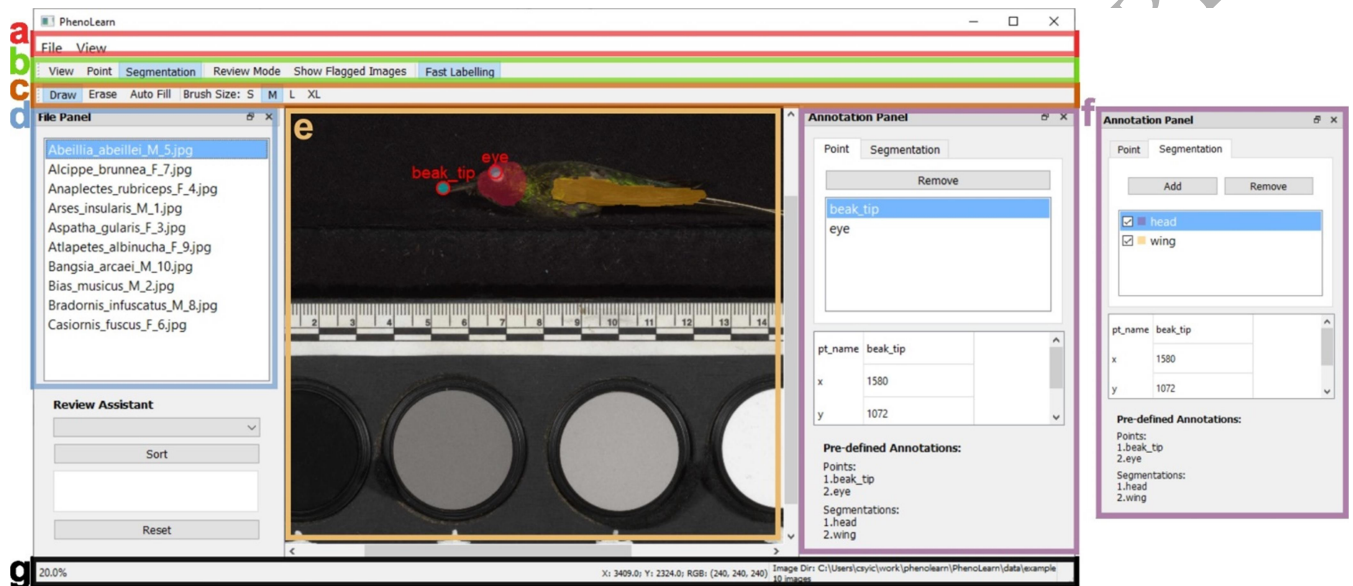
162 Users can place points or segmentations in the Main panel. For points, click the 'Point'
163 button on the Toolbar (**Figure 2b**) and left-click the image. Points can be named via a
164 dialogue box, either inputting a new name or selecting from a dropdown menu. Existing
165 points can be modified or deleted in the Annotation panel (**Figure 2f**). PhenoLearn
166 records vertical (y) coordinates from the top of the image downward. This is the
167 standard convention used in many Python-based image processing libraries, such as
168 opencv-python. In contrast, tools like tpsDig (Rohlf, 2006) and many R-based image
169 analysis tools typically record coordinates from the bottom up. Users working with
170 tpsDig datasets should be aware of this difference.

171 For segmenting, click the 'Segmentation' button on the Toolbar, and a Segmentation
172 Toolbar (**Figure 2c**) appears. Segmentation classes must be named using the 'Add'
173 button in the Segmentation tab (**Figure 2f**). Then by activating the 'Draw' button in the
174 Segmentation Toolbar and holding the left mouse button, users can use a paintbrush to
175 draw the region of interest (ROI). Each segmented ROI is automatically assigned a
176 distinct colour, allowing users to easily differentiate between them. Segmented areas
177 can be removed with the same operation with the 'Erase' button activated. To efficiently
178 segment large areas, users can outline a region and then use the 'Auto Fill' function to
179 fill the area within the outline. Four paintbrush sizes are available: S, M, L, and XL.

180 PhenoLabel's 'Fast Labelling' function automates annotation naming for cases that use
 181 consistent annotation names, eliminating repeated manual naming. This feature
 182 automatically creates annotation names for subsequent images using the annotation
 183 names from the current image. To ensure a newly placed point matches the preset point
 184 names, they need to be placed in the same order as the names displayed at the bottom
 185 of the Annotation Panel.

186 'Save' and 'Save As' in the File menu allow users to save their work in JSON format,
 187 which includes details on images and annotations. 'Open Labelling Progress' allows
 188 users to continue or review their labelling progress. Annotations can be exported to
 189 PhenoTrain in CSV or binary masks (for single-class segmentations). Two types of CSV
 190 exports are available: a point CSV file and a segmentation CSV file. Refer to Table 1 for
 191 the detailed structure of the JSON and CSV files.

192



193 **Figure 2. The PhenoLabel GUI.** (a) Menu bar: Provides functions for saving projects and
 194 loading files, (b) Toolbar: Tools for image annotation manipulation, (c) Segmentation Toolbar:
 195 Tools specifically designed for segmentation tasks, (d) File panel: Displays the loaded images
 196 and allows users to switch between images, (e) Main panel: The central workspace for image
 197 annotation, (f) Left: Annotation panel, Point tab for displaying details of point-based annotations;
 198 Right: Annotation panel, Segmentation tab for displaying details of segmentation-based
 199 annotations, (g) Status bar: Displays status such as image name, zoom level, and cursor
 200 position.
 201

202

203 **Table 1.** File structures used in PhenoLearn.

File	Description
Labelling progress file	A JSON file From: Created by the save function in PhenoLabel. Usage: Can be used to load the progress into PhenoLabel

Structure:

The file is a list of dictionaries.

- “file_name” stores the image name.
- “points” stores a list of dictionaries.
 - “name” stores the point name
 - “x” stores the x coordinate
 - “y” stores the y coordinate
 - “absence” stores if the point is missing
- “segmentations” stores a dictionary.
 - Dictionary keys are the names of the segmentations and dictionary values are the segmentations.

A segmentation is stored as a four-level nested list, which follows the format of segmentation contours extracted by OpenCV (Bradski, 2000). The format is:

- The first level corresponds to the segmentation itself.
- The second level is the contour level, where one segmentation may include one or more contours.
- The third and fourth levels pertain to the point level, with each contour having multiple points.

The example below shows a segmentation consisting of two contours. Contour 1 contains 'n' points, and Contour 2 contains 'm' points. Here <x₁₂> represents the x-coordinate of the second point in Contour 1.

Example:

```
[ {
  "file_name": "Abeillia_abeillei_M_5.jpg",

  "points":
  [
    {"name": "beak", "x": 1580, "y": 1072},
    {"name": "eye", "x": 1876, "y": 984}
  ],

  "segmentations":
  {"head":
    [
      [
        [ [ <x11>, <y11> ] ],
        [ [ <x12>, <y12> ] ], ...
        [ [ <x1n>, <y1n> ] ]
      ]
    ]
  }
}
```

ORIGINAL UNEDITED MANUSCRIPT

```

    [
      [[ <x_21>, <y_21> ]],
      [[ <x_22>, <y_22> ]], ...
      [[ <x_2m>, <y_2m> ]]
    ]
  ]
}

```

Output Point CSV file A CSV file
From: Exported by PhenoLabel or generated by PhenoTrain.
Usage: Can be imported into PhenoTrain as for training.

Structure:

The “file” column stores the image names.
 A “<point name>_x” column stores the x coordinate for a point.
 A “<point name>_y” column stores the y coordinate for a point.
 A value of -1 or an empty cell indicates the point is missing.

Example:

file	beak_x	beak_y	eye_x	eye_y
Abeillia_abeillei_M_5.jpg	1580	1072	1876	984

Output segmentation CSV file A CSV file
From: Exported by PhenoLabel or generated by PhenoTrain.
Usage: Can be imported into PhenoTrain for training.

Structure:

The “file” column stores the image names.
 The remaining columns store the segmentations.
 A segmentation is stored as a four-level nested list.
 The details and examples can be found in the “Labelling progress file” row. Here, the example only shows a four-level nested list placeholder for better readability.

Example:

file	head
Abeillia_abeillei_M_5.jpg	[[[[[]]]]]

Output binary mask A black and white image
From: Exported by PhenoLabel or generated by PhenoTrain.
Usage: Can be imported into PhenoTrain for training.

A grayscale image is saved under the same name as its input image, with background areas in black and segmentation areas in white. To prevent having the output masks replace the input images, ensure the input directory is not used as the output directory.

Property file A CSV file.

Usage: Import specific specimen properties into PhenoLabel to filter or sort images, allowing users to prioritise error-prone images first.

Structure:

The “file” column stores the image names.

Other columns store the properties.

- Categorical properties are stored as text strings.
- Numerical properties are stored as numbers.

Example:

file	id	sex
Abeillia_abeillei_M_5.jpg	5	M

204

205 Deep Learning

206 PhenoTrain allows users to train models and make predictions. This section
207 demonstrates how to set up model training and prediction in PhenoTrain.

208 Model Training

209 Before training, eleven settings are required via the Train tab of PhenoTrain (**Figure 3a**).
210 Some settings have default values derived from previous studies (Chen et al., 2017; K.
211 He et al., 2017; Y. He et al., 2022, 2023) and the PyTorch documentation (Paszke et al.,
212 2019). These defaults provide a solid starting point for various applications:

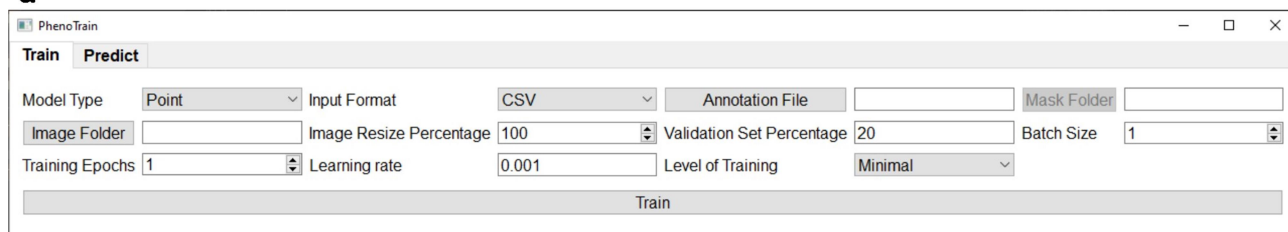
- 213 (1) *Model type*. Mask R-CNN (K. He et al., 2017) for point and DeepLabv3 (Chen et
214 al., 2017) for segmentation. Despite the availability of numerous new deep
215 learning architectures, we use Mask R-CNN and DeepLabv3 for their robust nature
216 and adaptability to various tasks. Being well-established models, there are many
217 tutorials available online that facilitate their implementation for users who want to
218 understand the detailed information.
- 219 (2) *Annotation Input format*. The default option is CSV. For single-class segmentations,
220 'Mask' option is also available for using binary masks as inputs. Please refer to
221 Table 1 for the details of the binary mask.
- 222 (3) *Annotation file*. The CSV annotation file from PhenoLabel (only applicable when
223 'CSV is selected for Setting 2).

- 224 (4) *Mask folder*. The folder of the binary masks (only applicable when 'Mask' is
 225 selected for Setting 2).
 226 (5) *Image folder*. The folder of training images.
 227 (6) *Image resize percentage*. Ranges from 1-100%, keeps aspect ratio, using nearest
 228 neighbour interpolation.
 229 (7) *Validation set percentage*. The percentage of validation images used for
 230 evaluating the model per epoch. A common split is 80/20 for training/validating.
 231 (8) *Batch size*. The number of images processed in one training iteration. The default
 232 is 1. A smaller batch size saves memory but may lead to less stable optimisation.
 233 Conversely, a larger batch size may provide better optimisation, but it uses more
 234 memory. Users need to test a set of batch sizes to find the optimal value.
 235 (9) *Training epochs*. The number of times the entire training set passes through the
 236 model. Training for more epochs may lead to better model performance. The
 237 default training epoch is set to 1. Users can estimate the training time by training
 238 for one epoch.
 239 (10) *Learning rate*. Controls the step size during the optimisation phase of training.
 240 The default learning rate for PhenoTrain is 0.001. A too-large learning rate may
 241 result in overly large steps, causing the model to miss the optimum. A too-small
 242 learning rate might lead to a very slow convergence towards the optimum.
 243 (11) *Level of training*. Controls the proportion of the model that is trained. The options
 244 are Minimal, Intermediate and Full. "Minimal" trains only the final layers,
 245 "Intermediate" trains half of the model layers, and "Full" trains the entire model.
 246 (12) *CPU/GPU*. Select whether to use the CPU or GPU for training. If GPU is selected
 247 but no GPU is available on the device, CPU will be used.

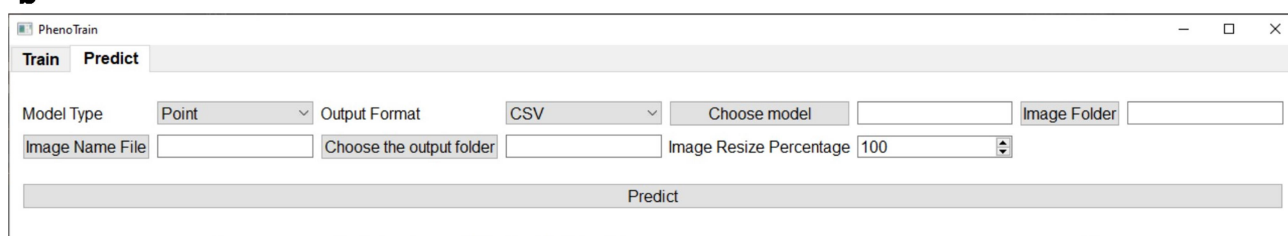
248 When the training is completed, a .pth file is saved in the 'saved_model' folder.

249

a



b



250

251 **Figure 3. The PhenoTrain GUI.** The interface has two tabs, (a) the Train tab and (b) the Predict
 252 tab. Settings for training and predicting can be specified in each tab.

253

254 The training level setting utilises transfer learning (Tan et al., 2018), focusing on training
255 with a pre-trained model. Transfer learning diverges from the approach of using
256 randomly initialised model weights, which generates poor initial predictions and can take
257 a longer training period. Instead, it leverages a pre-trained model, which effectively
258 gives the model prior knowledge gained from previous tasks. This approach can train on
259 parts of a model and achieve satisfactory results, saving both time and computational
260 resources. Both DeepLabV3 and Mask-R-CNN were pre-trained on the COCO dataset
261 (Lin et al., 2014), which is a large-scale image dataset for computer vision tasks such as
262 segmentation.

263 PhenoTrain integrates with TensorBoard (Martín Abadi et al., 2015) to visualise the
264 training progress. Logs are saved in the 'runs' folder. To view logs in TensorBoard, run
265 this command: `tensorboard --logdir=runs` in python. Upon execution, it can be
266 viewed in a web browser at <http://localhost:6006/>. Users can view and compare across
267 different training runs.

268 TensorBoard saves training and validation loss, along with evaluation metrics. Training
269 loss indicates the model's learning efficiency, while validation loss evaluates
270 performance on the validation set. Point accuracy is assessed using the pixel distance
271 (Euclidean distances between two points on an image). The Dice Score is used to
272 evaluate segmentations, based on the overlap between predicted and manual
273 segmentations. The Dice Score ranges from 0 (lowest) to 1 (highest). Average and
274 class-specific metrics for points or segmentations are stored.

275 **Generating Predictions**

276 Once a well-trained model is saved, users can generate predictions in the Predict tab
277 (**Figure 3b**) by configuring the following seven settings:

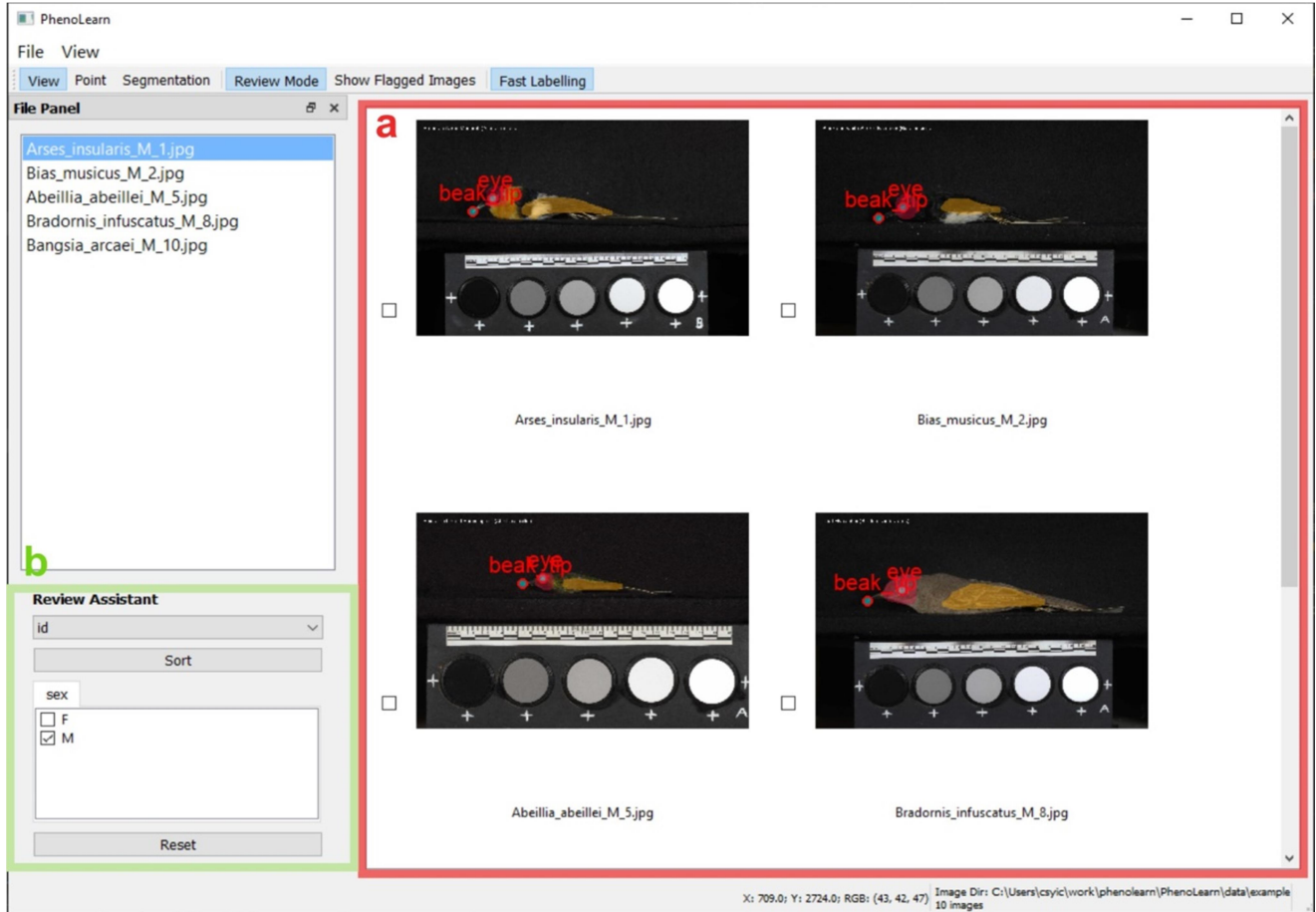
- 278 (1) *Model type*. Point or Segmentation.
- 279 (2) *Output format*. Options are CSV file or mask images (for single-class
280 segmentations only).
- 281 (3) *Choose model*. .pth file saved from training.
- 282 (4) *Image folder*. The folder of images for prediction.
- 283 (5) *Image name file*. CSV file with one column named 'file' for image names.
284 PhenoLabel can export an Image name file when no annotations are presented for
285 the images.
- 286 (6) *Choose the output folder*. A folder for the prediction file.
- 287 (7) *Image resize percentage*. Ranges from 1-100% and should be consistent with the
288 percentage used in training.
- 289 (8) CPU/GPU. Select whether to use the CPU or GPU for predicting. If GPU is
290 selected but no GPU is available on the device, CPU will be used.

291 PhenoTrain provides real-time updates during both training and prediction phases,
292 including a progress bar and elapsed time display.

293 **Reviewing predictions**

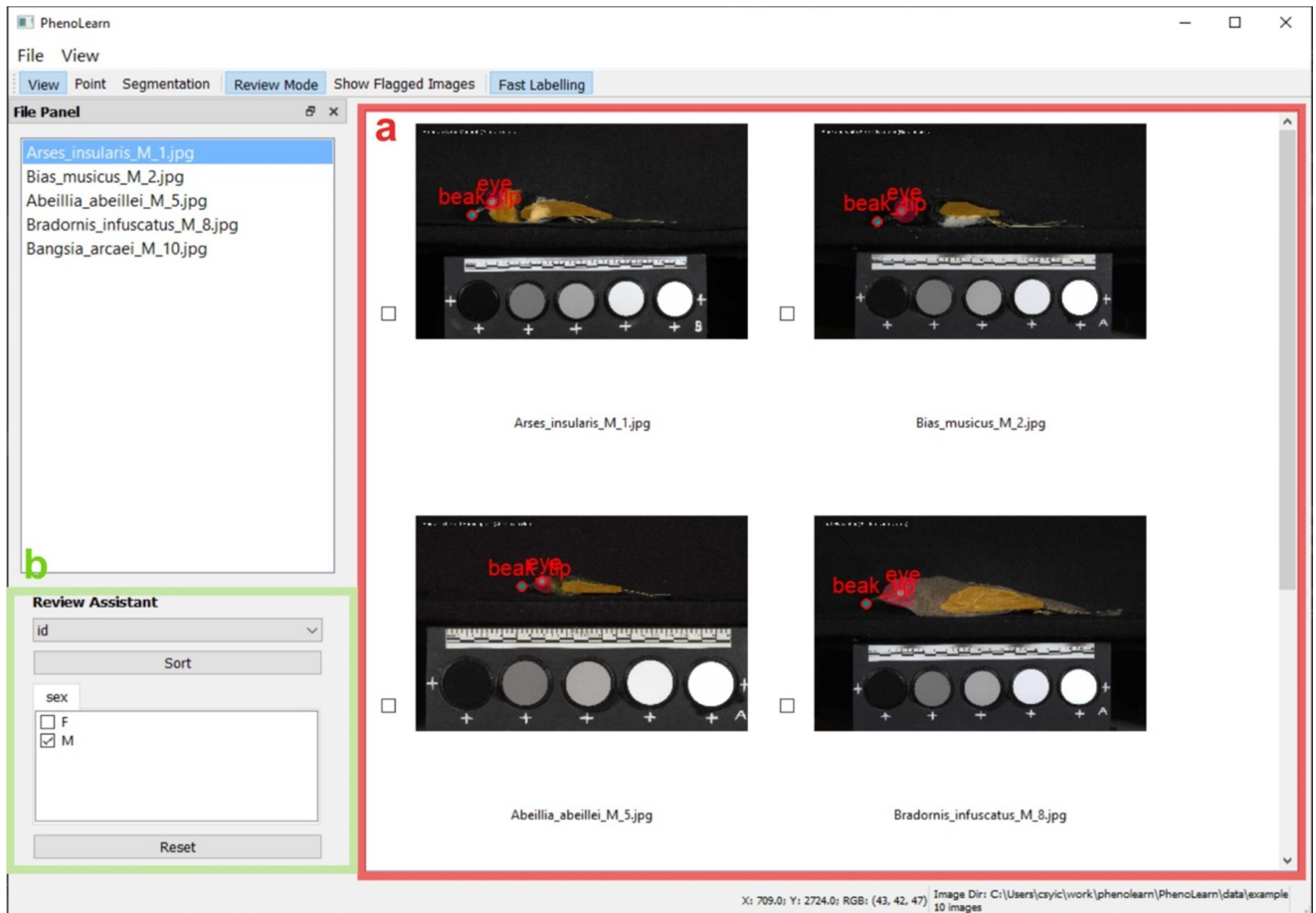
294 Deep learning predictions are not perfectly accurate, and reviewing predictions is often
295 necessary to confirm and/or improve accuracy for biological applications. To facilitate
296 this, we have incorporated two features within PhenoLabel: (1) Review Mode and (2)
297 Review Assistant to improve reviewing efficiency.

298 Users can open an image folder and import predictions (e.g., outputs from PhenoTrain) into
299 PhenoLabel, and subsequently review and improve these predictions. By activating the Review
300 Mode in the Toolbar, PhenoLabel displays multiple image thumbnails with annotations (**Figure 4**



301 **Figure 4a).** In this mode, users can quickly browse through images and flag any with
302 incorrect predictions by ticking adjacent checkboxes. After checking through thumbnails,
303 click 'Show Flagged Images' button to show only the flagged images for a more focused
304 review. Additionally, it is possible to export the predicted annotations for input into other
305 outlier detection methods and to create flagged images.
306

307 The Review Assistant improves review efficiency by leveraging specimen metadata. By
308 prioritising images with specific properties (e.g. a problematic species), users can optimise
309 accuracy and time efficiency. The Review Assistant facilitates this by offering options to sort or
310 filter images based on properties (**Figure 4**

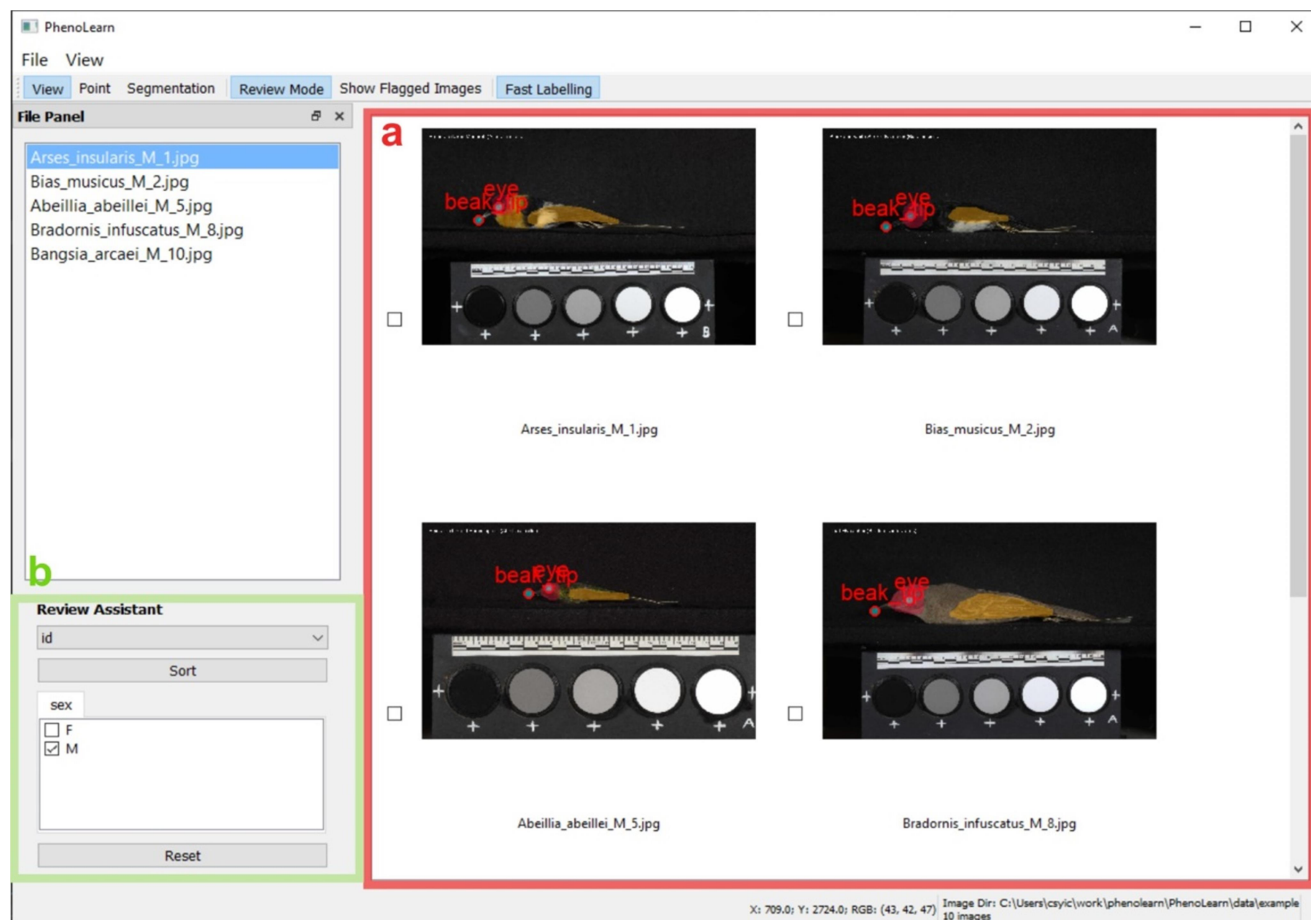


311

312 **Figure 4b**), which can be imported from a property file (**Table 1**). It can sort images by
 313 numerical properties (e.g. specimen length) and filter images by categorical properties
 314 (e.g. taxa). The 'Reset' button clears all filters and sorting.

315

ORIGINAL UNEDITED



316

317 **Figure 4. The PhenoLabel GUI with Review Mode activated.** (a) The Review panel, which
 318 replaces the Main panel, displays image thumbnails with annotations. (b) The Review Assistant.
 319 In this example, it is used to select male specimens and sort images by ID.

320 Examples

321 The examples described below were executed on a Windows 10 system featuring an
 322 Intel(R) Core(TM) i7-11800H CPU, 16 GB of RAM, and an RTX 3080 GPU with 16 GB
 323 of video memory (VRAM). For memory usage results, the highest memory allocation
 324 observed in Task Manager was recorded for CPU usage, while GPU memory usage
 325 was from the output of the nvidia-smi command.

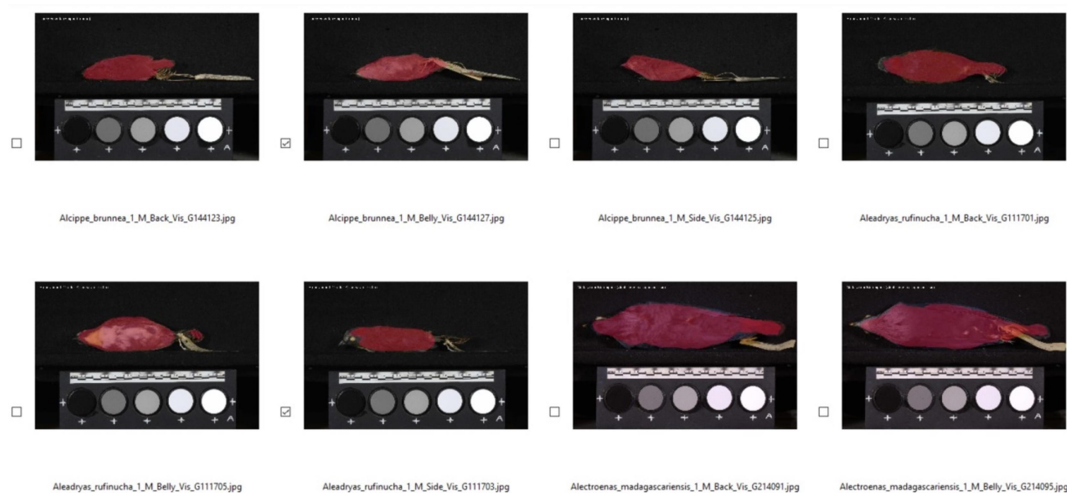
326 Segmenting with PhenoLearn

327 We tested PhenoLearn on a dataset of 220 bird images (4948 x 3280 pixels) to
 328 segment the whole plumage area. We used 120 images for training and the remaining
 329 100 images for prediction. The 120 training images were annotated in PhenoLabel for
 330 training. The DeepLabv3 model was trained for five epochs, with a 20% validation set,
 331 batch size of two, 0.001 learning rate, minimal training level, and an input resolution of
 332 494 x 328 pixels (10% downsampling).

333 The training process was faster with GPU, taking 3 minutes, compared to 13 minutes
 334 without it (CPU only). Predictions were generated in under a minute with GPU and 4
 335 minutes without it. Examples of the predictions can be found in Figure 5. One of the

336 authors (Y.H.) spent five minutes reviewing 100 images. An additional four minutes
 337 were used to correct predictions for these 18 images. In addition, we tested the training
 338 time, GPU usage and performance for using various configurations of GPU and CPU
 339 with different training levels to users with a comprehensive reference. The results are
 340 summarised in Table 2.

341



342

343 **Figure 5.** Examples of the segmentation predictions in the review mode.

344

345 **Table 2.** Training time, memory usage (RAM for CPU and VRAM for GPU), and
 346 performance across different hardware configurations and training levels on the
 347 segmentation test dataset.

<i>Training Level</i>	<i>Hardware</i>	<i>Training Time (Minutes)</i>	<i>Memory Usage (GB)</i>	<i>Average Dice Score</i>
<i>Minimal</i>	GPU	2	1	0.90
	CPU	13	1.2	
<i>Intermediate</i>	GPU	3	3.5	0.94
	CPU	26	4.1	
<i>Full</i>	GPU	3	3.9	0.93
	CPU	30	4.5	

348

349 [Placing points with PhenoLearn](#)

350 We evaluated PhenoLearn on a dataset of 220 *Littorina* images, each measuring 2592 x
 351 1944 pixels, with four points annotated on each image according to a 15-landmark
 352 scheme derived from Ravinet et al. (2016). For this study, 120 images were used for

353 training, while the remaining 100 served for prediction. Annotations for the training
354 images were performed using PhenoLabel.

355 We trained a Mask R-CNN model over five epochs, using a validation set comprising 20%
356 of the data, a batch size of two, a learning rate of 0.001, and an input resolution reduced
357 to 518 x 388 pixels (20% downsampling). We conducted experiments using both GPU
358 and CPU across various training levels. The best performance was an average pixel
359 distance of 21. Details on GPU usage and the performance of different runs can be
360 found in Table 3.

361

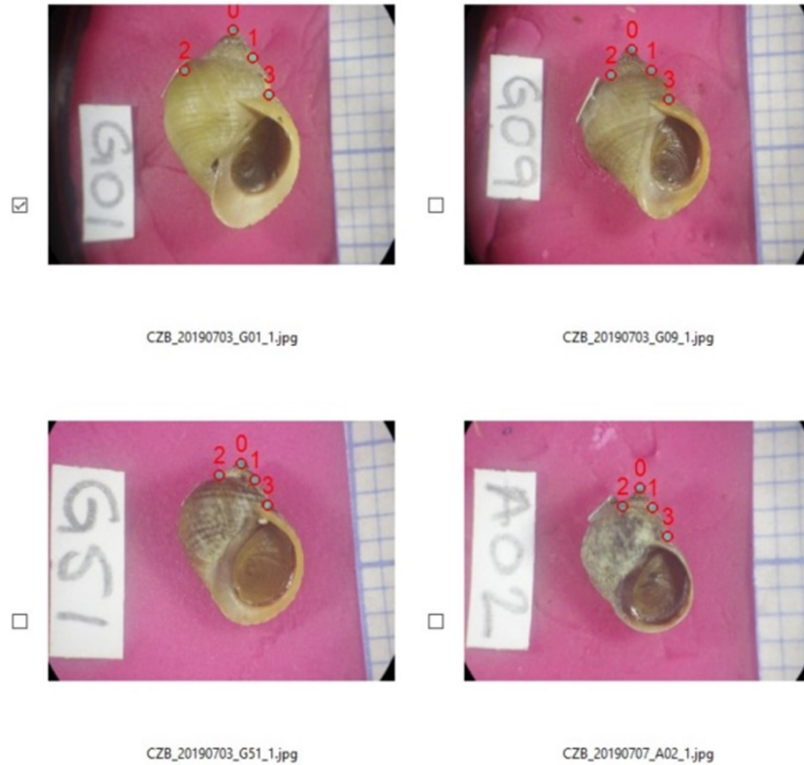
362 **Table 3.** Training time, memory usage (RAM for CPU and VRAM for GPU), and
363 performance across different hardware configurations and training levels on the point
364 test dataset.

<i>Training Level</i>	<i>Hardware</i>	<i>Training (Minutes)</i>	<i>Time</i>	<i>Memory Usage (GB)</i>	<i>Average Distance</i>	<i>Pixel</i>
<i>Minimal</i>	GPU	2		2.7	138	
	CPU	15		1.5		
<i>Intermediate</i>	GPU	2		3	126	
	CPU	25		1.9		
<i>Full</i>	GPU	3		4.5	37	
	CPU	29		3.5		

365

366 Examples of the predictions made using PhenoLearn are illustrated in Figure 6. One of
367 the authors (Y.H.) spent five minutes reviewing 100 *Littorina* images, during which 19
368 images with inaccurately placed points were flagged. An additional four minutes were
369 spent to correct these predictions.

370



371
 372 **Figure 6.** Examples of the point predictions in the review mode.

373
 374 The performance of PhenoTrain can vary with different datasets and training settings.
 375 As shown in the results, training from scratch is not guaranteed to outperform fine-
 376 tuning pre-trained models (see Table 2). The pre-trained models used in PhenoLearn
 377 are based on the ImageNet dataset (Deng et al., 2009), which provides a large and
 378 diverse set of features as a strong starting point. Pre-trained models are also less prone
 379 to overfitting and more capable of generalising to new datasets (Huh et al., 2016;
 380 Yosinski et al., 2014). This advantage makes fine-tuning a pre-trained network a reliable
 381 choice in many scenarios. However, the relative performance of these approaches can
 382 only be determined through testing. Based on our observations, we recommend starting
 383 with fine-tuning for most use cases and minimum computational cost.

384 Another important point is that the randomness inherent in the training process, such as
 385 random weight initialisation and data shuffling during batch creation, can lead to
 386 variability in results. Even with identical configurations and training data, different runs
 387 may yield slightly different outcomes. This variability should be considered when
 388 interpreting results.

389 Here are some other general guidelines:

- 390 • Test model performance with a small subset of your dataset (e.g., 20 images) to
- 391 quickly assess learning progress by monitoring if validation loss decreases and
- 392 the metrics on the validation set are increasing. extend the training to the full
- 393 dataset.

- 394
- 395
- 396
- 397
- 398
- 399
- 400
- 401
- 402
- 403
- 404
- 405
- 406
- 407
- 408
- 409
- 410
- Manage memory (either RAM or video memory) by starting with an input resolution of around 500 x 500 pixels. The resolution can be incrementally increased.
 - Carefully select the learning rate, as it significantly impacts model training. A learning rate that is too large may cause the model to diverge or produce unstable results. For example, using a learning rate of 0.1 on our point dataset caused the loss to become null, resulting in training failure. Conversely, a very small learning rate can result in slow learning and require a large number of epochs to converge. We recommend that users try multiple training runs with different learning rates and monitor performance to find an appropriate setting for their dataset.
 - Better performance may be achieved by increasing the input resolution, training set size, training epochs, and training level. Increasing these settings leads to longer training times. Results from runs with various configurations are provided in the Supplementary Material, where some performance differences can be observed across settings. However, we note that these comparisons are based on a small number of runs and should be interpreted with caution.

411 Users can change these settings to fit their datasets and research requirements.

412 Discussion

413 In summary, PhenoLearn provides a user-friendly, high-throughput data extraction
414 pipeline with fully integrated GUIs, enabling biologists without extensive computational
415 skills to effectively measure phenotypic traits from images. While tools like DeepLabCut
416 and Argos offer robust solutions for specific phenotyping tasks, they focus more deeply
417 on animal tracking, primarily supporting point-based annotations. In contrast,
418 PhenoLearn combines support for point annotations and segmentation tasks within a
419 single toolkit and has already been successfully applied for both annotation types in
420 previously published studies (Cooney et al., 2022; Y. He et al., 2022, 2023).
421 PhenoLearn also includes functions tailored specifically for handling 2D image datasets
422 of natural history collections. These features include 'Fast Labelling,' which streamlines
423 the annotation naming process, and 'Review Mode' and 'Review Assistant,' which
424 leverage specimen metadata to simplify the review process. These capabilities make
425 PhenoLearn particularly suited for natural history collections, which often include rich
426 metadata. Together, these features position PhenoLearn as a complementary tool for
427 phenotyping 2D images, offering unique advantages for researchers working with such
428 datasets.

429 As Lürig (2022) highlights, classic computer vision methods are more accessible to
430 biologists with only CPUs. To facilitate the wider application of deep learning among
431 biologists without GPU access, PhenoLearn leverages pre-trained models and partial
432 model training to shorten CPU training times. Moreover, small training sets can yield
433 accurate predictions for photographs with a highly consistent digitisation set-up, as
434 minimal variation among images may bring more efficient training (Mulqueeney et al.,
435 2024). From our results, it appears that CPU usage requires slightly more memory
436 compared to GPU usage. However, it is more cost-effective to upgrade system RAM
437 than to purchase GPUs with equivalent VRAM capacity. Additionally, most current

438 consumer-grade laptops are equipped with at least 8 GB of RAM, making it feasible for
439 a wide range of researchers to run PhenoLearn effectively on readily available CPU
440 hardware. These features make predicting annotations on digitised specimens possible
441 using only CPUs.

442 The modular design of PhenoLearn, comprising separate modules for image annotation
443 (PhenoLabel) and deep learning (PhenoTrain), offers flexibility to integrate with other
444 tools. This feature is particularly important in the fast-developing field of machine
445 learning, where new and powerful methods are continually being developed such as
446 Segment Anything (Kirillov et al., 2023), the foundation model for semantic
447 segmentation. Thus, with PhenoLearn, users have the option to export annotations from
448 PhenoLabel for other Deep Learning methods, and then use PhenoLabel again for
449 efficient prediction reviewing. PhenoLearn supports multiple output formats (CSV, JSON,
450 and image-based segmentation), making it compatible with other methods or toolkits.
451 These formats can be easily converted into target Python data structures commonly
452 used in deep learning pipelines. For example, regardless of the format, annotations can
453 be transformed into 2D tensors that represent segmentations or point heatmaps, which
454 are among the most used data structures for segmentation and point predictions.
455 PhenoLabel can also simply serve as a manual labelling tool for small datasets.

456 Taken together, PhenoLearn is a versatile toolkit that bridges the gap between
457 biological image datasets and downstream analysis, facilitating greater access for
458 researchers to deep learning tools for image processing and data extraction.

459 **Future Directions**

460 Future development of PhenoLearn will likely focus on four main areas: (1) Optimisation
461 of the user interface based on user feedback to increase usability. (2) Improvement of
462 software performance, such as integrating multi-threading for displaying thumbnails,
463 which will increase the efficiency of the review process. (3) Expansion of supported
464 annotation types based on future user requirements. Adding bounding box annotations,
465 for instance, could significantly broaden the toolkit's applications, including object
466 recognition tasks which can be used to identify specimen appearances in laboratory or
467 camera trap photographs. (4) Integrating alternative and newer models, such as
468 Segment Anything (Kirillov et al., 2023) and other state-of-the-art deep learning models,
469 to further enhance segmentation and landmark prediction capabilities.

470

471 **References**

- 472 Adams, D. C., & Otárola-Castillo, E. (2013). geomorph: An R package for the collection
473 and analysis of geometric morphometric shape data. *Methods in Ecology and*
474 *Evolution*, 4(4), 393–399.
- 475 Blagoderov, V., Kitching, I. J., Livermore, L., Simonsen, T. J., & Smith, V. S. (2012). No
476 specimen left behind: Industrial scale digitization of natural history collections.
477 *ZooKeys*, 209, 133.
- 478 Bradski, G. (2000). The OpenCV Library. *Dr. Dobb's Journal of Software Tools*.
- 479 Chang, J., & Alfaro, M. E. (2016). Crowdsourced geometric morphometrics enable rapid
480 large-scale collection and analysis of phenotypic data. *Methods in Ecology and*
481 *Evolution*, 7(4), 472–482. <https://doi.org/10.1111/2041-210X.12508>
- 482 Chen, L.-C., Papandreou, G., Schroff, F., & Adam, H. (2017). *Rethinking Atrous*
483 *Convolution for Semantic Image Segmentation* (arXiv:1706.05587). arXiv.
484 <http://arxiv.org/abs/1706.05587>
- 485 Cooney, C. R., He, Y., Varley, Z. K., Nouri, L. O., Moody, C. J. A., Jardine, M. D., Liker,
486 A., Székely, T., & Thomas, G. H. (2022). Latitudinal gradients in avian
487 colourfulness. *Nature Ecology & Evolution*, 6(5), 622–629.
488 <https://doi.org/10.1038/s41559-022-01714-1>
- 489 Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., & Fei-Fei, L. (2009). Imagenet: A large-
490 scale hierarchical image database. *2009 IEEE Conference on Computer Vision*
491 *and Pattern Recognition*, 248–255.
- 492 He, K., Gkioxari, G., Dollar, P., & Girshick, R. (2017). Mask R-CNN. *2017 IEEE*
493 *International Conference on Computer Vision (ICCV)*, 2980–2988.
494 <https://doi.org/10.1109/ICCV.2017.322>

- 495 He, Y., Cooney, C. R., Maddock, S., & Thomas, G. H. (2023). Using pose estimation to
496 identify regions and points on natural history specimens. *PLOS Computational*
497 *Biology*, 19(2), e1010933. <https://doi.org/10.1371/journal.pcbi.1010933>
- 498 He, Y., Varley, Z. K., Nouri, L. O., Moody, C. J. A., Jardine, M. D., Maddock, S., Thomas,
499 G. H., & Cooney, C. R. (2022). Deep learning image segmentation reveals
500 patterns of UV reflectance evolution in passerine birds. *Nature Communications*,
501 13(1), 5068. <https://doi.org/10.1038/s41467-022-32586-5>
- 502 Huh, M., Agrawal, P., & Efros, A. A. (2016). *What makes ImageNet good for transfer*
503 *learning?* (arXiv:1608.08614). arXiv. <https://doi.org/10.48550/arXiv.1608.08614>
- 504 John, A., Theobald, E. J., Cristea, N., Tan, A., & Hille Ris Lambers, J. (2024). Using
505 photographs and deep neural networks to understand flowering phenology and
506 diversity in mountain meadows. *Remote Sensing in Ecology and Conservation*,
507 10(4), 480–499. <https://doi.org/10.1002/rse2.382>
- 508 Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T.,
509 Whitehead, S., Berg, A. C., Lo, W.-Y., Dollár, P., & Girshick, R. (2023). *Segment*
510 *Anything* (arXiv:2304.02643). arXiv. <http://arxiv.org/abs/2304.02643>
- 511 Lin, T.-Y., Maire, M., Belongie, S. J., Bourdev, L. D., Girshick, R. B., Hays, J., Perona,
512 P., Ramanan, D., Dollár, P., & Zitnick, C. L. (2014). Microsoft COCO: Common
513 Objects in Context. *CoRR*, abs/1405.0312. <http://arxiv.org/abs/1405.0312>
- 514 Lürig, M. D. (2022). *phenotype*: A phenotyping pipeline for Python. *Methods in Ecology*
515 *and Evolution*, 13(3), 569–576. <https://doi.org/10.1111/2041-210X.13771>

516 Maia, R., Gruson, H., Endler, J. A., & White, T. E. (2019). pavo 2: New tools for the
517 spectral and spatial analysis of colour in r. *Methods in Ecology and Evolution*,
518 10(7), 1097–1107. <https://doi.org/10.1111/2041-210X.13174>

519 Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro,
520 Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat,
521 Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Jia, Y., Rafal
522 Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, ... Xiaoqiang Zheng. (2015).
523 *TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems*.
524 <https://www.tensorflow.org/>

525 Mathis, A., Mamidanna, P., Cury, K. M., Abe, T., Murthy, V. N., Mathis, M. W., & Bethge,
526 M. (2018). DeepLabCut: Markerless pose estimation of user-defined body parts
527 with deep learning. *Nature Neuroscience*, 21(9), 1281–1289.
528 <https://doi.org/10.1038/s41593-018-0209-y>

529 Mulqueeney, J. M., Searle-Barnes, A., Brombacher, A., Sweeney, M., Goswami, A., &
530 Ezard, T. H. G. (2024). How many specimens make a sufficient training set for
531 automated three-dimensional feature extraction? *Royal Society Open Science*,
532 11(6), rsos.240113. <https://doi.org/10.1098/rsos.240113>

533 Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z.,
534 Gimelshein, N., Antiga, L., Desmaison, A., Kopf, A., Yang, E., DeVito, Z., Raison,
535 M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., ... Chintala, S. (2019).
536 PyTorch: An Imperative Style, High-Performance Deep Learning Library. In H.
537 Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, & R. Garnett
538 (Eds.), *Advances in Neural Information Processing Systems 32* (pp. 8024–8035).

539 <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high->
540 [performance-deep-learning-library.pdf](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf)

541 Pennekamp, F., & Schtickzelle, N. (2013). Implementing image analysis in laboratory-
542 based experimental systems for ecology and evolution: A hands-on guide.
543 *Methods in Ecology and Evolution*, 4(5), 483–492.

544 Porto, A., & Voje, K. L. (2020). ML-morph: A fast, accurate and general approach for
545 automated detection and landmarking of biological structures in images. *Methods*
546 *in Ecology and Evolution*, 11(4), 500–512. [https://doi.org/10.1111/2041-](https://doi.org/10.1111/2041-210X.13373)
547 [210X.13373](https://doi.org/10.1111/2041-210X.13373)

548 Ravinet, M., Westram, A., Johannesson, K., Butlin, R., André, C., & Panova, M. (2016).
549 Shared and nonshared genomic divergence in parallel ecotypes of *L. ittorina*
550 *saxatilis* at a local scale. *Molecular Ecology*, 25(1), 287–305.

551 Ray, S., & Stopfer, M. A. (2022). Argos: A toolkit for tracking multiple animals in
552 complex visual environments. *Methods in Ecology and Evolution*, 13(3), 585–595.
553 <https://doi.org/10.1111/2041-210X.13776>

554 Rohlf, F. J. (2006). tpsDig, version 2.10. [Http://Life. Bio. Sunysb. Edu/Morph/Index. Html.](http://Life.Bio.Sunysb.Edu/Morph/Index.Html)

555 Schwartz, S. T., & Alfaro, M. E. (2021). *Sashimi*: A toolkit for facilitating high-throughput
556 organismal image segmentation using deep learning. *Methods in Ecology and*
557 *Evolution*, 12(12), 2341–2354. <https://doi.org/10.1111/2041-210X.13712>

558 Shedrawi, G., Magron, F., Vigga, B., Bosserelle, P., Gislard, S., Halford, A. R., Tiitii, S.,
559 Fepuleai, F., Molai, C., Rota, M., Jalam, S., Fatongiatau, V., Sami, A. P., Nikiari,
560 B., Sokach, A. H. M., Joy, L. A., Li, O., Steenbergen, D. J., & Andrew, N. L.
561 (2024). Leveraging deep learning and computer vision technologies to enhance

562 management of coastal fisheries in the Pacific region. *Scientific Reports*, 14(1),
563 20915. <https://doi.org/10.1038/s41598-024-71763-y>

564 Tan, C., Sun, F., Kong, T., Zhang, W., Yang, C., & Liu, C. (2018). A Survey on Deep
565 Transfer Learning. In V. Kůrková, Y. Manolopoulos, B. Hammer, L. Iliadis, & I.
566 Maglogiannis (Eds.), *Artificial Neural Networks and Machine Learning – ICANN*
567 *2018* (Vol. 11141, pp. 270–279). Springer International Publishing.
568 https://doi.org/10.1007/978-3-030-01424-7_27

569 Yosinski, J., Clune, J., Bengio, Y., & Lipson, H. (2014). How transferable are features in
570 deep neural networks? *Advances in Neural Information Processing Systems*, 27.

571 Zelditch, M. L., Swiderski, D. L., Sheets, H. D., & Fink, W. L. (2004). Geometric
572 morphometrics for biologists: A primer. *Elsevier*, 457.
573 <https://doi.org/10.1016/B978-0-12-386903-6.00001-0>

574
575
576

ORIGINAL UNEDITED MANUSCRIPT