

1 **Bornean Orangutan Nest Identification Using**
2 **Computer Vision and Deep Learning Models to**
3 **Improve Conservation Strategies**

6 Donna Simon^{1,2}, Keeyen Pang³, Rayner Bili⁴, Song-Quan Ong^{2*}, Henry Bernard^{2*}

8 ¹WWF-Malaysia, Sabah Office, Centre Point Complex, Jalan Centre Point, Kota Kinabalu,
9 Sabah, Malaysia.

10 ²Institute for Tropical Biology and Conservation, Universiti Malaysia Sabah, Jalan UMS, Kota
11 Kinabalu, Sabah, Malaysia.

12 ³Intrajasa Sdn.Bhd, Lot 9, Harapan Baru Light Industrial Estate. Mile 8, Jalan Labuk, 90009
13 Sandakan, Sabah, Malaysia

14 ⁴Sabah Forestry Department, Mile 6, Labuk Road, Sandakan, Sabah, Malaysia.

16 Corresponding Author:

17 Song-Quan Ong¹

18 Jalan UMS, Kota Kinabalu, Sabah, Malaysia.

19 Email address: songquan.ong@ums.edu.my

21 **Abstract**

22 **Background.** Regular population surveys are crucial for the evaluation of conservation measures
23 and the management of critically endangered species such as the Bornean orangutans. Uncrewed
24 aerial vehicles (UAV) are useful for monitoring orangutans by capturing images of the canopy,
25 including nests, to monitor their population. However, manually detecting and counting nests from
26 UAV imagery is time-consuming and requires trained experts. Computer vision and deep learning
27 (DL) for image classification offer an excellent alternative for orangutan nest identification.

28 **Methods.** This study investigated DL for nest recognition from UAV imagery. A binary dataset
29 (“with nest” and “without nest”) was created from UAV imagery from Sabah, Malaysian Borneo.
30 The images were captured using a fixed-wing UAV with a complementary metal-oxide
31 semiconductor camera. After image augmentation, 1624 images were used for the dataset and
32 further split into 70% training, 15% testing and 15% validation for model performance evaluation,
33 i.e. accuracy, precision, recall and F1-score. Four DL models (InceptionV3, MobileNetV2,
34 VGG19 and Xception) were trained to learn from the labeled dataset and predict the presence of
35 nests in new images.

36 **Results.** The results show that out of... (how many variants you had at the beginning??)
37 InceptionV3 has the best model performance with more than 99% accuracy and precision, while
38 VGG19 has the lowest performance. In addition, gradient-weighted class activation maps were
39 used to interpret the results, allowing visualization of the regions used by InceptionV3 and VGG19

40 for classification. This study demonstrates the potential of integrating DL into orangutan
41 conservation and suggests that future research should focus on automatic nest detection to improve
42 UAV-based monitoring of orangutans.

43

44

45 **Introduction**

46 All three orangutan species (it is worth to know - please mention them) living on Borneo and
47 Sumatra have been listed as 'Critically Endangered' on the International Union for Conservation
48 of Nature (IUCN) Red List since 2016, due to significant population declines (Ancrenaz et al.
49 2023). These population declines are primarily driven by habitat loss, degradation, and
50 fragmentation, along with retaliation killings due to conflicts with humans (Ancrenaz et al. 2023).
51 In Sabah, Malaysian Borneo, several measures have been introduced to protect orangutans,
52 including forest restoration in degraded areas (Mansourian et al. 2020), expanding totally protected
53 areas to 30%, and committing to sustainable timber production (Simon et al. 2019). Additionally,
54 the 10-year Sabah Orangutan Action Plan (2020-2029) was developed to ensure the species' long-
55 term viability in the region (Sabah Wildlife Department, 2020). Continuous monitoring is
56 crucial to assess population trends and evaluate the effectiveness of these conservation efforts
57 (Piel et al. 2022).

58 Orangutans are primarily found in lowland tropical rainforests (how many elevations?) ,
59 where they spend most of their time in the forest canopy (Manduell et al. 2012). They construct
60 new nests each day, with juveniles relying on their mothers to build them (Permana et al. 2024).
61 These nests are used for both night-time sleeping and daytime resting (Casteren et al. 2012). Since
62 observing orangutans directly is difficult due to the dense canopy and their elusive nature,
63 researchers often monitor populations by counting nests, which serve as reliable indicators of their
64 presence (Kuhl et al. 2008; Santika et al. 2019). Population estimates are derived from nest
65 densities (nests per km²), which are converted into orangutan numbers using established statistical
66 methods (Ancrenaz et al. 2005; Kuhl et al. 2008; Pandong et al. 2018).

67 Orangutan nests are distinct from those of other animals. Orangutans typically build their
68 nests in the upper canopy, around 11-20 meters above the ground (Casteren et al. 2012), and the
69 nests are about 100 cm wide to accommodate their large body size (Kamaruszaman et al. 2018).
70 The nest's base is made from thick branches, with thinner branches twisted and bent but not fully
71 broken. This partial break, known as a "greenstick fracture," is unique to orangutan nests (Casteren
72 et al. 2012). Leaves are added to form a flat sleeping platform. Orangutan nests are usually oval
73 and asymmetrical, with the long axis oriented towards the tree trunk (Biddle et al., 2014). While
74 most nests are built in the upper canopy, they can also be found at branch ends or close to the main
75 tree stem (Rayadin and Saitoh 2009).

76 Various methods are used to count orangutan nests, including ground-based nest surveys
77 (Pandong et al. 2018; Santika et al. 2019), helicopter surveys (Ancrenaz et al. 2005; Payne 1988;
78 Simon et al. 2019), and the latest technology involving uncrewed aerial vehicles (UAVs) or drones
79 (Hanggito 2020; Milne et al. 2021; Wich et al. 2015). Among these methods, drones are becoming

80 increasingly important as they are relatively inexpensive compared to helicopters and can capture
81 images or time-lapse video from the forest canopy, allowing many hard-to-access areas to be
82 studied (Wich and Koh 2018). In contrast to ground and helicopter surveys, where nests are
83 detected through direct field observations, drone imagery requires careful examination of each
84 image on a computer to identify nests. As nests decay, the fresh green foliage withers and turns
85 brown, making them stand out more clearly against the surrounding green canopy in the images
86 (Figure 1). During manual nest identification, each nest is marked or labelled and then counted
87 across all images. This allows researchers to calculate nest density, which can be used to estimate
88 the orangutan population size.

89 To classify the images, it is important to consider the canopy classification perspective.
90 Although nests made of branches and leaves can be distinguished from healthy trees as they decay
91 over time (Casteren et al. 2012), a key challenge in using drone imagery to explore orangutan nests
92 is that labeling nests from large volumes of image data still relies heavily on human experts,
93 making the process tedious and time-consuming (Milne et al. 2021; Wich et al. 2015). Therefore,
94 there is a need for an alternative method to identify nests from drone imagery that is as effective
95 as, if not more effective than, human expertise for nest detection.

96
97
98
99 **Figure 1 Example** images for a drone image with orangutan nests circled in red

100
101
102
103 The integration of artificial intelligence (AI) into nest detection is a possible alternative to improve
104 the efficiency of counting nests in drone image surveys. Machine learning (ML), is a branch of AI
105 that enables computers to learn from a diverse array of data, recognize patterns and make decisions
106 with minimal human intervention (Chahal and Gulia 2019). This study uses supervised learning, a
107 category of ML in which algorithms are trained on labelled datasets to predict outcomes and
108 recognize patterns. In contrast to unsupervised learning, supervised learning algorithms are trained
109 with labelled data to learn the relationship between the inputs, i.e., features such as colour, texture,
110 and shape of objects in aerial images, and the outputs i.e., labels indicating the presence or absence
111 of orangutan nests in those images (Wang et al. 2016). The term "annotation" used to label the
112 presence of orangutan nests on aerial images, is similar to the term "data labelling" in supervised
113 ML, where the images of the nests need to be labelled and used as training data for model
114 development.

115 It is also important to note that the matrices used to evaluate the context of the ecological study
116 and the ML model may be similar, e.g. accuracy and precision, but they differ in context. In an
117 ecological study, accuracy is the difference between the sample estimates and the true population
118 value (Hellmann and Fowler 1999). For example, the accuracy of species richness is the difference
119 between the estimate of species richness based on sample data and the true species richness of the

120 population or community being sampled. Whereas, precision is the difference between an estimate
121 of species richness based on sample data and the significance of all possible estimates of species
122 richness based on all possible samples of the same size from the sampled population or community
123 (Hellmann and Fowler 1999). In calculations, accuracy is measured by the mean square error of
124 the estimator and precision by the variance of the estimator. On the other hand, the ML model is
125 evaluated by the true value [the actual number of true positives (TP), true negatives (TN), false
126 positives (FP) and false negatives (FN) from a test set of the prediction] (Lebovitz et al. 2021). In
127 most cases, the result of the evaluation of the model can be expressed in a layout table, the so-
128 called confusion matrix, in which the proportion of TP and TN can be calculated. Accuracy, for
129 example, is the proportion of all classifications that were correct, whether positive or negative, and
130 precision is the proportion of all positive classifications of the model that are actually positive.
131 Since the model is calculated on the basis of ground truth, further evaluation matrices can be
132 calculated, e.g. recall or True Positive Rate (TPR) and False Positive Rate (FPR), which is crucial
133 for the evaluation of a model with an unbalanced data set. For example, where the number of
134 instances in one class (e.g., positive cases) is significantly lower than that in another class (e.g.,
135 negative cases).

136 Another subset of ML, often used in computer vision and image classification, is known as deep
137 learning (DL) or deep convolutional neural networks (DCNN). DL has been used extensively for
138 aerial images classification. Pearse et al. (2021), for example, have shown how DL models can
139 classify tree species by learning complex visual features of tree species from aerial images with an
140 accuracy of 92%, a sensitivity of 91% and a specificity of 94%. For monitoring orangutan
141 population, Davies et al. (2019) combined LiDAR and behavioural data to reveal relationships
142 between tree canopy structure and nest choice of orangutans in disturbed forests.

143 DL models are well suited for image classification as the architecture uses multiple layers of neural
144 networks consisting of perceptions to model complex data (e.g. images with different colour
145 channels) by learning features from images and making predictions (Smith et al. 2018). Further
146 details on how the DL model works can be found in Wang et al. (2016), Purwono et al. (2022), the
147 protocol paper by Isawasan et al. (2023) and Madhavan & Jones (2024). In image processing, DL
148 is widely used for image classification and object detection in ecological studies, such as species
149 identification, animal behavior classification and species diversity estimation from camera traps,
150 video and audio recordings (Christin et al. 2019). For orangutan studies, Guo et al. (2020)
151 developed Tri-AI, an automatic recognition system that identifies 41 primates and four carnivores
152 with 94% accuracy. In addition, Desai et al. (2023) developed an annotated database of apes in
153 different poses which enables object recognition for behavioral studies of apes in zoos.

154 Studies on orangutan recognition through computational methods to detect and count
155 orangutan nests remain limited. Nest building, a unique daily behavior of orangutans for sleeping,
156 offers valuable data for ecological monitoring, and by integrating DL techniques, it could enhance
157 population monitoring efforts. Amran et al. (2023) initiated the study on the use of ML – Support
158 Vector Machine (SVM) - in classifying the objects on the aerial images into branches, buildings
159 and orangutan nests; Teguh et al. (2024) provided the most recent study (at the time of writing this

160 manuscript) on orangutan detection using DL model, the You Only Look Once (YOLO) version 5
161 with 414 labelled orangutan nests and achieved a precision of 0.973 and a recall of 0.949. However,
162 Teguh et al. (2024) applied an object detection algorithm and demonstrated the effectiveness of a
163 DL model, but this raises additional questions. For instance, YOLO typically identifies and
164 classifies objects in a single step, but alternative classification algorithms may offer improved
165 performance. As biologists and ecologists, it is crucial not to treat these tools as black boxes. This
166 study focuses on interpreting the outputs to gain insight into how DL models 'visualize' image
167 patterns and identify the features utilized by neural network layers to classify tree canopy patterns
168 as 'with nest' or 'without nest.' Understanding this process is essential for accurate ecological
169 interpretation.

170 Therefore, this study aims to evaluate the effectiveness of different DL models in detecting
171 orangutan nests from aerial images captured from two orangutan habitats in Sabah, Malaysia. More
172 importantly, this study will visualize the model layers to understand how the features and
173 characteristics of orangutan nests are 'learned' by the model. Specifically, the aim of this study is
174 to create a labelled dataset of drone images containing the presence and absence of orangutan nests,
175 and finally to develop and compare four DL models for detecting and predicting nest presence
176 from drone images. Additionally, gradient-weighted class activation maps (Grad-CAM) were
177 presented which can visualize the activation region used by the models to distinguish orangutan
178 nests from the tree canopy.

179

180

181 **Materials & Methods**

182 **Study site**

183 These drone surveys were conducted in Sepilok Virgin Jungle Reserve (VJR) and Bukit Piton
184 Forest Reserve (FR) in Sabah, Malaysia (Fig 2.). Both reserves are under the management of the
185 Sabah Forestry Department and are known habitats for orangutans. It is estimated that there are
186 about 200 (100-300) orangutans in Sepilok (Ancrenaz et al. 2005) and 176 (119-261) orangutans
187 in Bukit Piton (Simon et al. 2019). The Sepilok VJR covers an area of approximately 40 km² and
188 is characterized by lowland dipterocarp and heath forests (Ball et al. 2023). The reserve has been
189 designated as a protected area where logging is strictly prohibited to keep the forest canopy intact.
190 In contrast, Bukit Piton FR, which consists mainly of dipterocarp lowland rainforest and is about
191 120 km² in size, is severely degraded due to heavy logging and forest fires in the past. In 2008, a
192 large-scale project was initiated to restore the forest for orangutans and the area was declared as a
193 protected forest in 2012. Since then, the forest has slowly regenerated, with fast-growing tree
194 species being used by the orangutans for nesting just three years after planting (Mansourian et al.
195 2020).

196

197

198

199 Fig 2. Location of Sepilok Forest Reserve and Bukit Piton Forest Reserve

200

201 **Study duration**

202 The Sepilok survey was conducted in July 2015 and covered an area of approximately 0.5 km². A
203 total of three flight missions were conducted to complete the survey with 1720 number of images.
204 The Bukit Piton survey was conducted in January 2016 and covered an area of approximately 0.5
205 km² resulting in 1911 images. A total of 4 missions were flown to survey the area in January 2016.
206 Both surveys were conducted in the morning on a sunny day (Temperature, Relative Humidity?).
207

208 **Equipment**

209 This study utilizes UAV imagery captured by a fixed-wing drone assembled by
210 ConservationDrones.org using an FX-71 frame. A Canon Power Shot S100 digital camera with
211 RBG CMOS sensor (type number usu. with code, made in what country?) was installed in the
212 drone. The drone was flown at least 100 meters from the highest point, which was determined
213 using the Digital Elevation Model (DEM). The time-lapse recordings were made in 3-second
214 intervals. The DL models compared in this study solve a classification problem in which the
215 models process the entire image as a target object instead of recognizing different objects from
216 one image (Sharma 2019). The datasets were created by combining images from both locations
217 and having four human experts examine them for nests, annotate them and categorize them into
218 two binary classes, i.e. images with nests and images without nests. This binary classification is
219 needed to train the model and determine whether an image contains an orangutan nest or not. The
220 field study and the use of the drone for aerial images **were conducted in** 2014 with the permission
221 of the Sabah Forestry Department under reference number (JPHTN/PP 100-22/4/KLT.11(44)).
222

223 **Pre-processing of the data and categorization**

224 Using the image classification task, the entire images were classified either into "with nest" or
225 "without nest". For images with multiple nests, the images were pre-processed by cropping out
226 the nest and labelling it as "with nest". The total number of aerial drone images from both
227 Sepilok and Bukit Piton is 406 images, which were further classified into two classes, i.e., with
228 nest (162 images) and without nest (244 images) (Table 1).

229 Nests from drone images have been identified by six orangutan field specialists, with more
230 than two years of field experience in conducting ground and helicopter nest surveys. The
231 identification of the orangutan nest at the same sites where drone images were captured is also
232 consistent with the ground survey data which confirmed the presence of nests through direct
233 observations. Then, the total number of images in each class was divided into three parts, also
234 known as data splitting, with 70% of the total images used for training, 15% for validation and
235 15% for testing or a 70:15:15 ratio (Figure 2). The ratio of data splitting is based on the amount of
236 data used for training and evaluation, and reducing the size of the training dataset tends to result
237 in a poorly performing model. Therefore, an international standard of computer vision and DL
238 competition (Fei-Fei et al. 2009) was referenced, along with insights from previous studies (Khan
239 and Ullah 2022; Ong and Hamid 2022). Data splitting enables the machine to use the training set

240 to obtain the weights and biases for classification. The validation set helped to better generalize
241 the models to new, unseen data and prevent over-fitting while the testing set is to assess the model's
242 performance. As the number of images was relatively small, each image was subjected to a rotation
243 expansion of 0°, 90°, 180° and 270° and finally the number of images was increased by a factor of
244 four (Ong et al. 2022; Chen et al. 2021), totaling to 1624 images used for the model development.
245

246

247 **Models development**

248 *Model build-up*

249 To develop the DL models, the convolutional blocks of the pre-trained convolutional neural
250 networks (CNNs) were unfrozen for retraining purposes (a process in which the weights and biases
251 that the model learns from the ImageNet are unlocked for a customized task, i.e., orangutan nest
252 classification). This was done for four DL architectures – InceptionV3, MobileNetV2, VGG19 and
253 Xception – to optimize them for the specific task of identifying nests from aerial images, as
254 described in Ong et al. (2022). The Keras DL Framework on an NVIDIA Tesla A100 Google
255 Compute Engine (GPU) platform was used to train and evaluate the models. The models were
256 trained with the Adaptive Moment Estimation (ADAM) optimizer, which improves the stability
257 and efficiency of the training process and enables efficient learning (Shao 2024). Three learning
258 rates (0.01, 0.001 and 0.0001) with 32 batches were analyzed. The training process was set to 50
259 epochs, meaning that the model performed 50 complete iterations through the training dataset
260 (Wang et al. 2016). Increasing the number of epochs allows the model to refine its parameters and
261 could improve its performance. After developing the models, the performance of these models
262 were evaluated using the four metrics of accuracy, precision, recall and F1-score (Table 2) (Hosin
263 and Sulaiman 2015; Kumar 2020). In addition, the mean accuracy (number of correct
264 predictions/total number of images) was compared between the models to test the significance of
265 the four DL models. The code that used for the model development was publicly available at github
266 with the link <https://github.com/songguan26/Bornean-Orangutan-Nest->
267

268 *Activation map to distinguish orangutan nests from aerial images*

269 To gain further insight into how the neural network in the DL models can recognize the orangutan
270 nest, Grad-CAM was used to visualize the area used by the neural network to classify the orangutan
271 nest with a variety of normal tree canopy backgrounds. In general, one layer at a time was retrieved
272 to extract low- and high-level features. The code that used for the model development was publicly
273 available at github with the link <https://github.com/songguan26/Bornean-Orangutan-Nest->

274 **Results**

275 *Model performance*

276 Four DL models were attempted, and the images were trained, tested and validated for image
277 classification tasks by classifying UAV images into “without nests” and images “with nests”
278 categories. Fig. 3, shows the performance of the four models in predicting the images with presence

279 or absence of nests. It can be seen that VGG19 performs lower than the other models. InceptionV3,
280 MobileNetV2 and Xception were ranked first, second and third. The Shapiro-Wilk normality test
281 was performed to assess the normality of the accuracy values for the models across three learning
282 rates. The results are as follows: InceptionV3 ($W = 0.75$, $p = 0.0000009$), MobileNet ($W = 0.99$,
283 $p = 0.99$), VGG19 ($W = 0.95$, $p=0.566$) and Xception ($W=0.89$, $p=0.37$). Based on these results,
284 only InceptionV3 is not normally distributed ($p < 0.05$). Therefore, a non-parametric test, the
285 Kruskal Wallis H-test, was used to compare the models based on their accuracy values across three
286 LRs. The result of the Kruskal Wallis H-test shows no significant difference (i.e., at the threshold
287 p-value <0.05) in the accuracy of the four models at three learning rates ($H (3) = 6.751$, $p = 0.087$).
288 Additionally, as most of the models are normally distributed except InceptionV3 with a very small
289 p-value, the model performance is presented in Fig.4 using the mean value to better represent the
290 data.

291 To assess the generalization capabilities of the model — its ability to make accurate
292 predictions on new data (Caro et al. 2022) the training validation accuracy (TVA) and training
293 validation loss (TVL) of the models across three learning rates on the test set were evaluated and
294 presented in Table 2. The new data was validation splits (15%, in section methodology) that were
295 never used in the model development. Although the epochs were set at 50, the early-stopping-
296 method was employed — to prevent overfitting and underfitting (Cai et al. 2022) causing the model
297 computation to halt early once the validation accuracy did not improve (epochs indicated in X-
298 axis). The results of TVA and TVL (**Table 3**) show that LR 0.001 generally achieves a balance
299 between efficient training and robust generalization across the models. Whereas, LR 0.01 risks
300 instability and overfitting, which occurs when the model fits the training data too closely and failed
301 to generalize to new data (Charilaou and Battat 2022). Meanwhile, LR 0.0001 results in slow or
302 failed convergence and underfitting is shown by the poor performance of VGG19 model, which is
303 incapable of learning the patterns in the training data (Jabbar and Khan 2015).

304 In addition, the confusion matrix for each model is shown in **Table 4** to visualize how well
305 the classification model works by showing the correct and incorrect predictions made by the
306 model, in comparison with the actual answer. The confusion matrix in binary classification consists
307 of four components i.e. True positives (TP) is when the model correctly predicts the positive class;
308 True negatives (TN) is when the model correctly predicts the negative class; False positives (Type-
309 1 error) is when the model incorrectly predicts the positive class and False negative (Type-2 error)
310 when the model incorrectly predicts the negative class (Saito and Rehmsmeier 2015). InceptionV3
311 at LR 0.01, LR 0.001 and Xception at LR 0.0001 have made all correct predictions. Meanwhile,
312 InceptionV3 and MobileNetV2 at LR 0.0001, Xception at LR 0.01 and LR 0.001, as well as
313 VGG19 at all LR, have a Type-1 error in nest prediction. Whereas MobileNetV2 at LR 0.001 has
314 a Type-2 error in nest prediction.

315
316
317
318

319 **Identification and visualization of input features**

320 Heatmaps illustrate which parts of an image the model considers important by highlighting them
321 in warm colors such as yellow, orange and red. Due to the superior overall performance of
322 InceptionV3, five convolutional layers of the InceptionV3 architecture covering the low- and high-
323 level features were used to visualize how the neural network identified the orangutan nest. Table
324 5 shows some examples of the convolutional layers of InceptionV3 compared to the original image
325 of a human. The most common 2D convolutional layer “Conv2d” (Khan 2019) is used to visualize
326 the region used by the model for classification. The heatmaps derived from Conv2d_89 and
327 Conv2d_90 highlighted the corners of the images and underlined subtle colors on the nest itself.
328 In contrast, the nest was emphasized in the Conv2d_91 and Conv2d_92 heatmaps. In addition, the
329 upper right corner of the image was emphasized in the heatmap derived from Conv2d_93. Based
330 on the result, the neural network was able to identify the features of the nest – edge, shape and
331 texture – reflected in the different intensities of warm color. As mentioned by LeCun et al. (2015),
332 there were blocks of low and high feature extraction in InceptionV3. Fig. 5 shows an example of
333 the original image used to extract the feature for classification.

334
335
336

337 **Discussion**

338 The increasing use of drones to monitor orangutan populations since when?? This would be
339 interesting to know ..could be an excellent alternative to improve the monitoring and protection of
340 orangutan populations. However, the enormous amount of data generated by UAV imagery, which
341 needs to be identified and annotated by trained experts, poses a major time and labor-intensive
342 challenge. What about the supporting facilities (equipment) as well as internet stability, where this
343 is sometimes challenging for some developing countries? Therefore, this study was conducted with
344 the aim of evaluating the feasibility of using computer vision and DL to classify orangutan nests
345 from UAV imagery.

346 This study is focused on image classification rather than on object detection (Sharma 2019).
347 Specifically, it supports the second stage of the two-stage object recognition algorithm which in
348 this case involves identifying the orangutan nest. The concept of two-stage detection consists of
349 the first stage of detecting the object of interest (usually with the YOLO or SSD algorithm) and
350 the second stage of a classifier by a DL algorithm (the DL models investigated in this study).
351 Although many data scientists or ML engineers have proposed only the YOLO algorithm, which
352 can solve both localization (detecting the position of the object of interest on an image) and
353 classification in one step, detecting and classifying an orangutan nest on aerial images of tree
354 canopies is a great challenge in reality (due to the very similar patterns of tree canopies) and
355 requires a large number of aerial images as training data.

356 The result of this orangutan nest recognition study is consistent with that of Chen et al.
357 (2014), who integrated various AI methods, including ML, optimization algorithms and adaptive
358 decision-making systems, to develop intelligent systems capable of performing complex orangutan

359 nest detection tasks from UAV imagery. In addition, the current study on the use of DL
360 architectures with feature extraction from the images has continued the study of Amran et al.
361 (2023) who used hand-crafted feature extraction and multi-class classification with Support Vector
362 Machines (SVM) for orangutan nest in Borneo. Although Teguh et al. (2024) attempted to use
363 YOLOv5 and achieved a precision of 0.973 and a recall of 0.949 when recognizing the orangutan
364 nest from the drone images, this study has shown that orangutan nest recognition can achieve
365 higher accuracy and precision when using lower computational power (and focusing only on the
366 classification task). In addition, this study has shown that unlike YOLO (single-stage recognition
367 algorithm), the use of transfer learning (transferring weights and bias in the classification of
368 ImageNet images to another classification task) also helps to overcome the problem of data scarcity
369 associated with the lack of sufficient training examples. While counting nests from the ground is
370 easier than locating and counting individual orangutans, drone surveys capture only a fraction of
371 nests in aerial views. Nests under the canopy in dense forests are often missed, and fresh green
372 nests or those in advanced decay stages are harder to detect in drone images. As a result, this may
373 cause insufficient training data for model training. Please highlight the challenges when using the
374 UAV imagery vs. with manual observation incl. counting??

375 So far, this study was the first to compare four state-of-the-art pre-trained DL models -
376 InceptionV3, MobileNetV2, VGG19 and Xception. The data was further augmented and the
377 hyperparameters were refined by training for nest recognition from UAV imagery, resulting in
378 high accuracies (>96%). The model performance result is in line with Ong and Hamid (2022) and
379 Ong et al. (2022), where InceptionV3 is the best model for this task, while VGG19 performs the
380 worst. When comparing between the three learning rates, the learning rate (LR) of 0.001 achieved
381 the optimal performance, with fewer problems related to overfitting and underfitting. InceptionV3
382 with LR 0.001 performed well and delivered all correct predictions.

383 It is worth noting that VGG19 performs the worst in this study, in contrast to other studies
384 which showed that VGG19 performs better than InceptionV3 and MobileNet. A look at the layouts
385 of VGG19 (Table 6) compared to InceptionV3 (Table 5) shows that VGG19 is not able to
386 recognize the features of the orangutan nest, which could be the main reason for the poor
387 performance. Nevertheless, there are previous studies that also show that VGG19 performs worse.
388 This emphasizes the need to compare DL models for a specific task.

389 To interpret the result of the computer vision system for the orangutan nest, the layers of the
390 architecture with Grad-CAM were visualized, which to our knowledge is also the first report.
391 Using Grad-CAM, the region of biases and weights defined by the perceptron within the DL
392 architecture was able to highlight the shape and texture of the orangutan nest, which was later used
393 in the classification block for classification. Considering the similarity of the present study to the
394 task of classifying the canopy of a forest, this study result was compatible with that of Nezami et
395 al. (2020), who used a multilayer perceptron (MLP) to classify tree species using aerial images
396 generated from RGB and hyperspectral (HS) images and achieved an accuracy of 99.6% with the
397 best 3D CNN classifier. Moreover, the result of this study in classifying tree canopy with and
398 without orangutan nests is consistent with that of Huang et al. (2023), who used ResNet,

399 ConvNeXt, ViT and Swin Transformer and achieved at least 96% accuracy in classifying tree
400 species from aerial images.

401 However, there are still many aspects that require further investigation and improvement.
402 One of these is the quality of aerial images. As mentioned by Huang et al. (2023), the degradation
403 of image quality and aerial images at different altitudes needs to be explored further. The key
404 question for future study is to determine what altitude achieves the ideal balance between drone
405 flight feasibility and image quality. In this study, for example, a fixed-wing drone with a Canon
406 Power Shot S100 RGB CMOS sensor was used, which was flown at the highest point of the
407 treetops at an altitude of 100 meters. The image quality could be improved by using a multi-rotor
408 UAV with better camera control. Image quality could also be improved by flying at a lower altitude
409 where the camera is closer to the canopy and can capture more detail. However, this depends on
410 the feasibility of the flight, where many factors determine the closest distance between the drone
411 and the tree canopy, such as the availability of the crash sensor. With better image quality, further
412 exploration can be conducted, such as classifying the nest decay stage of nests and increasing the
413 ability to detect fresh green nests. Additionally, there is a need to augment both the quantity and
414 diversity of aerial imagery to increase the robustness and subsequent generalization of the model.
415 The diversity of the data could also include false positives and negatives in the training data to
416 further improve the generalization of the model. Another important consideration is the
417 deployment of the model to ensure its practical applicability. In the field for detecting and counting
418 the number of orangutan nests. Additionally, building a model by local or regional dataset was
419 always facing a challenge in generalizing good results for other similar datasets (e.g. by using the
420 DL model in this study to predict aerial images from Indonesia).

421 Many future studies will aim to improve the model, software and hardware. However, it is
422 vital to ensure that these improvements consistently contribute to orangutan conservation.
423 Streamlining orangutan survey and monitoring processes to be more cost and time efficient,
424 alongside leveraging computer vision and DL models for automatic annotation of orangutan nests
425 from aerial images, could significantly advance orangutan monitoring efforts.

426

427

428 **Conclusions**

429 The present study encourages further development of DL models for the automatic detection of
430 orangutan nests from aerial UAV images. Further research and refinement in this area could lead
431 to more accurate and efficient methods for identifying nests. Nevertheless, additional data sets,
432 especially from different forest types used by orangutans, such as forest patches within plantations,
433 timber plantations, logged and unlogged forests, are crucial to improve the generalization of the
434 model in the field. In the future, other remote sensing data such as through partnerships with other
435 agencies could be incorporated to obtain more imagery and make significant improvements in this
436 area.

437

438

439 **Ethics statement**

440 The drone was deployed in the primary protection forest where no residents lived, and only images
441 of the canopy were collected, so there was no risk to people's privacy. The field study and the use
442 of the drone for aerial photography were conducted in 2014 with permission from the Sabah
443 Forestry Department under reference number (JPHTN/PP 100-22/4/KLT.11(44)).

444

445

446 **Acknowledgements**

447 This research was funded in part by WWF-UK, WWF-Malaysia and ETIKA. The drone survey
448 activities were funded by AREAS. We would like to express our thanks to Viveca Soripin, Middle
449 Kapis, Tinrus Tindok and William Joseph from the WWF- Orangutan team for their expertise in
450 nest annotation. Additionally, we thank the Sabah Forestry Department for permitting and
451 supporting the survey activities.

452

453

454 **References**

455

456 Amran AA, On CK, Hung LP, Rossdy M, Simon D, See CS (2023) Bornean orangutan nests
457 classification using Multiclass SVM. In 2023 IEEE Symposium on Computers &
458 Informatics (ISCI):1-6

459 Ancrenaz M, Gimenez O, Ambu L, Ancrenaz K, Andau P, Goossens B, Payne J, Sawang, A,
460 Tuuga A, Lackman-Ancrenaz I (2005) Aerial surveys give new estimates for orangutans in
461 Sabah, Malaysia. PLoS Biology 3(1). <https://doi.org/10.1371/journal.pbio.0030003>

462 Ancrenaz M, Gumal M, Marshall A J, Meijaard E, Wich S, Husson S (2023) *Pongo pygmaeus*
463 (amended version of 2016 assessment). The IUCN Red List of Threatened Species 2023:
464 e.T17975A247631797

465 Ball J, Hickman S, Jackson T, Jing K X, Hirst J, Jay W M, and Coomes, D. A. (2023). Accurate
466 delineation of individual tree crowns in tropical forests from aerial RGB imagery using
467 mask r-cnn. Remote Sensing in Ecology and Conservation, 9(5), 641-655.
468 <https://doi.org/10.1002/rse2.332>

469 Biddle L, Deeming D, Goodman A (2014) Morphology and biomechanics of the nests of the
470 common black bird *Turdus merula*. Bird Study 62(1):87-95.
471 <https://doi.org/10.1080/00063657.2014.988119>

472 Casteren AV, Sellers WI, Thorpe SK, Coward S, Crompton RH, Myatt JP, Ennos AR (2012)
473 Nest-building orangutans demonstrate engineering know-how to produce safe, comfortable
474 beds. Proceedings of the National Academy of Sciences, 109(18): 6873-6877

475 Cai Y, Wang Z, Yao L, Lin T, Zhang J (2022) Ensemble dilated convolutional neural network
476 and its application in rotating machinery fault diagnosis. Computational Intelligence and
477 Neuroscience 2022: 1-14. <https://doi.org/10.1155/2022/6316140>

478 Chahal A, Gulia P (2019) Machine learning and deep learning. International Journal of
479 Innovative Technology and Exploring Engineering. 8(12), 2778-3075

480 Charilaou P, Battat R (2022) Machine learning models and over-fitting considerations. World
481 Journal of Gastroenterology, 28(5), 605–607. <https://doi.org/10.3748/wjg.v28.i5.605>

482 Chen H, Guo S, Hao Y, Fang Y, Fang Z, Wu W, Li S (2021) Auxiliary diagnosis for Covid-19
483 with deep transfer learning. Journal of Digital Imaging 34(2): 231-241.
484 <https://doi.org/10.1007/s10278-021-00431-8>

485 Chen Y, Shioi H, Montesinos CF, Koh LP, Wich S, Krause A (2014). Active detection via
486 adaptive submodularity. In Proceedings of the 31st International Conference on Machine
487 Learning: 55–63

488 Christin S, Hervet É, Lecomte N (2019) Applications for deep learning in ecology. Methods in
489 Ecology and Evolution 10(10): 1632–1644. <https://doi.org/10.1111/2041-210X.13256>

490 Davies AB, Oram F, Ancrenaz M, Asner GP (2019) Combining behavioral and LiDAR data to
491 reveal relationships between canopy structure and orangutan nest site selection in disturbed
492 forests. Biological conservation 232: 97-107.

493 Desai N, Bala P, Richardson R, Raper J, Zimmermann J, Hayden B (2023) Open Ape Pose, a
494 database of annotated ape photographs for pose estimation. Elife: 12

495 Fei-Fei L, Deng J, Li K (2009) ImageNet: Constructing a large-scale image database. Journal of
496 vision, 9(8): 1037-1037

497 Guo S, Xu P, Miao Q, Shao G, Chapman CA, Chen X, Li B (2020) Automatic identification of
498 individual primates with deep learning techniques. Iscience 23(8)

499 Hanggito MS (2020) Development of an unmanned aerial vehicle-based orangutan population
500 assessment and monitoring method for the multifunctional landscape of East Kalimantan,
501 Indonesia Open Access Theses and Dissertations
502 https://scholarworks.utep.edu/open_etd/3166

503 Hellmann JJ, Fowler GW (1999) Bias, precision, and accuracy of four measures of species
504 richness. Ecological applications, 9(3), 824-834

505 Huang Y, Wen X, Gao Y, Zhang Y, Lin G (2023) Tree Species Classification in UAV Remote
506 Sensing Images Based on Super-Resolution Reconstruction and Deep Learning. Remote
507 Sensing, 15(11): 2942.

508 Isawasan P, Abdullah ZI, Ong SQ, Salleh KA (2023). A protocol for developing a classification
509 system of mosquitoes using transfer learning. *MethodsX*, 10, 101947.

510 Jabbar HK, Khan ZR (2015) Methods to avoid over-fitting and under-fitting in supervised
511 machine learning (Comparative study). Computer Science, Communication and
512 Instrumentation Devices.

513 Kamaruszaman SA, Nik Fadzly, Mutalib AH, Muslim AM, Atmoko SSU, Mansor M, Mansor A,
514 Rupert N, Zakaria R, Hashim ZH, Sah ASR, Jamsari FF, Azman NM (2018). Measuring
515 Orangutan nest structure using Unmanned Aerial Vehicle (UAV) and Image J. BioRxiv.
516 <https://doi.org/10.1101/365338>

517 Khan MK, Ullah MO (2022) Deep transfer learning inspired automatic insect pest recognition.
518 In Proceedings of the 3rd International Conference on Computational Sciences and
519 Technologies. Jamshoro, Pakistan: Mehran University of Engineering and Technology :
520 17-19.

521 Khan AA (2019) What is Keras Conv2D https://medium.com/@arif_ali/what-is-keras-conv2d-f234e48fde6 Accessed 3 June 2024

522

523 Kühl H, Maisels F, Ancrenaz M, Williamson EA (2008) Best practice guidelines for surveys and
524 monitoring of great ape populations occasional paper of the IUCN Species Survival
525 Commission 36 (36). www.iucn.org/themes/ssc.

526 Lebovitz S, Levina N, Lifshitz-Assaf H (2021) Is AI ground truth really true? The dangers of
527 training and evaluating AI tools based on experts 'know-what.' *MIS quarterly* 45, no. 3

528 LeCun Y, Bengio Y, Hinton G (2015) Deep learning. *Nature* 521(7553): 436-444

529 Li J, Xu Q, Shah N, Mackey TK (2019) A machine learning approach for the detection and
530 characterization of illicit drug dealers on instagram: model evaluation study. *Journal of*
531 *medical Internet research*, 21(6), e13803.

532 Madhavan S, Jones TM (2024). Deep learning architectures. The rise of artificial intelligence.
533 IBM Developer. <https://developer.ibm.com/articles/cc-machine-learning-deep-learning-architectures/> Accessed on 18 September 2024

534

535 Manduell K L, Harrison M E, and Thorpe, S. K. S. (2012). Forest structure and support
536 availability influence orangutan locomotion in Sumatra and Borneo. *American Journal of*
537 *Primateology*, 74(12), 1128-1142. <https://doi.org/10.1002/ajp.22072>

538 Mansourian S, Vallauri D, France W (2020) Lessons learnt from 12 years restoring the
539 orangutan's habitat: the Bukit Piton Forest Reserve in the Malaysian State of Sabah.
540 <https://www.researchgate.net/publication/343609942>

541 Milne S, Martin JGA, Reynolds G, Vairappan CS, Slade EM, Brodie JF, Wich SA, Williamson
542 N, and Burslem DFRP (2021). Drivers of Bornean orangutan distribution across a
543 multiple-use tropical landscape. *Remote Sensing*, 13(3), 1–16.
544 <https://doi.org/10.3390/rs13030458>

545 Nezami S, Khoramshahi E, Nevalainen O, Pölönen I, Honkavaara E (2020) Tree species
546 classification of drone hyperspectral and RGB imagery with deep learning convolutional
547 neural networks. *Remote Sensing* 12(7):1070

548 Ong SQ, Ahmad H, Majid AHA. (2021). Development of a deep learning model from breeding
549 substrate images: a novel method for estimating the abundance of house fly (*Musca*
550 *domestica* L.) larvae. *Pest management science*, 77(12), 5347-5355.

551 Ong SQ, Nair G, Yusof UK, Ahmad H (2022) Community-based mosquito surveillance: an
552 automatic mosquito-on-human-skin recognition system with a deep learning algorithm.
553 *Pest Management Science* 78(10):4092–104. <https://doi.org/10.1002/ps.7028> PMID:
554 35650172

555 Ong SQ, Hamid SA (2022) Next generation insect taxonomic classification by comparing
556 different deep learning algorithms. Plos One 17(12), e0279094.
557 <https://doi.org/10.1371/journal.pone.0279094>

558 Pandong J, Gumal M, Alen L, Sidu A, Ng S, Koh LP (2018) Population estimates of Bornean
559 orangutans using Bayesian analysis at the greater Batang Ai-Lanjak-Entimau landscape in
560 Sarawak, Malaysia. Scientific Reports 10:1–11. <https://doi.org/10.1038/s41598-018-33872-3>

561 Payne J (1988). Orang-utan Conservation in Sabah (Report No. 3754). WWF-Malaysia, Kuala
562 Lumpur. 274 pp.

563 Pearse GD, Watt MS, Soewarto J, Tan AY (2021). Deep learning and phenology enhance large-
564 scale tree species classification in aerial imagery during a biosecurity response. Remote
565 Sensing, 13(9): 1789.

566 Permana AL, Permana JJ, Nellissen L, Prasetyo D, Wich SA, Schaik CPV, Schuppli C (2024)
567 The ontogeny of nest-building behaviour in Sumatran orang-utans, *Pongo abelii*. Animal
568 Behaviour, 211:53-67 .ISSN 0003-3472, <https://doi.org/10.1016/j.anbehav.2024.02.018>.

569 Piel AK, Crunchant A, Knot IE, Chalmers C, Fergus P, Mulero-Pazmany M, Wich SA (2022)
570 Non-invasive technologies for primate conservation in the 21st century. Int J Primatol
571 43:133–167. <https://doi.org/10.1007/s10764-021-00245-z>

572 Purwono P, Ma'arif A, Rahmani W, Fathurrahman HIK, Frisky AZK, ul Haq QM (2022)
573 Understanding of convolutional neural network (cnn): A review. International Journal of
574 Robotics and Control Systems, 2(4), 739-748

575 Rayadin Y, Saitoh T (2009) Individual variation in nest size and nest site features of the Bornean
576 orangutans (*Pongo pygmaeus*). American Journal of Primatology 71(5): 393-399.
577 <https://doi.org/10.1002/ajp.20666>

578 Riedler B, Millesi E, Pratje P (2010) Adaptation to forest life during the reintroduction process
579 of immature *Pongo abelii*. International Journal of Primatology, 31(4): 647-663
580 <https://doi.org/10.1007/s10764-010-9418-2>

581 Saito T, Rehmsmeier M (2015) The precision-recall plot is more informative than the ROC Plot
582 When evaluating binary classifiers on imbalanced datasets. PLOS ONE 10(3): e0118432.
583 <https://doi.org/10.1371/journal.pone.0118432>

584 Simon D, Davies G, Ancrenaz M (2019) Changes to Sabah's orangutan population in recent
585 times: 2002–2017. PLoS ONE 14(7): 1–14

586 Santika T, Wilson K, Meijaard E, Ancrenaz M (2019) The power of mixed survey
587 methodologies for detecting decline of the Bornean orangutan.
588 <https://doi.org/10.1101/775064>

589 Sharma P (2019). Image classification vs. object detection vs. image segmentation. Analytics
590 Vidhya. Available via Online. <https://medium.com/analytics-vidhya/image-classification-vs-object-detection-vs-image-segmentation-f36db85fe81>. Accessed 14 October 2024

591 Sabah Wildlife Department (2020) Orangutan action plan for Sabah 2020-2029. Kota Kinabalu,
592 Sabah, Malaysia.

593

594

595 Smith J, Legg P, Matovic M, Kinsey K (2018) Predicting user confidence during visual decision
596 making. Acm Transactions on Interactive Intelligent Systems 8(2), 1-30.
597 <https://doi.org/10.1145/3185524>

598 Teguh R, Maleh IMD, Sahay AS, Pratama MP, Simon O. Object detection of the Bornean
599 orangutan nests using drone and YOLOv5. Int J Artif Intell ISSN, 2252(8938), 1641.

600 Wang H, Lei Z, Zhang X, Zhou B, Peng J (2016). Machine learning basics. Deep learning, 98-
601 164.

602 Wich S, Dellatore D, Houghton M, Ardi R, Koh L P (2015) A preliminary assessment of using
603 conservation drones for Sumatran orangutan (*Pongo abelii*) distribution and density.
604 Journal of Unmanned Vehicle Systems, 4(1): 45–52. <https://doi.org/10.1139/juvs-2015-0015>

606 Wich SA, Koh LP (2018) Conservation Drones: Mapping and Monitoring Biodiversity. Oxford
607 University Press.

608