

1 **Interpretable Wildlife Classification by Coupling Genetics, Scoring**  
2 **Systems, and Computer Vision**

3 Sara Gonzalez<sup>a\*</sup>, Andrew Lensen<sup>b</sup> and Philip Lavretsky<sup>a</sup>

4 *<sup>a</sup>Department of Biology, University of Texas at El Paso. 500 W. University, El Paso, TX 79968,*  
5 *United States;*

6 *<sup>b</sup>School of Engineering and Computer Science, Victoria University of Wellington, PO Box 600,*  
7 *Wellington 6140, New Zealand*

8 **\*Corresponding Author:**

9 Sara Gonzalez

10 Department of Biology

11 University of Texas at El Paso

12 500 W. University Ave

13 El Paso, TX 79968

14 United States

15 sgonzalez28@miners.utep.edu

16

17

18 **Abstract**

19 Despite increasing interest in using machine learning to improve field-based species  
20 identification, key challenges remain, particularly limited model interpretability and the large  
21 datasets typically required for reliable performance. Additionally, model accuracy often depends  
22 on expert-labelled training data, which may reflect the same assumptions about trait  
23 diagnosability that constrain traditional methods. We develop interpretable, trait-specific  
24 machine learning protocols designed to perform effectively with relatively small datasets. We  
25 incorporate heatmaps to assess whether poor model performance arises from data limitations,  
26 modeling artifacts, or irrelevant image features. Interpretability is further enhanced through a  
27 numerical scoring system applied to individual diagnostic traits, with final classifications based  
28 on aggregated trait values rather than a single model output. Using transfer learning with a  
29 pretrained ResNet34 architecture, we developed ten neural network models trained on images  
30 representing diagnostic phenotypic traits. Models were trained on photographs of genetically  
31 characterized individuals, reducing reliance on expert identification alone. We applied this  
32 framework to distinguish between two closely related and historically difficult-to-identify  
33 species, the Mexican Duck (*Anas diazi*) and the Mallard (*Anas platyrhynchos*). Trait-based  
34 models achieved balanced accuracies exceeding 90%, accounting for an unbalanced training  
35 dataset. However, when evaluated against expert and non-expert classifications for the same 30  
36 individuals, overall AI performance more closely resembled that of non-experts. Heatmap  
37 analyses revealed inconsistent attention to background features, indicating sensitivity to image  
38 artifacts. This study demonstrates how integrating genomics, interpretable AI, and multi-trait  
39 analyses can reduce bias in phenotypic identification and improve classification performance  
40 with limited datasets.

41 **Keywords:** Artificial Intelligence; Computer Vision; Convolutional Neural Networks; Binary  
42 Classification; Image Segmentation; Wildlife Conservation; Waterfowl

Preprint not peer reviewed

## 43 **1. Introduction**

44 Application of artificial intelligence (AI) in biology has become more accessible by harnessing  
45 the growing power of computer vision and machine learning with technological advances in  
46 Graphic Processing Units (GPUs) (Liu & Du, 2024). Generally, machine learning is the process  
47 by which computers conduct pattern recognition through specified algorithms on provided data  
48 and is a tool that has been widely used throughout research practices, including towards species  
49 identification (Nadimpalli, 2005, Brenowitz & Larson, 2015, Atanbori et al, 2016). For wildlife,  
50 machine learning has been used to identify individuals from data ranging from sound pattern  
51 recognition (Brandes, 2008; Incze et al, 2018; Mehyadin et al, 2021) to color histograms (Marini  
52 et al, 2013; Miao et al, 2019; Moallem, 2021). These studies demonstrate the ability to not only  
53 automate the process but help avoid individual-based biases that often arise (Mendes et al, 2024).  
54 However, machine learning still requires sufficient and high-quality data vetted by experts in the  
55 field (Picard et al, 2020; Gong et al, 2023). While such a model framework has been well  
56 integrated into large species databases like Merlin Bird (i.e., ebird; Chu, 2012). and iNaturalist  
57 (Ueda & Loarie, 2008), these remain contingent on those expert's initial identification and  
58 constant improvements based on the processing of hundreds of thousands of images constantly  
59 being fed by end users. Consequently, such applications are not often feasible for most  
60 researchers and their study systems (Farley et al, 2018; Hampton et al, 2013). Furthermore, the  
61 identification of phenotypically similar species can be highly complex even for the most trained  
62 individuals (Müller et al, 2007), which becomes even more convoluted if such species make  
63 viable hybrid offspring in the wild. Here, we attempt to circumvent such limitations by  
64 integrating methods and processes with datasets of genetically rather than expertise verified  
65 individuals into a machine learning algorithm. In addition to understanding the potential for

66 identifying parental and hybrid individuals, we assess any limitations of such pipelines for  
67 datasets of moderate sizes that are more realistic for most wildlife systems. Importantly, unlike  
68 many previous approaches that rely on widely available public image sets (Man & Chahl, 2022),  
69 this study employs a newly assembled training dataset in which all images are molecularly  
70 verified. Doing so has the potential to remove the need for trained experts to verify and remove  
71 potentially unknown biases when solely dealing with phenotypic traits (Adolph & Hardin, 2007;  
72 Van der Sluis et al, 2010). In addition, the training dataset captures a broad range of phenotypic  
73 variation within each species, minimizing the risk that unaccounted plumage variation biases  
74 model performance. Although AI has been used to study genic-disease correlations in human  
75 diseases (Duong et al, 2022; Reyes & Sanchez, 2024), this is the first attempt to determine trait-  
76 genetic associations from pictures using artificial intelligence and computer vision.

77 Here, we aim to develop methods that standardize species identification with measurable  
78 confidence for non-experts, while accounting for data limitations. Critically, this study  
79 contributes to both ecological and AI research by emphasizing model interpretability by  
80 designing architectures that mimic expert decision-making processes and applying visual  
81 explainability tools such as activation heat maps to reveal regions influencing model predictions  
82 within images. Together, these efforts move beyond simple classification accuracy toward a  
83 biologically meaningful understanding of how models “see,” ensuring that automated systems  
84 make ecologically valid classifications rather than relying on spurious cues. By integrating  
85 interpretability and biological reasoning, this work advances the development of transparent,  
86 reliable, and scalable AI approaches for ecological trait detection.

87 To build interpretable classification models, we deployed convolutional neural networks  
88 (CNNs), which identify patterns in images through layered filtering operations (LeCun et al.,

89 2015). Specifically, we implemented a CNN architecture based on a ResNet34 backbone trained  
90 through transfer learning, allowing the model to leverage pre-trained visual features while fine-  
91 tuning to recognize biologically relevant plumage traits. Although deeper networks tend to yield  
92 more robust performance, their increasing complexity also reduces interpretability, producing  
93 hyperdimensional representations that are difficult to link to human logic and ultimately lead to  
94 the “black box” effect (Murdoch et al., 2019). To mitigate this, we constrained the model  
95 architecture to mimic the process biologists use when identifying species by evaluating the  
96 presence or absence of individual diagnostic traits within specialized models. This approach  
97 allows causes of model decisions to be extracted more transparently. We further enhanced  
98 interpretability using heat maps of the final convolutional layer to visualize which image regions  
99 drive classification decisions (Lipton, 2018; Lucas, 2020), enabling an assessment of whether the  
100 model relied on biologically relevant cues or on spurious image artifacts. Although constraining  
101 the model in this way may reduce its maximal predictive power, we considered this an  
102 appropriate tradeoff as it substantially improves interpretability and allows us to evaluate the  
103 biological validity of the model’s decisions.

104 We applied these models to a case study involving two closely related duck species, the Mexican  
105 Duck (*Anas diazi*) and Mallard (*Anas platyrhynchos*). Both belong to the Mallard Complex, a  
106 radiation of 14 species that diverged over the last two million years (Lavretsky, 2021), with  
107 Mexican Ducks and Mallards representing the most recent split (~500,000 years; Lavretsky et  
108 al., 2015; Brown et al., 2022a), and population genomic studies have repeatedly shown that  
109 many phenotypic traits used in traditional field keys fail to reliably distinguish parental  
110 individuals from hybrids (Lavretsky et al., 2019; Lavretsky et al., 2021; Brown et al., 2022b). For  
111 example, keys designed by Kirby et al. (2000) correctly classified only 80% of parental birds and

112 60% of hybrids (Lavretsky et al., 2019). Because misidentification can introduce substantial bias  
113 into management metrics such as sex-age structure estimates (McCartney et al., 2019; Perry et  
114 al., 2002), improved diagnostic approaches are needed. Recent genomic-first trait analyses have  
115 demonstrated >98% accuracy when distinguishing sex-age cohorts of Mexican Ducks, Mallards  
116 and their hybrids (Brown et al., 2022b). We build directly on this work by training models using  
117 photographs of genetically vetted individuals, enabling classification of parental species with  
118 hybrids without reliance on phenotype-driven expert judgement. This approach also allows us to  
119 evaluate the capabilities and limitations of model AI-based methods independently of observer  
120 bias and to determine whether classification challenges arise from phenotypic similarity, data  
121 quality, or algorithmic constraints.

## 122 **2. Methods**

123 A total of 11 phenotypic traits were considered across images (Table 1). Although initial  
124 testing used whole-bird photographs, we found that the potential of computer vision to score traits  
125 increased by partitioning each trait by body region including, the dorsal and ventral sides of the  
126 head, body, and wings (Figure 1). This ensured that the neural network was trained on images that  
127 are directly relevant to each trait, reducing the likelihood of misclassification due to irrelevant  
128 plumage or background features. For example, classification accuracy for detecting green  
129 coloration on the head increased markedly when using head-specific partitions rather than full-  
130 body images (Figure 2).

131 Once corresponding trait folders were prepared, we performed data augmentation using  
132 the Python Albumentations package (Buslaev et al, 2020), applying randomized rotation,  
133 transpositions, crops, and small perturbations in hue, saturation and value while remaining within  
134 the RGB color space. Alternate color spaces (e.g. HSV) were not used at this stage due as RGB is

135 the native color format of the photographic dataset and is the default input representation for  
136 ResNet-based architectures, and thus, models were trained directly on RGB tensors. This avoided  
137 unnecessary color-space transformations while maintaining compatibility with established  
138 computer-vision pipelines. These alterations increased the effective size of training sets. Moreover,  
139 using models like ResNet34 (He et al, 2015) that is pretrained on approximately 1.3 million  
140 ImageNet images reduced the total amount of training required through transfer learning. During  
141 training, we monitored performance using trained and validation loss, error rate, balanced accuracy  
142 and F1 score. Balanced accuracy was used due to an uneven distribution in quantity of images  
143 between Mexican Ducks (N = 2,467) and Mallards (N = 834) (Brodersen et al, 2010), while the  
144 F1 score used to determine accuracy using precision and recall in machine learning (Sørensen,  
145 1948).

146 After training, transfer learning was done using the *fast.ai* package in which the full data  
147 pipeline takes the augmented images and separates 20% of them for validation. Following *fast.ai*'s  
148 protocol, a learning rate was determined for each learning run (Howard & Gugger 2020), which  
149 ensured that the learning rate wasn't too slow or fast that could cause the model to take too long  
150 to learn or the model to overshoot and likewise fail to train properly, respectively. The model is  
151 first trained for 50 cycles (called 'epochs') using the 1Cycle Policy (Smith & Topin, 2019). The  
152 1Cycle Policy is a method designed to accelerate neural network training by dynamically adjusting  
153 the learning rate. An early stopping criterion was applied to halt training if no improvement in  
154 training or validation loss was observed for three consecutive epochs, helping to prevent  
155 overfitting. During this phase, all layers of the neural network except the last two were frozen,  
156 allowing only these last layers to be updated. This approach leverages the fact that earlier  
157 convolutional layers typically capture low-level features such as edges and textures, while deeper

158 layers encode more complex visual patterns. These weights were saved as “st1” (step one) once  
159 the training completed. Next, the model is completely unfrozen, and a new learning rate is found;  
160 once again, the network is trained for 50 epochs (with the same early-stop criterion) with the new  
161 learning rate though the rate is severely reduced for the topmost layers as these are already trained  
162 on finding basic shapes and patterns that do not need to be re-learned. This second pass is saved  
163 separately as “unfrozen.” Note all model training was done with a RTX 4000 ADA Generation  
164 Graphics Processing Unit (GPU), on a Windows Subsystem for Linux (WSL) using an Ubuntu  
165 24.04.2 distribution. Our pipeline also readies the transfer learning architecture to have the desired  
166 amount of output nodes at the final layer.

167         With a binary classification system, our final layers had two possible outputs as either  
168 Mexican Duck or Mallard. In short, once training and learning was complete across models, results  
169 were aggregated into a dictionary in python where they were assigned to their respective body part  
170 during testing (Table 1, Figure 2). In the end, the developed application assigned binary numerical  
171 values across traits for any new specimen for which images that were not used as a training set.  
172 For each trait, the most frequently occurring value (mode) across all images in the folder is  
173 calculated. Mode values are stored in a Python dictionary as one entry per trait. The final output  
174 includes the mode for each model and a sex prediction. Thus, the species identity is determined by  
175 summing the results for every model except the Sex classifier, with a minimum score of 0 and a  
176 maximum score of 9. If the final score is  $\leq 3$  or  $\geq 7$ , the duck would be classified as a Mexican  
177 Duck or Mallard, respectively, with in-between scores labeled as Unknown. These thresholds were  
178 chosen to represent conservative boundaries that minimize misclassification, recognizing that  
179 individuals with intermediate scores likely reflect either hybrid ancestry or ambiguous phenotypic  
180 expression. Finally, sex was determined using the bill color model with an output of 0 or 1

181 indicating female or male, respectively. Overall, we evaluated classification performance at the  
182 individual level by combining predictions from the 10 trait-specific classifiers post-hoc using the  
183 consensus scoring procedure as a sum across traits and using the class with the highest cumulative  
184 score to assign the final label.

185 Finally, to calculate the overall confidence of the classification models, we used the raw  
186 prediction outputs (logits) generated by running each image through the trained *fast.ai* visual  
187 learner models saved as *.pkl* files. These logits were converted into class probabilities using the  
188 softmax function, which normalizes the outputs so that the sum across all classes equals one. For  
189 each prediction, we determined the predicted class by applying the argmax function that selects  
190 the index of the class with the highest probability of classification. The probability corresponding  
191 to the predicted class was then extracted and treated as a confidence score for that prediction.  
192 Because multiple predictions were made for each model, a list of confidence scores was collected  
193 for each model. From these lists, the highest confidence value per model was selected to represent  
194 the maximum certainty achieved by that model. To obtain a single measure of confidence across  
195 all models, the maximum confidence values from each model were averaged. This average of per-  
196 model maximum confidences was then expressed as a percentage (a.k.a., overall accuracy) to  
197 provide an interpretable summary of the overall model confidence.

198 Finally, we attempted to test model performance by having three non-experts and two  
199 experienced field biologists independently classify the same 30 images evaluated by the CNN for  
200 species, sex, and full identification (species + sex). Classification accuracies were recorded to  
201 contextualize the CNN's performance relative to expert identification.

202 *2.1 Heatmaps and trait validation*

203           Once models were trained, a set of nine randomly chosen images of the same body region  
204 that the trait is testing on are printed to show the prediction and true label for that image (Figure  
205 3). Doing so provides clarity into whether the computer vision is correctly scanning the trait of  
206 interest, as well troubleshoot when scores and predictions are inaccurate (Figure 1). To do this, we  
207 generate an occlusion map using the Captum's IntegratedGradients package (Kokhlikyan et al,  
208 2020) that was specifically developed for visualization of ResNet models. This occlusion map  
209 takes a test picture and creates a heatmap of where in the image the model focuses most for its  
210 prediction (Figure 2), and where the darker color on the image demarcated highest model focus.  
211 Finally, a confusion matrix is generated to inspect the amount of type 1 and type 2 error(s). This  
212 protocol is repeated for each individual trait, using the previously reported body region folder as  
213 the training input for each trait model.

### 214 **3. Results**

215           Each trait model was trained and validated on a unique body-specific images of genetically  
216 known Mexican Ducks or Mallards ( $N_{\text{Mallard}} = 1268$  images;  $N_{\text{Mexican Duck}} = 2824$  images) with an  
217 80%-20% train-validation split (362 – 958 images for training, 74 – 239 for validation; see  
218 Supplementary Table S1 for per-body region counts). Across the 10 trait-specific classifiers,  
219 balanced accuracies during internal validation ranged from 78% to 100% (Table 2), where  
220 balanced accuracy evaluates the network's performance when dealing with imbalanced datasets,  
221 with most models performing extremely well. Models trained on high-contrast or structurally  
222 distinct traits achieved >97% balanced accuracy. For example, Speculum Color model reliably  
223 classified the presence of iridescent green or blue (balanced accuracy = 100%, see Figure 3A), and  
224 the Overall Breast and Belly plumage model effectively distinguished between solid and mottled  
225 patterns (balanced accuracy = 100%, see Figure 3B). In contrast, models trained on presence of

226 green on the head and Greater Secondary Covert Patterns produced lower balanced accuracies  
227 (range: 78–92%; see Table 2 and Supplementary Table S2). Finally, the sex classification model  
228 achieved a balanced accuracy of 93%, but sex was not used as part of the final score calculation  
229 for species classification.

230 To evaluate generalization beyond the training dataset, the complete set of trait models  
231 were validated with 30 new genetically-known individuals that were entirely unseen during  
232 training. When combined to classify species for these individuals, the integrated trait-based system  
233 achieved an overall balanced accuracy of 72% with precision, recall, and F1 scores of 50%, 67%,  
234 and 57% respectively. The model resulted in a full identification accuracy (where both sex and  
235 species were correct) of 40.00%, sex accuracy of 53.33% and species accuracy of 66.67% (Table  
236 3). AI model accuracies were closer to non-expert people achieved 43-56% accuracy in full  
237 identification and 66.7-76.7% species accuracy as compared to experts who reached accuracies of  
238 86.7% full identification, 90% for sex, and 86% for species (Table 4, Supplementary Tables S3-  
239 S8).

240 Next, we applied IntegratedGradients attribution heatmaps (Kokhlikyan et al., 2020) across  
241 trait models to understand regions of model activity. Although, heatmaps consistently highlighted  
242 biologically relevant regions in several models (e.g., sex, belly models; Figure 3A), others showed  
243 activation outside the primary trait region (e.g., the speculum model also included adjacent primary  
244 feathers; Figure 3B). For example, whereas the sex model concentrated activation on the bill and  
245 head plumage as desired (Figure 3C), background or non-biological elements (e.g., hands, tags,  
246 shoes, flooring) were detected among the top attribution regions in six (of 10) trait-specific  
247 heatmaps and despite image segmentation (Supplementary Figure S2).

#### 248 **4. Discussion**

249 The application of convolutional neural networks towards species conservation is  
250 increasingly inevitable as ecological research shifts towards data-intensive, image-based  
251 monitoring (Lamba et al, 2019; Ullah et al, 2025). The exponential growth of digital imagery  
252 from community science platforms, camera traps, and museum collections provides a natural  
253 foundation for machine learning applications, while the need for rapid, scalable, and objective  
254 species identification makes deep learning approaches particularly well suited for conservation  
255 science. However, appropriate dataset size and accuracy, as well as model development, remain  
256 key limitations before their full potential can be realized. Existing tools such as Merlin &  
257 iNaturalist rely on extensive visual datasets requiring substantial expert annotation, an effort  
258 often infeasible for rare or cryptic species, and may embed unknown observer biases. Reported  
259 classification accuracies in generalist ecological CNN systems vary widely, often ranging from  
260 70-97% for species-level identification when trained on very large and visually redundant  
261 datasets (Binta Islam et al., 2023; Norouzzadeh et al., 2018). More specialized bird-identification  
262 models trained only on external phenotypes similarly achieve high accuracy primarily when  
263 dataset sizes exceed tens or hundreds of thousands of labeled images. In comparison, the  
264 individual model balanced accuracies obtained here (>90%) using a genetically vetted dataset of  
265 only 3,301 images demonstrate that biologically grounded training labels can substantially  
266 compensate for smaller sample sizes. Although direct comparison is imperfect, given that studies  
267 do not generally incorporate underlying genetic structure, the results indicate competitive  
268 performance relative to large-scale visual-only systems while offering greater taxonomic rigor.  
269 In most applications, model optimization focuses primarily on expanding dataset size to  
270 overcome variability within smaller samples. Still, interpretability-based approaches offer an

271 alternative path, reducing the need for such massive datasets by identifying how and where  
272 models visualize patterns (Yosinski et al, 2015).

273 Here we introduce a novel and interpretable CNN framework specifically designed for  
274 ecological application. Rather than training models on visually identified individuals, we use  
275 genetically verified images as the foundation for trait recognition, ensuring that learned features  
276 reflect true biological difference rather than superficial phenotypic variability. This approach  
277 directly addresses a major limitation in previous species identification efforts that rely solely on  
278 external appearance which is particularly problematic for closely related species or hybrids  
279 where plumage traits can be ambiguous (Brown et al., 2022b). By pairing genetically partitioned  
280 datasets with trait-specific CNN classifiers, we demonstrate that high accuracy (balanced  
281 accuracies >90%; Table 2) comparable to expert-based identification can be achieved even with  
282 limited dataset of 3,301 images. When compared to human ability, although the CNN does not  
283 yet reach expert-level performance, it matches or exceeds non-expert capabilities and thus  
284 provides a promising baseline for automated identification – particularly in scenarios where  
285 skilled field experts are unavailable (Table 4).

286 Among improvements, the application of heatmaps proved to identify innate challenges  
287 during trait learning by the simple nature of image interpretability among convolutional neural  
288 networks (Figure 3). Models struggled to focus on specific traits even when individuals were  
289 partitioned by parts of the birds first (Figure 2), with models often marking background or  
290 foreground objects unrelated to the bird of interest. For example, we found that due to many  
291 Mexican Duck images taken with a white background surrounding the duck caused test images  
292 with white background to be identified as a Mexican Duck with high confidence regardless of  
293 subject identity. Understanding this, we was able to include steps for further image partitioning

294 through background segmentation and removal to circumvent this issue. Finally, although the  
295 multi-network approach undertaken to emphasize interpretability may limit performance  
296 achieved from a single, standard CNN, we argue that the tradeoff is necessary when working  
297 with relatively small to moderate datasets. Specifically, application of outlined data  
298 augmentation steps of cropping, hue and value shifts, and rotations to our modest dataset of  
299 2,824 and 1,268 Mexican Duck and Mallard images, respectively, proved to overcome the large  
300 datasets often required for training of CNNs (Uchida et al, 2016); however, we acknowledge that  
301 while predictive accuracies were promising, further improvements is expected with larger, more  
302 diverse datasets.

303 Beyond qualitative interpretability, formal analysis of model error and uncertainty  
304 represents an important opportunity for future optimization. Misclassification patterns in this  
305 study, such as systematic confusion driven by background artifacts, suggest that probabilistic  
306 outputs (i.e. the softmax-derived confidence scores for each class) could be leveraged to identify  
307 image regions or trait categories that require targeted augmentation, improved segmentation, or  
308 rebalancing. Incorporating uncertainty-aware methods, such as Monte Carlo dropout  
309 (Asgharnezhad et al., 2025) or deep ensembles (Zhou et al., 2024), would enable quantification  
310 of prediction confidence and help distinguish between errors driven by insufficient data versus  
311 genuinely ambiguous phenotypes. Such approaches could guide iterative dataset expansion by  
312 highlighting which regions of phenotype space require additional genetically verified training  
313 images.

#### 314 *4.1 Application in Wildlife Conservation & Future Work*

315 Wildlife conservation is increasingly challenged by accelerating biodiversity loss, habitat  
316 fragmentation, and the growing need for efficient monitoring at scale. Traditional field methods

317 remain invaluable but are often labor intensive, time consuming, and limited in geographic scope.  
318 Consequently, artificial intelligence is being recognized as a transformative tool in wildlife  
319 conservation for species identification (Sharma et al, 2022; Soni et al, 2023) and population  
320 monitoring (Hedge et al, 2024; Kumar & Jakhar, 2022). Generally, field key development can  
321 achieve sufficiently high accuracies, but depend on manual measurements and notes that require  
322 some expertise, and thus, remain limiting for more rapid and large-scale deployment (Brown et al,  
323 2022b). Rather, our approach no longer requires field expertise by coupling a field key-like scoring  
324 system with computer vision that enables near-instantaneous identification of individuals from  
325 photos; though the initial scoring system that the models are trained on still require expertise to  
326 ensure high-quality training. Importantly, validation of accuracy is not limited to model scores,  
327 but the inclusion of genetic information provides a secondary validation step. In short, building  
328 models based on genetically vetted individuals not only ensures that any nuance of phenotypic  
329 variability due to hybridization does not further complicate models up front, but also will  
330 eventually permit for a scoring system to output genetic probabilities. However, hybrid  
331 classification remains outside the present scope of this study, as current trait-based classifiers are  
332 not yet sufficiently robust to reliably identify mixed and intermediate phenotypes.

333 We acknowledge that further training of our models with additional individuals will be  
334 required to increase accuracy scores, as well as training models with longer epochs on dedicated  
335 GPU hardware (which is currently out of our budget's scope) will further improve classification  
336 accuracy. Nevertheless, our approach can be extrapolated to any species that uses trait-specific  
337 classifiers and a moderate, but high-quality dataset of images to train with that can facilitate real-  
338 time species differentiation in the field. Moreover, developing a system grounded in genetically  
339 confirmed data would elevate the reliability of training datasets, ensuring that ground-truth labels

340 are accurate rather than reliant on potentially ambiguous visual identification. In addition,  
341 establishing a high-quality reference dataset of individuals with known genetic lineages would  
342 allow for true categorical classification of parental taxa, providing stronger alignment between  
343 genetic and image-based datasets. Next, integrating advanced image segmentation approaches  
344 (i.e., Mask R-CNN; He et al, 2017) with a ResNet backbone can help isolate relevant bird features  
345 prior to classification. Doing so would minimize the influence of background noise and sharpen  
346 model attention on diagnostically important traits. Model architecture also presents opportunities  
347 for performance gains. Originally developed for natural language processing, transformer-based  
348 models are now used in computer vision operate by modeling long-range relationships across an  
349 image through self-attention mechanisms (Khan et al., 2022), enabling them to capture subtle,  
350 spatially distributed features that standard CNNs may overlook. Because their architecture excels  
351 at fine-grained pattern recognition, models such as Vision Transformers (ViT) (Dosovitskiy et al.,  
352 2021) and Swin Transformers (Liu et al., 2021), alongside advanced CNN architectures like  
353 EfficientNet (Tan & Le, 2020), may be particularly well suited for the trait-specific and fine-scale  
354 discrimination needed for species identification in similar-looking taxa. Their ability to extract  
355 multi-scale features and integrate information across the entire image suggests strong potential for  
356 improving classification of intermediate, ambiguous or hybrid phenotypes. Similarly, lightweight  
357 architectures such as MobileNetV3 (Howard et al., 2019) and ConvNeXt (Liu et al., 2022) could  
358 support deployment in field applications with limited hardware while maintaining high accuracy.  
359 Future work should evaluate these architectures using the genetically confirmed dataset established  
360 here to determine whether more expressive backbones can further close the performance gap with  
361 expert observers. Finally, the potential of extracting numerical values for traits for individuals with

362 genetic information allows the possibility of understanding complex trait evolution and expression  
363 through genome-wide association studies (Sella & Barton, 2019).

Preprint not peer reviewed

365 **Literature Cited**

366

367 Adolph, S. C., & Hardin, J. S. (2007). Estimating phenotypic correlations: correcting for bias due

368 to intraindividual variability. *Functional Ecology*, 178-184.

369 Asgharnezhad, H., Shamsi, A., Alizadehsani, R., Mohammadi, A., & Alinejad-Rokny, H. (2025).

370 *Enhancing Monte Carlo Dropout Performance for Uncertainty Quantification*

371 (arXiv:2505.15671). arXiv. <https://doi.org/10.48550/arXiv.2505.15671>

372 Atanbori, J., Duan, W., Murray, J., Appiah, K., and Dickinson, P (2016). “Automatic

373 classification of flying bird species using computer vision techniques,” *Pattern Recognit.*

374 *Lett.*, vol. 81, pp. 53–62, doi: 10.1016/j.patrec.2015.08.015.

375 Binta Islam, S., Valles, D., Hibbitts, T. J., Ryberg, W. A., Walkup, D. K., & Forstner, M. R. J.

376 (2023). Animal Species Recognition with Deep Convolutional Neural Networks from

377 Ecological Camera Trap Images. *Animals*, 13(9), 1526.

378 <https://doi.org/10.3390/ani13091526>

379 Brandes, T (2008). Automated sound recording and analysis techniques for bird surveys and

380 conservation. *Bird Conservation International*, 18(S1), S163-S173.

381 Brenowitz, E. A., and Larson, T. A. “Neurogenesis in the adult avian song-control system.” *Cold*

382 *Spring Harbor perspectives in biology* vol. 7,6 a019000. 1 Jun. 2015,

383 doi:10.1101/cshperspect.a019000

384 Brodersen, K. H., Ong, C. S., Stephan, K. E., & Buhmann, J. M. (2010, August). The balanced

385 accuracy and its posterior distribution. In 2010 20th international conference on pattern

386 recognition (pp. 3121-3124). *IEEE*.

387 Brown, J. I., Harrigan, R. J., & Lavretsky, P. (2022). Evolutionary and ecological drivers of local  
388 adaptation and speciation in a North American avian species complex. *Molecular*  
389 *Ecology*, 31(9), 2578-2593.

390 Brown, J. I., Hernández, F., Engilis Jr, A., Hernández-Baños, B. E., Collins, D., & Lavretsky, P  
391 (2022). Genomic and morphological data shed light on the complexities of shared  
392 ancestry between closely related duck species. *Scientific reports*, 12(1), 10212.

393 Buslaev, A., Iglovikov, V. I., Khvedchenya, E., Parinov, A., Druzhinin, M., & Kalinin, A. A  
394 (2020). Alumentations: fast and flexible image augmentations. *Information*, 11(2), 125.

395 Chu, M (2012). Merlin: Online bird identification with human learning and machine learning.

396 Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T.,  
397 Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., & Houlsby, N. (2021).  
398 *An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale*  
399 (arXiv:2010.11929). arXiv. <https://doi.org/10.48550/arXiv.2010.11929>

400 Duong, D., Waikel, R. L., Hu, P., Tekendo-Ngongang, C., & Solomon, B. D (2022). Neural  
401 network classifiers for images of genetic conditions with cutaneous  
402 manifestations. *Human Genetics and Genomics Advances*, 3(1).

403 Farley, S. S., Dawson, A., Goring, S. J., & Williams, J. W. (2018). Situating ecology as a big-  
404 data science: Current advances, challenges, and solutions. *BioScience*, 68(8), 563-576.

405 Gong, Y., Liu, G., Xue, Y., Li, R., & Meng, L (2023). A survey on dataset quality in machine  
406 learning. *Information and Software Technology*, 162, 107268.

407 Hampton, S. E., Strasser, C. A., Tewksbury, J. J., Gram, W. K., Budden, A. E., Batcheller, A. L.,  
408 Duke, C. S., & Porter, J. H. (2013). Big data and the future of ecology. *Frontiers in*  
409 *Ecology and the Environment*, 11(3), 156–162.

410 He, K., Gkioxari, G., Dollár, P., & Girshick, R. (2017). Mask R-CNN. In Proceedings of the  
411 *IEEE International Conference on Computer Vision* (pp. 2961-2969).

412 He, K., Zhang, X., Ren, S., & Sun, K (2015). Deep Residual Learning for Image Recognition.  
413 *arXiv preprint arXiv:1512.03385*.

414 Hegde, N. R., & Bargavi, D. S. (2024). Smart Conservation: Integrating AI for enhanced wildlife  
415 monitoring. *International Research Journal of Computer Science*.

416 Howard, A., Sandler, M., Chu, G., Chen, L.-C., Chen, B., Tan, M., Wang, W., Zhu, Y., Pang, R.,  
417 Vasudevan, V., Le, Q. V., & Adam, H. (2019). *Searching for MobileNetV3*  
418 (arXiv:1905.02244)

419 Howard, J., & Gugger, S. (2020). Fastai: a layered API for deep learning. *Information*, 11(2),  
420 108.

421 Incze, A., Jancsó, H. B., Szilágyi, Z., Farkas, A., & Sulyok, C. (2018). Bird sound recognition  
422 using a convolutional neural network. In 2018 IEEE 16th international symposium on  
423 intelligent systems and informatics (SISY) (pp. 000295-000300). IEEE.

424 Khan, S., Naseer, M., Hayat, M., Zamir, S. W., Khan, F. S., & Shah, M. (2022). Transformers in  
425 Vision: A Survey. *ACM Computing Surveys*, 54(10s), 1–41.  
426 <https://doi.org/10.1145/3505244>

427 Kirby, R. E., Reed, A., Dupuis, P., Obrecht, H. H. III, & Quist, W. J. (2000). Description and  
428 identification of American black duck, mallard, and hybrid wing plumage (p. 26).  
429 *Biological Science Report. U.S. Geological Survey, Biological Resources Division*.

430 Kokhlikyan, N., Miglani, V., Martin, M., Wang, E., Alsallakh, B., Reynolds, J., Melnikov, A.,  
431 Kliushkina, N., Araya, C., Yan, S., & Reblitz-Richardson, O. (2020). *Captum: A unified*  
432 *and generic model interpretability library for PyTorch* (arXiv:2009.07896).

433 Kumar, D., & Jakhar, S. D. (2022). Artificial intelligence in animal surveillance and  
434 conservation. *Impact of artificial intelligence on organizational transformation*, 73-85.

435 Lamba, A., Cassey, P., Segaran, R. R., & Koh, L. P. (2019). Deep learning for environmental  
436 conservation. *Current Biology*, 29(19), R977-R982.

437 Lavretsky, P (2021). Population genomics provides key insights into admixture, speciation, and  
438 evolution of closely related ducks of the mallard complex. *Population Genomics:  
439 Wildlife*, 295-330.

440 Lavretsky, P., DaCosta, J. M., Hernández-Baños, B. E., Engilis, A., Sorenson, M. D., & Peters, J.  
441 L. (2015). Speciation genomics and a role for the Z chromosome in the early stages of  
442 divergence between Mexican ducks and mallards. *Molecular Ecology*, 24(21), 5364–  
443 5378. <https://doi.org/10.1111/mec.13402>

444 Lavretsky, P., DaCosta, J. M., Sorenson, M. D., McCracken, K. G., & Peters, J. L (2019).  
445 ddRAD-seq data reveal significant genome-wide population structure and divergent  
446 genomic regions that distinguish the mallard and close relatives in North  
447 America. *Molecular Ecology*, 28(10), 2594-2609.

448 Lavretsky, P., Duenez, E., Kneece, M., & Kaminski, R. M (2021). Population genetics of a  
449 translocated population of Mottled Ducks and allies. *The Journal of Wildlife  
450 Management*, 85(8), 1616-1627.

451 LeCun, Y., Bengio, Y., & Hinton, G (2015). Deep learning. *Nature*, 521(7553), 436–444.  
452 <https://doi.org/10.1038/nature14539>

453 Lipton, Z. C (2018). The mythos of model interpretability: In machine learning, the concept of  
454 interpretability is both important and slippery. *Queue*, 16(3), 31-57.

455 Liu, L., & Du, K. (2024). A perspective on computer vision in biosensing. *Biomicrofluidics*,  
456 18(1), 011301. <https://doi.org/10.1063/5.0185732>

457 Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., & Guo, B. (2021). *Swin*  
458 *Transformer: Hierarchical Vision Transformer using Shifted Windows*  
459 (arXiv:2103.14030).

460 Liu, Z., Mao, H., Wu, C.-Y., Feichtenhofer, C., Darrell, T., & Xie, S. (2022). *A ConvNet for the*  
461 *2020s* (arXiv:2201.03545).

462 Lucas, T. C (2020). A translucent box: interpretable machine learning in ecology. *Ecological*  
463 *Monographs*, 90(4), e01422.

464 Man, K., & Chahl, J (2022). A review of synthetic image data and its use in computer  
465 vision. *Journal of Imaging*, 8(11), 310.

466 Marini, A., Facon, J., and Koerich, A. L (2013). “Bird Species Classification Based on Color  
467 Features,” *2013 IEEE International Conference on Systems, Man, and Cybernetics*, pp.  
468 4336–4341, doi: 10.1109/SMC.2013.740.

469 McCartney, K. R., Kumar, S., Sing, S. E., & Ward, S. M (2019). Using invaded-range species  
470 distribution modeling to estimate the potential distribution of *Linaria* species and their  
471 hybrids in the US northern Rockies. *Invasive Plant Science and Management*, 12(2), 97-  
472 111.

473 Mehyadin, A. E., Abdulazeez, A. M., Hasan, D. A., & Saeed, J. N (2021). Birds sound  
474 classification based on machine learning algorithms. *Asian Journal of Research in*  
475 *Computer Science*, 9(4), 1-11.

476 Mendes, E.D., Wooley, J., Pi, Y., Tao, J., & Tedeschi, L.O (2024). 419 Automated individual  
477 animal identification and feeding bunk scoring: a computer vision approach for beef  
478 cattle at Calan gate feeding system. *Journal of Animal Science*.

479 Miao, Z., Gaynor, K. M., Wang, J., Liu, Z., Muellerklein, O., Norouzzadeh, M. S., McInturff, A.,  
480 Bowie, R. C. K., Nathan, R., Yu, S. X., & Getz, W. M. (2019). Insights and approaches  
481 using deep learning to classify wildlife. *Scientific Reports*, 9(1), 8137.

482 Moallem, G. (2021). Applications of image processing and machine learning techniques in  
483 wildlife monitoring and cancer cell characterization (Doctoral dissertation).

484 Müller, T., Philippi, N., Dandekar, T., Schultz, J., & Wolf, M. (2007). Distinguishing  
485 species. *Rna*, 13(9), 1469-1472.

486 Murdoch, W. J., Singh, C., Kumbier, K., Abbasi-Asl, R., & Yu, B. (2019). Interpretable machine  
487 learning: definitions, methods, and applications. *arXiv preprint arXiv:1901.04592*.

488 Nadimpalli, U. D., (2005). "Image processing techniques to identify predatory birds in  
489 aquacultural settings" LSU Master's Theses. 276.

490 Norouzzadeh, M. S., Nguyen, A., Kosmala, M., Swanson, A., Palmer, M. S., Packer, C., &  
491 Clune, J. (2018). Automatically identifying, counting, and describing wild animals in  
492 camera-trap images with deep learning. *Proceedings of the National Academy of*  
493 *Sciences*, 115(25), E5716–E5725. <https://doi.org/10.1073/pnas.1719367115>

494 Perry, W. L., Lodge, D. M., & Feder, J. L (2002). Importance of hybridization between  
495 indigenous and nonindigenous freshwater species: an overlooked threat to North  
496 American biodiversity. *Systematic Biology*, 51(2), 255-275.

497 Picard, S., Chapdelaine, C., Cappi, C., Gardes, L., Jenn, E., Lefèvre, B., & Soumarmon, T  
498 (2020). Ensuring dataset quality for machine learning certification. In *2020 IEEE*

499            *international symposium on software reliability engineering workshops (ISSREW)*. (pp.  
500            275-282). IEEE.

501 Reyes, D., & Sánchez, J (2024). Performance of convolutional neural networks for the  
502            classification of brain tumors using magnetic resonance imaging. *Heliyon*, 10(3).

503 Sella, G., & Barton, N. H. (2019). Thinking about the evolution of complex traits in the era of  
504            genome-wide association studies. *Annual review of genomics and human genetics*, 20(1),  
505            461-493.

506 Sharma, S., Sato, K., & Gautam, B. P. (2022). Bioacoustics Monitoring of Wildlife using  
507            Artificial Intelligence: A Methodological Literature Review. *2022 International*  
508            *Conference on Networking and Network Applications (NaNA)*, 1–9.  
509            <https://doi.org/10.1109/NaNA56854.2022.00063>

510 Smith, L. N., & Topin, N (2019, May). Super-convergence: Very fast training of neural networks  
511            using large learning rates. In *Artificial intelligence and machine learning for multi-*  
512            *domain operations applications* (Vol. 11006, pp. 369-386). SPIE.

513 Soni, P., Dhavale, S., Yenishetti, S., Panat, L., & Karajkhede, G. (2023, October). Medicinal  
514            plant species identification using AI. In *2023 IEEE 11th Region 10 Humanitarian*  
515            *Technology Conference (R10-HTC)* (pp. 662-668). IEEE.

516 Sørensen, T. (1948). A method of establishing groups of equal amplitude in plant sociology  
517            based on similarity of species and its application to analyses of the vegetation on Danish  
518            commons. *Biol Skrifter/Kongelige Danske Videnskabernes Selskab.*, 5, 1.

519 Szegedy, C., Ioffe, S., Vanhoucke, V., & Alemi, A (2017, February). Inception-v4, inception-  
520            resnet and the impact of residual connections on learning. In *Proceedings of the AAAI*  
521            *conference on artificial intelligence* (Vol. 31, No. 1).

522 Tan, M., & Le, Q. V. (2020). *EfficientNet: Rethinking Model Scaling for Convolutional Neural*  
523 *Networks* (arXiv:1905.11946). arXiv. <https://doi.org/10.48550/arXiv.1905.11946>

524 Uchida S., Ide, S., Iwana, B. K., and Zhu, A (2016). A Further Step to Perfect Accuracy by  
525 Training CNN with Larger Data. *2016 15th International Conference on Frontiers in*  
526 *Handwriting Recognition (ICFHR)*. 405-410, doi: 10.1109/ICFHR.2016.0082.

527 Ueda, K., & Loarie, S (2008). iNaturalist. Retrieved from <https://www.inaturalist.org/>

528 Ullah, F., Saqib, S., & Xiong, Y. C. (2025). Integrating artificial intelligence in biodiversity  
529 conservation: bridging classical and modern approaches. *Biodiversity and Conservation*,  
530 34(1), 45-65.

531 Van Der Sluis, S., Verhage, M., Posthuma, D., & Dolan, C. V. (2010). Phenotypic complexity,  
532 measurement bias, and poor phenotypic resolution contribute to the missing heritability  
533 problem in genetic association studies. *PloS one*, 5(11), e13929.

534 Wells, C. P., Lavretsky, P., Sorenson, M. D., Peters, J. L., DaCosta, J. M., Turnbull, S., ... &  
535 Engilis Jr, A (2019). Persistence of an endangered native duck, feral mallards, and multiple  
536 hybrid swarms across the main Hawaiian Islands. *Molecular Ecology*, 28(24), 5203-5216.

537 Yosinski, J., Clune, J., Nguyen, A., Fuchs, T., & Lipson, H. (2015). Understanding Neural  
538 Networks Through Deep Visualization (No. arXiv:1506.06579). *arXiv*.  
539 <https://doi.org/10.48550/arXiv.1506.06579>

540 Yoss, A. (2020). Transfer learning using pre-trained AlexNet model and fashion-MNIST.  
541 *Towards Data Science*. [https://towardsdatascience.com/the-w3h-of-alexnet-vggnet-resnet-](https://towardsdatascience.com/the-w3h-of-alexnet-vggnet-resnet-and-inception-7baaaecccc96)  
542 [and-inception-7baaaecccc96](https://towardsdatascience.com/the-w3h-of-alexnet-vggnet-resnet-and-inception-7baaaecccc96)

543 Zhou, Y., Zhu, H., Chai, Y., Jiang, Y., & Liu, Y. (2024). *Towards Trustworthy Web Attack*  
544 *Detection: An Uncertainty-Aware Ensemble Deep Kernel Learning Model*  
545 (arXiv:2410.07725). arXiv. <https://doi.org/10.48550/arXiv.2410.07725>

**Table 4.** Accuracy of identification of the same 30 individuals, compared by AI models, non-experts and experts.

	AI CORRECT	AI INCORRECT	AI % ACCURACY						
<b>FULL ID</b>	12	18	40.00%						
<b>SEX</b>	16	14	53.33%						
<b>SPECIES</b>	20	10	66.67%						
	NON EXPERT 1 CORRECT	NON EXPERT 1 INCORRECT	NON EXPERT 1 % ACCURACY	NON EXPERT 2 CORRECT	NON EXPERT 2 INCORRECT	NON EXPERT 2 % ACCURACY	NON EXPERT 3 CORRECT	NON EXPERT 3 INCORRECT	NON EXPERT 3 % ACCURACY
<b>FULL ID</b>	13	17	43.33%	15	15	50.00%	17	13	56.67%
<b>SEX</b>	20	10	66.67%	20	10	66.67%	20	10	66.67%
<b>SPECIES</b>	20	10	66.67%	23	7	76.67%	26	4	86.67%
	EXPERT 1 CORRECT	EXPERT 1 INCORRECT	EXPERT 1 % ACCURACY	EXPERT 2 CORRECT	EXPERT 2 INCORRECT	EXPERT 2 % ACCURACY			
<b>FULL ID</b>	26	4	86.67%	25	5	83.33%			
<b>SEX</b>	27	3	90.00%	30	0	100.00%			
<b>SPECIES</b>	26	4	86.67%	24	6	80.00%			

**Table 3.** Aggregate AI model classification results for species and sex identification, along with the according score and final model confidence. Text in red highlights an incorrect identification.

Sample	Correct Species	Identified Species	Correct Sex	Identified Sex	Score	Confidence
1	Mallard	Mexican Duck	Female	Female	4	77%
2	Mallard	Mexican Duck	Male	Female	3	70%
3	Mallard	Mallard	Male	Female	8	71%
4	Mexican Duck	Mexican Duck	Male	Male	4	79%
5	Mexican Duck	Mexican Duck	Male	Male	3	82%
6	Mexican Duck	Mexican Duck	Female	Male	4	87%
7	Mallard	Mallard	Female	Female	7	75%
8	Mallard	Mallard	Male	Male	8	71%
9	Mexican Duck	Mexican Duck	Male	Female	3	70%
10	Mexican Duck	Mallard	Male	Female	9	78%
11	Mexican Duck	Unknown	Male	Male	6	81%
12	Mexican Duck	Mexican Duck	Female	Male	1	75%
13	Mexican Duck	Mexican Duck	Female	Female	2	81%
14	Mallard	Mallard	Male	Female	3	74%
15	Mallard	Unknown	Male	Female	3	80%
16	Mallard	Mallard	Male	Male	3	71%
17	Mallard	Mallard	Male	Male	3	75%
18	Mallard	Unknown	Female	Male	0	83%
19	Mallard	Unknown	Female	Male	0	77%
20	Mallard	Mallard	Male	Male	3	77%
21	Mallard	Mexican Duck	Female	Male	0	77%
22	Mallard	Mallard	Male	Male	3	79%
23	Mallard	Mallard	Female	Male	3	79%
24	Mexican Duck	Mexican Duck	Male	Male	0	76%
25	Mexican Duck	Mexican Duck	Male	Male	0	72%
26	Mexican Duck	Unknown	Female	Female	3	78%
27	Mexican Duck	Unknown	Male	Male	0	77%
28	Mexican Duck	Mexican Duck	Male	Female	0	79%
29	Mexican Duck	Mexican Duck	Female	Male	0	72%
30	Mexican Duck	Mexican Duck	Male	Male	0	71%

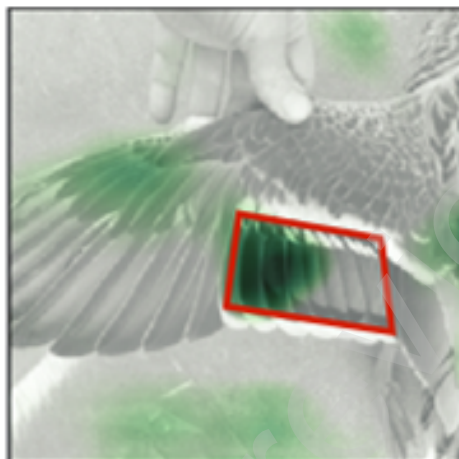
Full Species		
Correct	12	40.00%
Incorrect	18	60.00%
Sex		
Sex Correct	16	53.33%
Sex Incorrect	14	46.67%
Species		
Species Correct	20	66.67%
Species Incorrect	10	33.33%

**Figure 2.** Graphical representation of the body regions used for training individual trait-specific models. *Art credit to August Konvalin.*

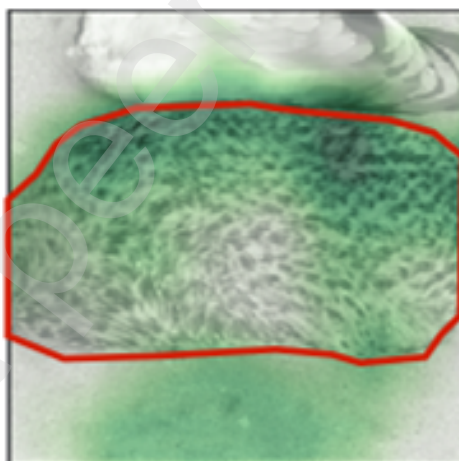
Preprint not peer reviewed

**Figure 3.** Example of heatmap integration within the (A) speculum, (B) belly, and (C) sex models by using the IntegratedGradients package, demonstrating interpretability through visualization of regional activation of each model. Note that while traits of interest were focused on, activation of backgrounds and non-trait regions were evident across models

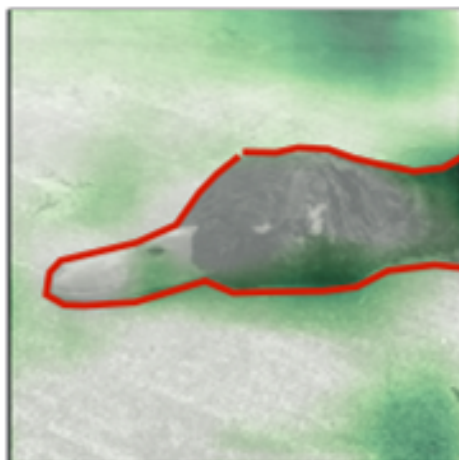
**(A) Speculum**



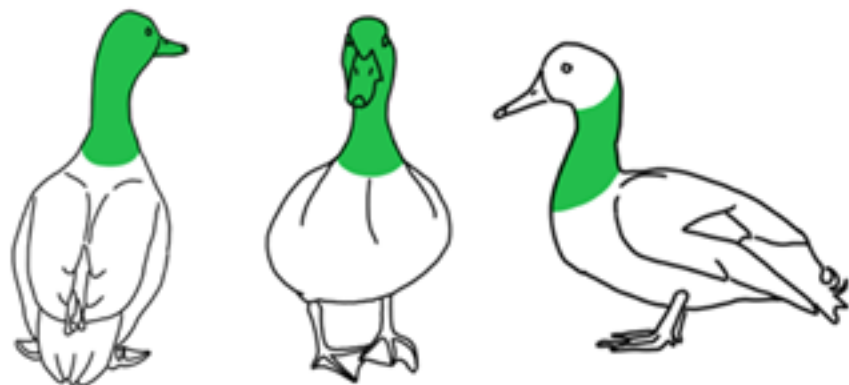
**(B) Breast & Belly**



**(C) Sex**



DORSAL VENTRAL SIDE



= % GREEN  
= FACE/NECK  
= SEX

DORSAL

VENTRAL



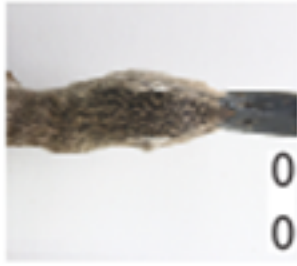
= PRIMARY COVERT  
= LESSER COVERT  
= GREATER SECONDARY  
= SPECULUM COLOR



= BACK PATTERN  
= OUTER TAIL FEATHERS  
= BREAST/BELLY



Preprint not peer reviewed



Preprint not peer-reviewed

**Figure 1.** Per trait-specific model Confusion matrices for the models, here showing the matrix for the sex classifier. 9 images were selected from the test pool with the model's predictions as well as the correct label printed under the prediction. In this case, out of the 9 images, only the bottom middle was incorrectly labelled as male (1) when it should have been female (0).