


RESEARCH ARTICLE

A convolutional neural network for detecting sea turtles in drone imagery

Patrick C. Gray¹  | Abram B. Fleishman² | David J. Klein² | Matthew W. McKown² |
Vanessa S. Bézy³ | Kenneth J. Lohmann³ | David W. Johnston¹

¹Division of Marine Science and Conservation, Nicholas School of the Environment, Duke University Marine Laboratory, Beaufort, North Carolina

²Conservation Metrics, Inc., Santa Cruz, California

³Department of Biology, University of North Carolina at Chapel Hill, Chapel Hill, North Carolina

Correspondence

Patrick C. Gray

Email: patrick.c.gray@duke.edu

Funding information

National Science Foundation, Grant/Award Number: IOS-1456923

Handling Editor: Oscar Gaggiotti

Abstract

1. Marine megafauna are difficult to observe and count because many species travel widely and spend large amounts of time submerged. As such, management programmes seeking to conserve these species are often hampered by limited information about population levels.
2. Unoccupied aircraft systems (UAS, aka drones) provide a potentially useful technique for assessing marine animal populations, but a central challenge lies in analysing the vast amounts of data generated in the images or video acquired during each flight. Neural networks are emerging as a powerful tool for automating object detection across data domains and can be applied to UAS imagery to generate new population-level insights. To explore the utility of these emerging technologies in a challenging field setting, we used neural networks to enumerate olive ridley turtles *Lepidochelys olivacea* in drone images acquired during a mass-nesting event on the coast of Ostional, Costa Rica.
3. Results revealed substantial promise for this approach; specifically, our model detected 8% more turtles than manual counts while effectively reducing the manual validation burden from 2,971,554 to 44,822 image windows. Our detection pipeline was trained on a relatively small set of turtle examples ($N = 944$), implying that this method can be easily bootstrapped for other applications, and is practical with real-world UAS datasets.
4. Our findings highlight the feasibility of combining UAS and neural networks to estimate population levels of diverse marine animals and suggest that the automation inherent in these techniques will soon permit monitoring over spatial and temporal scales that would previously have been impractical.

KEYWORDS

convolutional neural networks, deep learning for ecology, marine megafauna, marine population monitoring, object detection, sea turtles, unoccupied aircraft systems

1 | INTRODUCTION

Accurate and efficient population estimates are crucial for ecological studies and wildlife management (Cohen, Jonsson, & Carpenter,

2003; Krebs, 1978). For many marine megafauna species, these data are difficult to collect because the animals spend much of their time under water, move rapidly over large areas and occupy remote habitats. As a result, aerial surveys are commonly used to

collect population data for these largely inaccessible species, and in recent years, researchers have turned to unoccupied aircraft systems (UAS, or drones) for these tasks (Johnston, 2019). Surveying populations using UAS can be less logistically challenging than traditional methods, and can also reduce costs and human risk (Arona, Dale, Heaslip, Hammill, & Johnston, 2018) without sacrificing data quality (Hodgson et al., 2018; Johnston et al., 2017). Such surveys have been successfully undertaken with a number of animals, including dugongs (Hodgson, Kelly, & Peel, 2013), seals (Johnston et al., 2017; Seymour, Dale, Hammill, Halpin, & Johnston, 2017), sea turtles (Sykora-Bodie, Bezy, Johnston, Newton, & Lohmann, 2017) and several seabird species (Hodgson, Baylis, Mott, Herrod, & Clarke, 2016).

Globally, six of the seven marine turtle species are listed on the IUCN Red List of Threatened Species under various categories of extinction risk. Estimating the abundance of sea turtle populations is important for conservation efforts, as is developing robust estimates of density in specific breeding, foraging and nesting areas where negative interactions may occur (James, Ottensmeyer, & Myers, 2005). This may be especially true for species like olive ridley sea turtles that exhibit mass nesting, and which aggregate in extraordinarily dense concentrations in coastal areas. While UAS-based methods can facilitate these population assessments (Rees et al., 2018), an essential part of surveys is analysing the resulting images and videos to determine the number of turtles present. Until now, analyses of this type have typically been carried out by trained observers who carefully view each image and count the number of animals present (Sykora-Bodie et al., 2017), but because these analyses are time-consuming and labour-intensive, they place a significant constraint on UAS surveys.

One possible way to overcome problems with analysing images from drones is to automate methods for the detection, localization and enumeration of target animals. Computer vision techniques have the potential to greatly increase the efficiency, repeatability and precision of image assessments and overcome bottlenecks posed by large imagery datasets (Weinstein, 2017). Indeed, a variety of computer vision and machine learning techniques have been applied to assess wildlife populations using data collected not only by UAS imagery (Seymour et al., 2017), but also by camera traps (Schneider, Taylor, & Kremer, 2018; Weinstein, 2018), traditional aerial imagery (Chabot, Dillon, & Francis, 2018) and satellites (Fretwell, Staniland, & Forcada, 2014; Lynch & Schwaller, 2014; Moxley et al., 2017).

Several of these studies have applied modern object-based image analysis and conventional machine learning methods, but interest in deep-learning techniques for more complex detection of specific objects within images and video has been growing. Convolutional neural networks (CNNs), a prominent category of deep-learning classifier inspired by the neural connections in the human brain, are a fundamental source of recent computer vision advances and allow efficient discrimination of objects in noisy and complex environments (Lecun, Bengio, & Hinton, 2015). Although CNN models have typically been applied to large-scale computer vision and image recognition problems, such as efficiently differentiating millions of images into thousands of classes of objects from standardized image

libraries such as ImageNet (Krizhevsky, Sutskever, & Hinton, 2012), they have also been applied to other domains such as image and video data collected for ecological analysis.

Camera traps have been an early testbed for CNNs applied to ecological data. Camera traps are easy to operate and generate high-resolution imagery, but typically collect many unwanted frames due to false camera triggering. CNNs have shown considerable ability to detect objects of interest, such as birds and mammals, from these sources (Gupta & Verma, 2018; Schneider et al., 2018; Yousif, He, & Kays, 2018). Progress has also been made in CNN-based detection of animals in UAS imagery. For example Borowicz et al. (2018) successfully counted Adélie penguins using a CNN (Szegedy et al., 2015). Beyond enumeration, CNN approaches can be used to identify between many species. One recent study was capable of differentiating nearly 600 common North American bird species with only 4% error rates in classification (Van Horn et al., 2015). Additionally, these methods can potentially be applied in the acoustic realm, inasmuch as a CNN exists for monitoring a population of bats through automated detection of their echolocation signals (Aodha et al., 2018). The reduction in required human effort from applying these systems can be considerable. Norouzzadeh et al. (2017) found their CNN could identify animals in 99.3% of the 3.2 million image Snapshot Serengeti dataset with the same accuracy as their crowdsourced identifications, saving volunteers the equivalent of 8.4 years of human labelling effort.

Although promising, CNNs can be prohibitively complex to implement. Moreover, they are computationally intensive and may require more data than is practical for most ecological studies. For example, Merlin used upwards of 50,000 images generated through human annotation to distinguish among bird species (Van Horn et al., 2015). Although details vary widely across studies, and although mitigation strategies such as transfer learning do exist, the application of these techniques to real-world problems clearly poses significant technical challenges.

In this study, we explored the feasibility of using sophisticated yet accessible deep-learning techniques to increase the efficiency of an aerial-image-based population assessment for sea turtles. Specifically, we used a CNN to detect and enumerate olive ridley sea turtles in UAS-generated imagery from at-sea surveys conducted during a mass-nesting event in Ostional, Costa Rica. To our knowledge, this is the first use of CNNs for detecting sea turtles in aerial imagery and demonstrates the broad applicability of combining UAS-based data collection with neural networks for monitoring populations of marine animals.

2 | MATERIALS AND METHODS

2.1 | Study area

At-sea surveys of marine turtles were conducted in nearshore (<3 km from land) waters of the Pacific within the marine protected area at the Ostional National Wildlife Refuge on the Nicoya Peninsula of Costa Rica (Figure 1). The refuge extends 200 m inland from the high

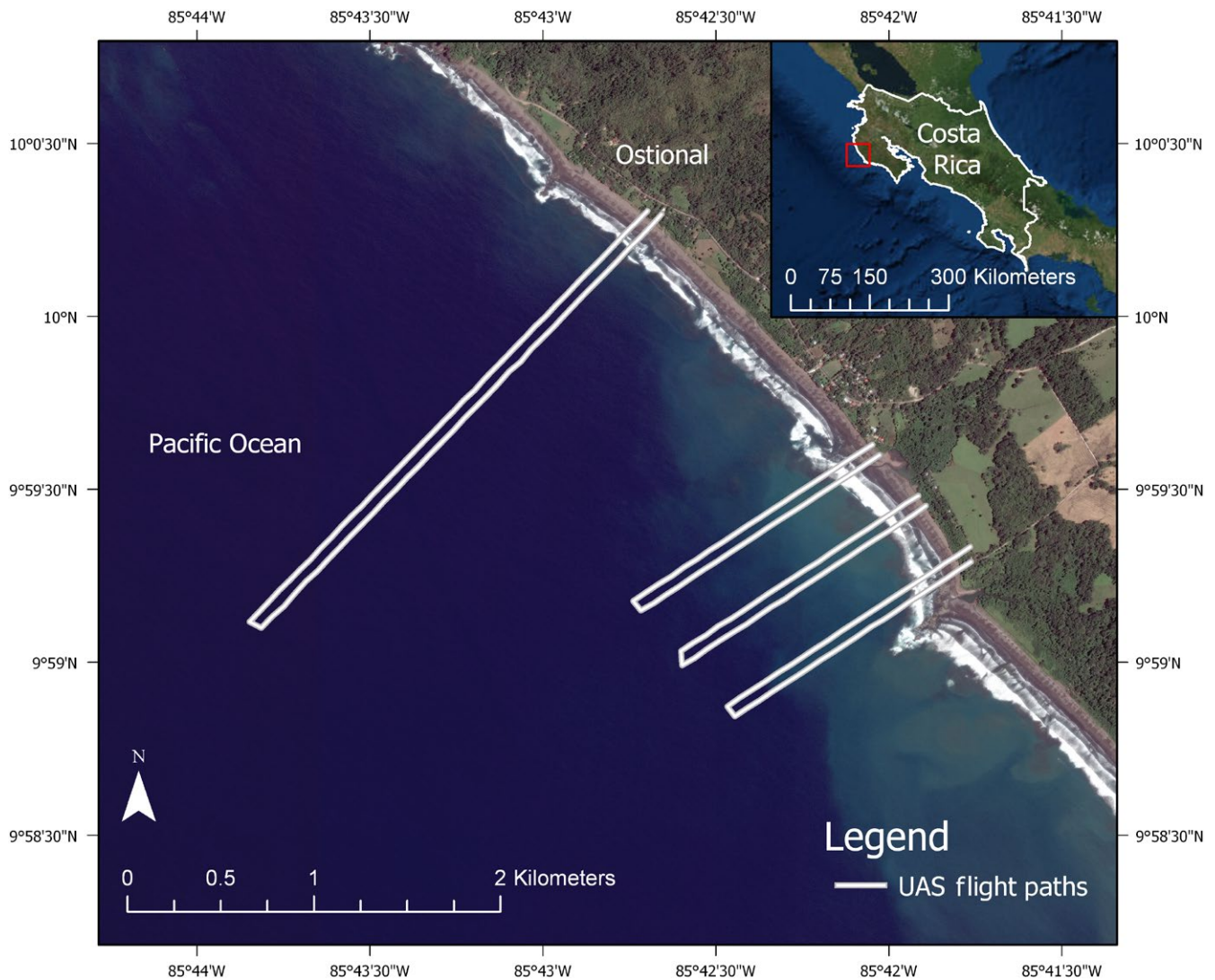


FIGURE 1 Map of the study site at Ostional, Costa Rica with an overview of unoccupied aircraft system (UAS) flight paths. Imagery Sources: Esri, DigitalGlobe

tide line and approximately 5.5 km offshore. Mass-nesting events of olive ridley sea turtles occur at Ostional Beach almost every month of the year. Peak nesting season coincides with the rainy season (May–November).

2.2 | UAS imagery and collection

Aerial surveys were conducted using an eBee (senseFly SA) fixed wing UAS, a modular UAS constructed of light-weight foam and powered by a single electric motor in push configuration (Sykora-Bodie et al., 2017). The UAS was outfitted with a Canon PowerShot S110 near-infrared (NIR) camera to capture aerial photographs. Initial tests with NIR and traditional red-green-blue imagery revealed that NIR imagery provides superior contrast that facilitates detection of turtles in surface waters. Images were collected during transects designed for estimating turtle densities in nearshore waters (Figure 1, and see Sykora-Bodie et al., 2017). Flights were conducted opportunistically during daylight hours, regardless of tidal state or sun angle. A total

of 20 UAS flights were conducted along four transects perpendicular to the beach (five flights per transect) during August 6, 7, 8 and 9, generating a series of overlapping false-colour NIR jpeg images ($N = 1,059$, 12.1 megapixel, pixel dimensions = $4,048 \times 3,048$).

2.3 | Image processing and human counts of turtles

To generate the dataset of labelled turtle locations, 467 of the 1,059 UAS images were used. Using iTag (<https://sourceforge.net/projects/itagbiology/>; version 0.6), three independent reviewers tagged turtles in the image set, taking approximately 6 hr for each reviewer to go through all of the images (Sykora-Bodie et al., 2017). A subset of these manually reviewed images ($N = 275$) were used for model training and the remainder ($N = 192$) were used for direct comparison of counts between manual review and CNN detection. As described previously in Sykora-Bodie et al. (2017), reviewers used pre-set identification criteria to assign each possible turtle into one of two categories, “certain” or “probable.” Certain turtles were

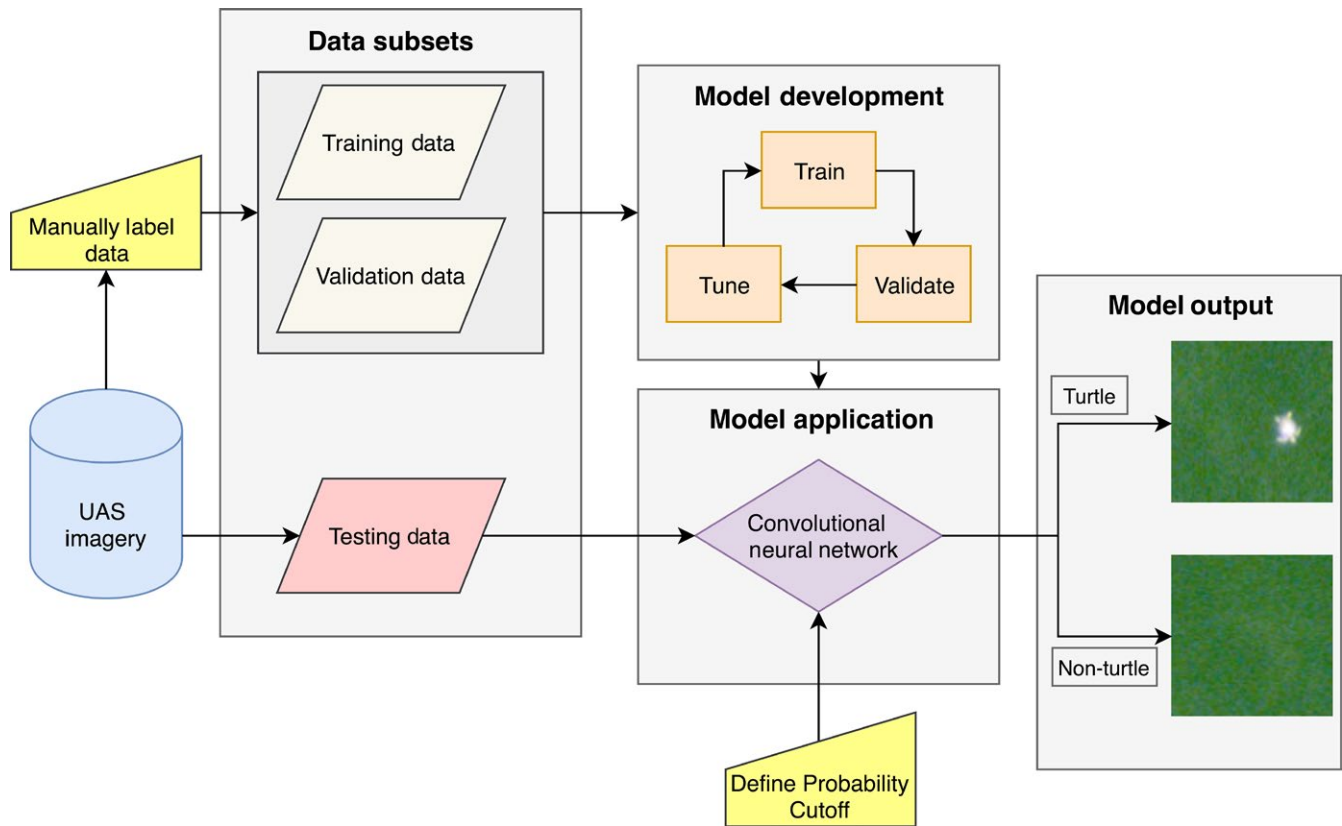


FIGURE 2 Overview of convolutional neural network development, application and output as applied to drone imagery of olive ridley turtles in the coastal waters of Ostional, Costa Rica

those in which appendages were visible and definitive identification was possible. Probable turtles were objects that resembled turtles in size, shape and colour, but could not be identified with certainty. For each photograph, the number of certain and probable turtles was determined. The location (x, y pixel coordinates) of each known or putative turtle within the image was also recorded in iTag. The GPS coordinates of each image were collected in WGS84 and were transformed into a Universal Transverse Mercator (UTM) Zone 16N projection, which is a square grid with constant distances in metres. Pixel locations in each image were converted to UTM coordinates for each turtle detection. After counts were completed, turtle detections were compiled into a text file with (x, y) coordinates of the centre of each identified object. The subset used for training had $N = 616$ certain turtle labels and $N = 328$ probable turtle labels. The subset used for comparison between manual counts and CNN detection had $N = 384$ certain turtle labels and $N = 253$ probable turtle labels. For training and validation, data from these two categories were grouped together into a single class to facilitate automated assessments of detection.

2.4 | Deep-learning model

An overview of the workflow applied to imagery acquired by the UAS is provided in Figure 2. Briefly, the imagery dataset was first partitioned into two components. One was manually labelled while the other was

reserved for testing the model. The labelled data were used to train and validate the CNN and assess its initial performance. Once the CNN was trained, the testing data were run-through the CNN to detect and enumerate turtles imaged during drone flights. Aspects of the CNN deployed in this workflow are visually expanded in Figures 3 and 4 below.

2.5 | Data input and cleaning

A visual inspection of turtles in the images revealed that each turtle could fit within a 50×50 pixel spatial region. Based on this observation, each image was decomposed into an array of 100×100 pixel windows, with 1/3 overlap in the x and y directions ($N = 2,806$ windows per image). The overlap was chosen so that turtle centres would fall within or near the interior region of only one window, to minimize instances of missed or duplicated turtle detections because of window centring. A binary classification approach was used in which these 100×100 pixel windows were denoted with the positive class (1) if the central 50×50 region contained the centre of a turtle. If a window did not contain a turtle in the centre, it was denoted with the negative class (0).

2.6 | Neural network architecture and training

Each training example for the CNN was comprised of a $100 \times 100 \times 3$ tensor ($x \times y \times \text{Green, Red, Near Infrared}$) and the accompanying label

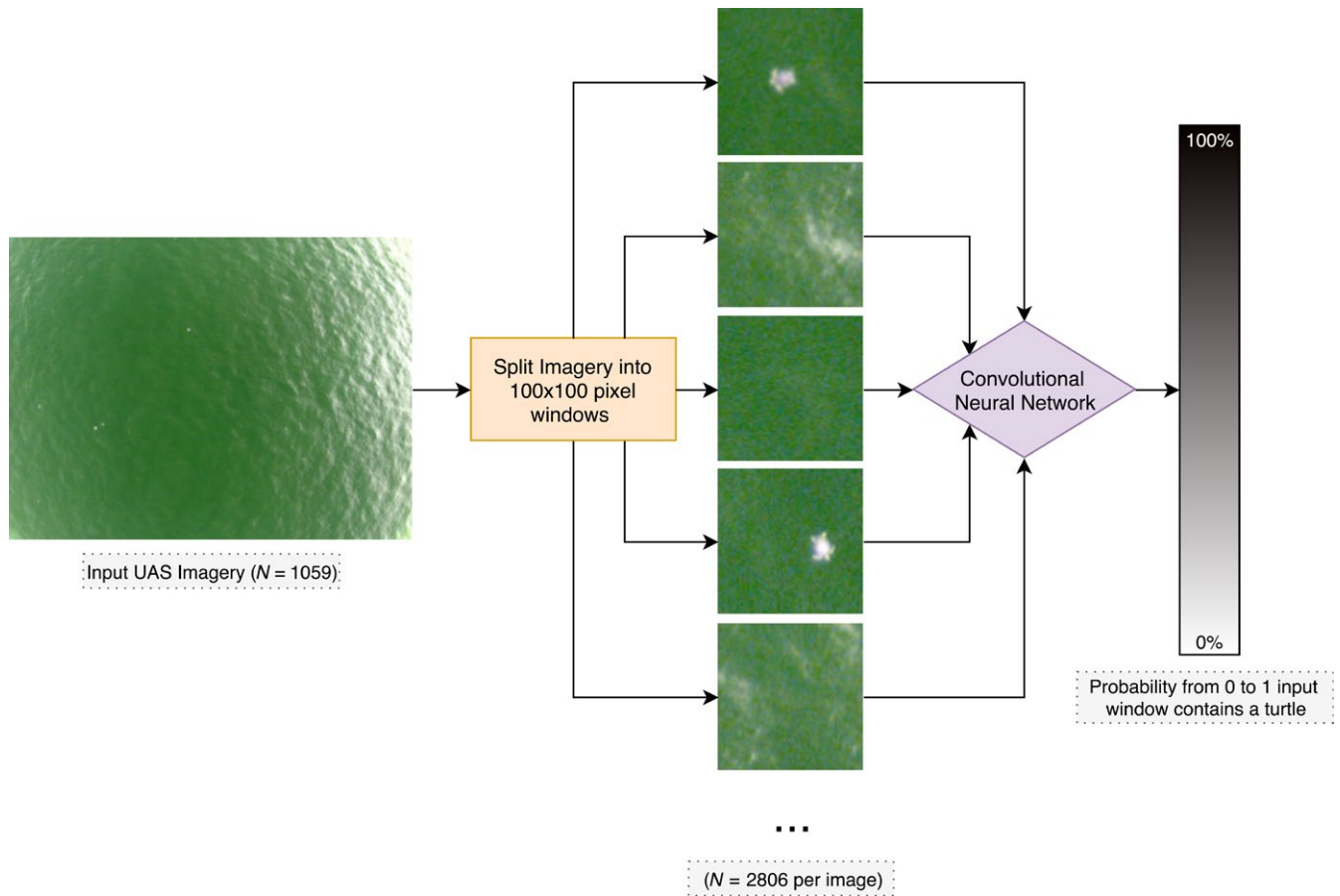


FIGURE 3 Overview of convolutional neural network (CNN) “Model Application” from Figure 2. Unoccupied aircraft system (UAS) images were subdivided into 100×100 pixel windows, which served as input to the CNN. The CNN then computed a probability that each window included a turtle

classifying the example as negative or positive (0 or 1 respectively). Examples of the positive class ($N = 944$) were assembled by cropping the images around each labelled $(x + dx, y + dy)$ turtle coordinate, where the variables (dx, dy) have a uniform random distribution between -25 and $+25$. Thus, the positive examples all had turtle centres randomly distributed in the interior 50×50 pixel spatial region. This was done to simulate the random positioning of turtles within windows extracted from new images in deployment. Although there were no instances of multiple turtles within the centre of a 50×50 pixel region in the training data, this approach would count multiple turtles in that 50×50 pixel region as one. Using a class balance of 10 negative examples for each positive example in our training data, examples of the negative class ($N = 9,440$) were assembled by cropping random 100×100 pixel windows from the image library, provided the centre of the window was at least 125 pixels away from the centre of a labelled turtle (in both the x and y directions). This led to a total of 10,384 labelled examples. To train the CNN, a random 85% of the positive and negative examples were chosen as training data and the remaining 15% used as validation data.

Given the training dataset was comprised of only $N = 8,826$ examples, a CNN of modest size was employed (Figure 4). For detailed information on general CNN architecture and training, see

Lecun et al. (2015) for a technical yet cogent overview. The CNN for this study was comprised of four convolutional layers with sixty-four 3×3 kernels interleaved with max pooling layers. These convolutional layers form the backbone of the CNN and slide over the image essentially outputting heatmaps of various features within the image, called feature maps. Our CNN with 64 kernels will output 64 feature maps at each convolutional layer. The interwoven max pooling layers slide over the convolutional layer's feature maps with a 2×2 window and output the maximum value in that window. Important features from the maps are typically still retained by this simple max operation while the overall size of the feature maps is reduced with each iteration. Thus, the CNN keeps the same effective “field of view” while considerably reducing dimensionality. Including multiple iterations of interleaved convolutional and max pooling layers, developing a smaller feature map each run-through, permits a CNN to ingest noisy and variable images, find useful features within them and condense that into a relatively small yet informative final feature map. The first layer of a CNN typically creates maps of features such as edges, curves and colour gradients. The feature maps created in deeper layers in a CNN are more abstract and aggregate the previous layer's feature maps; in our case, combining them into groups of curves and

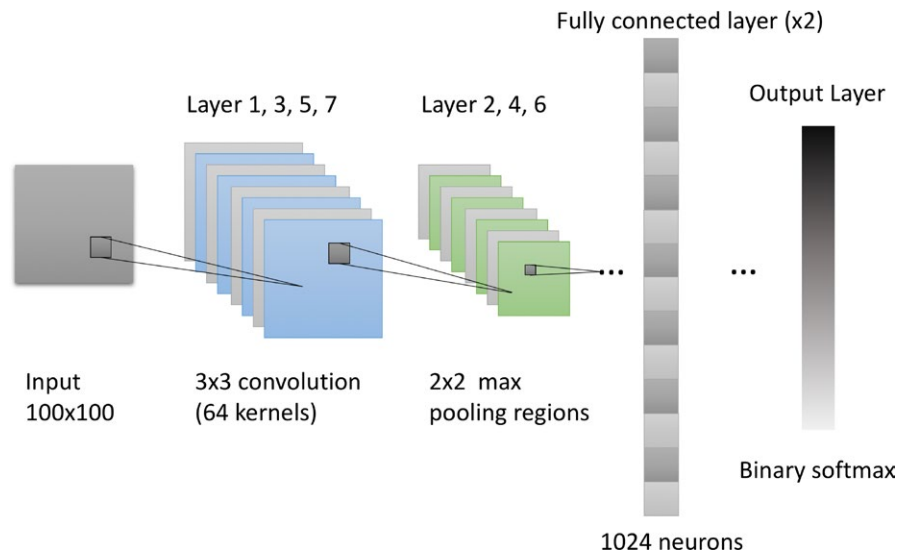


FIGURE 4 Overview of the convolutional neural network (CNN) architecture. The CNN had four convolutional layers alternating with max pooling layers, these layers perform the feature extraction for the CNN, effectively distinguishing which aspects of the image are informative for classification. These layers were followed by two fully connected layers of 1,024 neurons which combine the previously extracted features into meaningful combinations that ideally provide some predictive power for classification. The final layer employed a binary normalized exponential (softmax) function which ingests the final fully connected layer and its learned combinations of features, and returns a value between 0 and 1, with higher values signalling higher confidence of a turtle in the image window

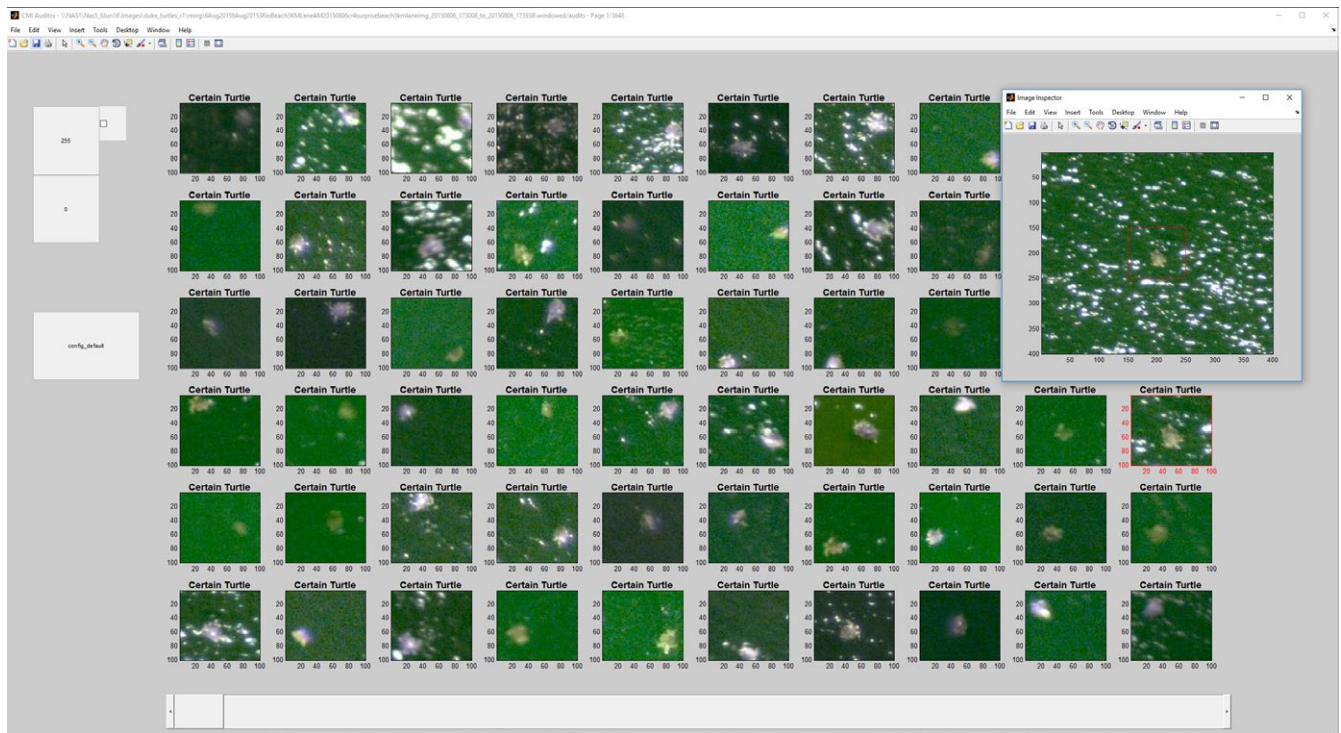


FIGURE 5 Screenshot of the custom application used to rapidly review turtle detections

edges that may indicate turtle flippers or shells. Through this process, the CNN extracts the distinguishing features that will permit effective classification. This process of feature extraction was followed by two fully connected layers of 1,024 neurons. In order to prevent overfitting of the model, 50% of the connections between the neurons of these two fully connected layers were randomly ignored (dropout ratio of 0.5) during training.

The fully connected layers take the final 64 feature maps, ideally representing useful and high-level image components, and learn a mapping from those feature maps to the output classes (turtle–non-turtle). We used a binary normalized exponential (softmax) function output layer, which takes the final, fully connected layer and provides a continuous value between 0 and 1 for each window. Values closer to 1 indicate a higher likelihood that a turtle is present in the central

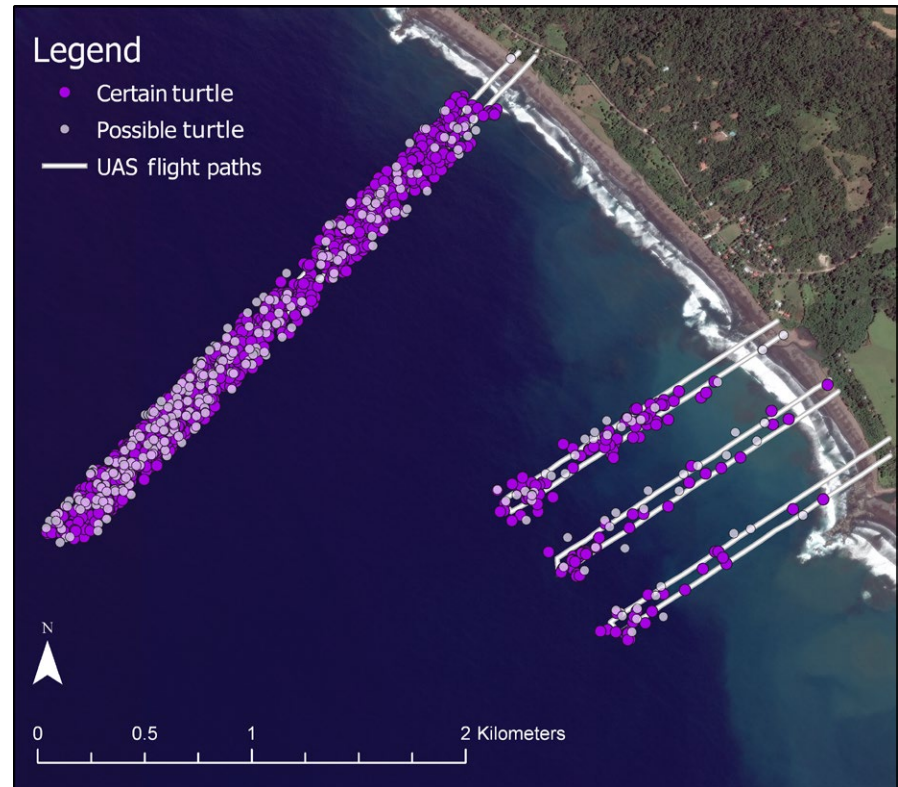


FIGURE 6 Distribution of certain and probable turtle detections from the convolutional neural network after the output was manually validated. Unoccupied aerial system (UAS) flight paths show distribution of turtle detections along survey transects

50 × 50 region of the window. The CNN was trained to optimize the performance of the classification model (via a categorical cross-entropy, or log loss function) using stochastic gradient descent, with a learning rate of 0.01 and a Nesterov momentum of 0.9. Training was continued over 20 epochs, after which the model parameters were selected from the epoch resulting in the lowest loss on the validation set of $N = 1,558$ ($N = 153$ of which contained turtles).

2.7 | Model validation

Detections were reviewed with an internally developed custom user interface, allowing quick review of the 100 × 100 pixel windows, sorted by the CNN's output turtle likelihood (Figure 5). The tool uses a variety of keyboard shortcuts to note and correct instances of duplicate detections, false positives or any false negatives encountered.

Because the windows overlapped by 30%, we detected some turtles more than once. When there was more than one detection of a turtle, we selected the window that included a larger proportion of the turtle or the turtle closest to the centre of the photo. Double detections were removed manually by sorting true detections in photo capture order and marking duplicates with a new label.

3 | RESULTS

3.1 | Model performance

The overall accuracy of our model to detect turtles was 99.83%, however, overall accuracy is not the best metric in classification problems with a sparse class of interest such as our study species,

where the vast majority of image windows do not contain turtles. Instead, we present recall and precision as metrics of model performance. Recall describes the number of true positives from the model divided by the total number of true positives in the data, or simply, what proportion of true positives were detected. Precision describes the number of true positives from the model divided by the total number of detections by the model, again more simply, what proportion of detections were actually correct. After training and validation, the model was tuned to prioritize recall in order to miss as few true positives as possible, at the cost of additional false positives, with a detection probability threshold of 0.93. Out of 2,971,554 total windows (1,059 images each broken into 2,806 windows of 100 × 100 pixels), at this threshold the model flagged 6,696 windows as having turtles. Of those 6,696, we manually reviewed all windows and found there were 874 with certain turtles, 218 with probable turtles and 5,402 with no turtles (Figure 6). Recall and precision metrics were calculated by combining the certain and probable turtle counts (Table 1). Accuracy, precision and recall were quantified without manually verifying all 2,971,554 windows by reviewing all windows with a detection probability of 0.93–0.05 ($N = 44,799$). This lower bound was determined empirically, while turtles could still exist in windows below that number, we found negligible true positives below a probability threshold of 0.15 and made the assumption that the total number of turtles detected, as well as our precision and recall statistics, would not change materially with further review below the 0.05 probability threshold. This assumption prevented the necessity to manually review the remaining 2,926,732 windows below the 0.05 threshold. At this lower threshold, 199 additional certain turtles and 137 probable turtles

TABLE 1 Confusion matrix for turtle detection results when the model is tuned to show detections above the 0.93 confidence threshold

	Predicted	
	Turtle	Non-turtle
Validated		
Turtle	1,092	336
Non-turtle	5,604	2,964,522 ^a

Note. Shaded cells represent windows that were correctly classified (true positives and true negatives).

^aValidated non-turtle numbers are based on manually verifying windows down to a 0.05 probability threshold.

were found, and 44,463 windows with no turtles. Assuming these are all of the turtles within our dataset, this leads to 16.3% precision and 76.5% recall when using the 0.93 probability threshold for this model (Figure 7). The recall for our testing dataset closely matched the validation data (Figure 7a) while precision was poorer in the testing dataset (Figure 7b) suggesting slight model overfitting.

3.2 | CNN vs. manual counts

Using the same review method that generated the training data, a comparison of verified CNN detections to manual counts shows

similar results but marginally better detection capability, approximately 8%–9%, for both certain and probable turtle detections with the CNN (Table 2).

3.3 | Model validation effort

It took 1 hr of analyst time to review detections from the 0.93 threshold (6,696 windows). It took an additional 11.8 hr to review all detections between probability thresholds of 0.93 and 0.05 (44,799 windows).

4 | DISCUSSION

Our study represents the first use of deep-learning methods to assess at-sea densities of sea turtles. Results of the CNN analysis are similar to manual counts of the same imagery in a previous density assessment (Sykora-Bodie et al., 2017). The CNN approach identified 8%–9% more turtles (Table 2), suggesting previous assessment should be viewed as conservative. Overall, our results demonstrate the feasibility of using neural networks to facilitate the analysis of images acquired for the purpose of monitoring animal populations. The general approach described here can also be applied to aerial surveys for other species large enough to be detected in coastal or pelagic marine environments. This study illustrates how CNNs can

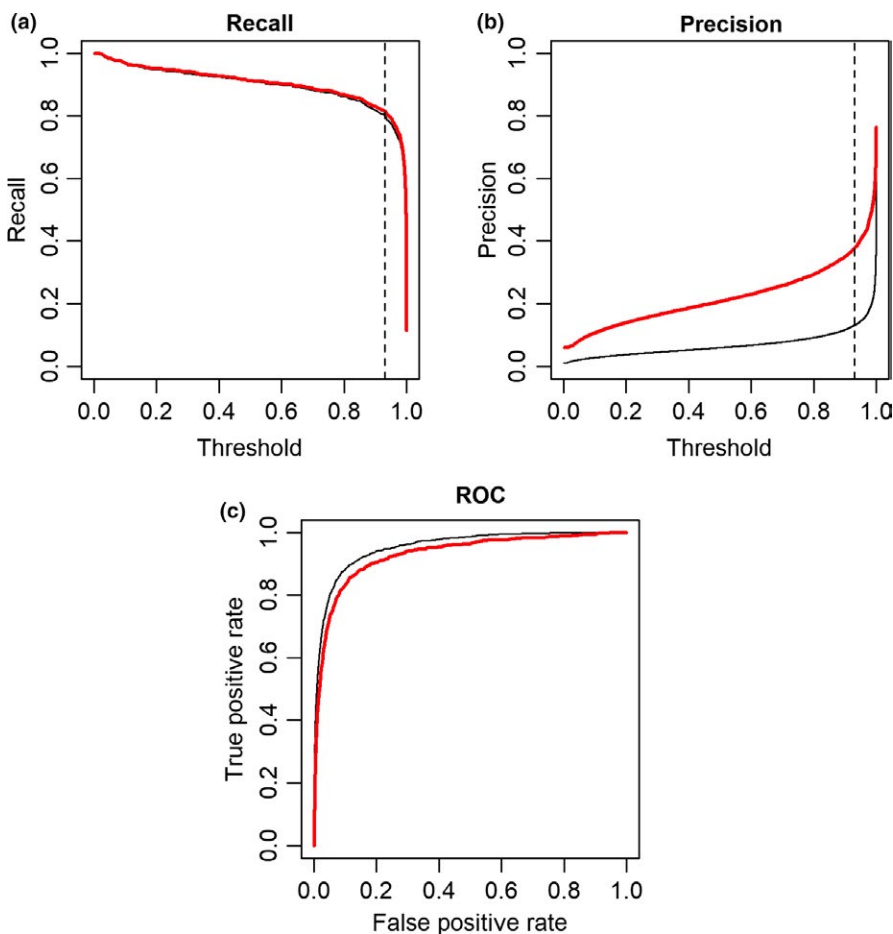


FIGURE 7 Convolutional neural network (a) recall, (b) precision, and (c) precision-recall curve for training (red) and full testing (black) data. Recall is the number of true positives from the model divided by the total number of true positives in the data. Precision is the number of true positives from the model divided by the total number of detections by the model. These two values represent model performance. Our model was tuned to optimize recall at the cost of precision

TABLE 2 A comparison of manual turtle counts, detections from the convolutional neural network (CNN), and the per cent difference from manual count to the CNN-based method

	Manual	CNN	% Difference
Certain	384	418	+8.9
Probable	253	274	+8.3
Total	637	692	+8.6

facilitate efforts to enumerate sea turtles in drone imagery and eliminate analyst biases introduced during manual counts.

In many conservation studies, researchers must balance the risks of committing Type I and Type II errors. Most statistical analyses of scientific data focus on reducing Type I errors in an effort to maximize the ability to reject null hypotheses with confidence. In conservation biology, however, the consequences of committing a Type II error are often far worse (Shrader-Frechette, 1994). Thus, researchers typically adopt approaches that minimize the likelihood of making Type II errors. When applying a machine learning approach to counting species at risk, researchers can tune the system to minimize either Type I or Type II errors by balancing recall (the number of true positives from the model divided by the total number of true positives in the data) and precision (the number of true positives from the model divided by the total number of detections by the model).

In the context of this study, our model reduced the manual analysis burden to 1.5% of the initial amount using a 0.93 probability threshold. This was achieved without a diminished ability to detect turtles, in fact the model identified more turtles in a subset of the imagery than were identified using manual counts. While many deep-learning systems must reach predictably high precision and recall (e.g., detecting a stoplight in a vision system for a self-driving car), models applied to ecological data can be tuned to achieve alternative outcomes and can be useful even when they have low precision. Our model probability threshold was intentionally set low, minimizing the likelihood of committing Type II errors (i.e., failing to detect a turtle when it is present) when assessing images. Here, we accept the financial cost of marginally more analyst time to review detections in an effort to generate the most robust density information for conservation purposes.

Although our model's precision is likely acceptable for many conservation applications, our testing data yielded lower precision than did the training and validation data (Figure 7). This difference is likely due to some overfitting of the CNN on the training data, perhaps caused by the limited number of training samples and differences between the training and testing datasets, which were captured on different days, under different weather and sun angle conditions and contained some land images which were not in the training dataset. CNNs are prone to overfitting in circumstances of limited training data and if the network is too large for the given relationships it is attempting to model (Srivastava, Hinton, Krizhevsky, Sutskever, & Salakhutdinov, 2014). Our training sample class imbalance (10:1 negative to positive), while improving precision and recall on the validation set, may have also decreased precision in our testing data.

A larger training dataset with examples of more varied image conditions might help mitigate this issue. The main sources of false detections in this study were large jellyfish, breaking waves and glare. Planning flights to reduce glare (e.g., during periods of lower sun angle), as well as image processing prior to CNN application, could substantially reduce these errors. For this study, our goal was to build a robust CNN, able to function well under variable conditions, in order to avoid imposing additional constraints on drone flights.

Major benefits of CNNs include: (a) they can be applied across images in different conditions, and (b) the model can be repeatedly upgraded using additional training data (Oquab, Bottou, Laptev, & Sivic, 2014). Thus, when new imagery is acquired during future deployments, the performance of the model can be improved by following the same initial process (lowering the probability threshold and manually inspecting detections), then adding the additional images of verified turtles into a new, larger training dataset. This iterative approach is likely to be useful across many similarly noisy and imbalanced datasets, increasing precision, improving detection of rarer image variants and saving considerable time relative to the brute-force manual inspection of all data.

The next iteration of this study's CNN could be trained to identify false positive generating classes in addition to the primary class of interest. This has been shown to force the neural network to better separate these objects within its internal representation and thus reduce false positives (A. B. Fleishman, unpublished data). While not an issue in our study, adjustments should be made to permit detection of multiple turtles within close proximity, given that our current model would count multiple turtles within the 50×50 pixel region as a single turtle, in order to increase usefulness on data with denser aggregations or mating behaviour. Beyond these two relatively minor changes, further understanding how fine-tuning a pre-trained CNN compares to our network built from scratch will be beneficial. Fine-tuning a CNN is the process of beginning with a pre-trained model, trained on a large dataset such as ImageNet (Deng et al., 2009) or the Common Objects in Context dataset (Lin et al., 2014), and then adding in additional training samples specific to the desired detection problem. This process, more generally called transfer learning, has been shown to allow fewer training samples (Razavian, Azizpour, Sullivan, & Carlsson, 2014) and reduce overfitting (Yosinski, Clune, Bengio, & Lipson, 2014), thus increasing the potential applicability of this method in real-world biology scenarios (Schneider et al., 2018). If the CNNs from our study were first trained on ImageNet, which would build out many of the class agnostic image process capabilities, such as edge detectors, and then trained on our turtle dataset, we hypothesize there would be a noticeable improvement in precision. A final avenue of inquiry is comparison of this custom-designed CNN architecture to models currently available open source and "off the shelf" that provide state-of-the-art results in other domains (Lin, Goyal, Girshick, He, & Dollár, 2017; Ren, He, Girshick, & Sun, 2017).

Open-source CNN implementations, combined with transfer learning, will allow CNNs to be improved upon with additional data

and fed back into the community, considerably enhancing our ability to detect and study wildlife. Improvements in CNN speed and efficiency may eventually permit detectors to run in real-time on UAS (Huang et al., 2017), facilitating autonomous monitoring and behavioural analysis. In order to build out these capabilities, we advocate for the creation of appropriately sized open-access training datasets of aerial imagery for all sea turtle species in various conditions, permitting rapid creation and improvement of CNNs for sea turtle population monitoring globally.

ACKNOWLEDGEMENTS

We thank Seth Sykora-Bodie for providing details on his study on UAS-based sea turtle population assessment. Thanks to Everette Newton for conducting flights and contributing to manual counts. Thanks to Walter Wright and Matt Elardo for manual counting of turtles. We greatly appreciate SINAC and employees at the Ostional National Wildlife Refuge for logistical support and our research being permitted by the Costa Rican government. This research was funded by the National Science Foundation grant (IOS-1456923) to K. J. Lohmann.

AUTHORS' CONTRIBUTIONS

M.W.M. and D.W.J. conceived the ideas and designed methodology; V.S.B. and K.J.L. collected the data; A.B.B. and D.J.K. analysed the data; P.C.G. led the writing of the manuscript. All the authors contributed critically to the drafts and gave final approval for publication.

CONFLICT OF INTEREST

The authors declare no conflicts of interest.

DATA ACCESSIBILITY

The authors have uploaded the convolutional neural network code to Github.com at url: https://github.com/patrickcgray/cnn_sea_turtle_detection and <https://doi.org/10.5281/zenodo.1973808>. All relevant imagery, including training and testing labels, has been stored on the Dryad Data Repository: <https://doi.org/10.5061/dryad.5h06vv2> (Gray et al., 2019).

ORCID

Patrick C. Gray  <https://orcid.org/0000-0002-8997-5255>

REFERENCES

- Aodha, O. M., Gibb, R., Barlow, K. E., Browning, E., Firman, M., Freeman, R., ... Jones, K. E. (2018). Bat detective—Deep learning tools for bat acoustic signal detection. *PLoS Computational Biology*, 14, e1005995. <https://doi.org/10.1371/journal.pcbi.1005995>
- Arona, L., Dale, J., Heaslip, S. G., Hammill, M. O., & Johnston, D. W. (2018). Assessing the disturbance potential of small unoccupied aircraft systems (UAS) on gray seals (*Halichoerus grypus*) at breeding colonies in Nova Scotia, Canada. *PeerJ*, 6, e4467. <https://doi.org/10.7717/peerj.4467>
- Borowicz, A., McDowall, P., Youngflesh, C., Sayre-McCord, T., Clucas, G., Herman, R., ... Lynch, H. J. (2018). Multi-modal survey of Adélie penguin mega-colonies reveals the Danger Islands as a sea-bird hotspot. *Scientific Reports*, 8, 1–9. <https://doi.org/10.1038/s41598-018-22313-w>
- Chabot, D., Dillon, C., & Francis, C. M. (2018). An approach for using off-the-shelf object-based image analysis software to detect and count birds in large volumes of aerial imagery. *Avian Conservation and Ecology*, 13, 15. <https://doi.org/10.5751/ACE-01205-130115>
- Cohen, J. E., Jonsson, T., & Carpenter, S. R. (2003). Ecological community description using the food web, species abundance, and body size. *Proceedings of the National Academy of Sciences of the United States of America*, 100, 1781–1786. <https://doi.org/10.1073/pnas.232715699>
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., & Fei-Fei, L. (2009). ImageNet: A large-scale hierarchical image database. 2009 *IEEE conference on computer vision and pattern recognition* (pp. 248–255). Miami, FL: IEEE. <https://doi.org/10.1109/cvprw.2009.5206848>
- Fretwell, P. T., Staniland, I. J., & Forcada, J. (2014). Whales from space: Counting southern right whales by satellite. *PLoS ONE*, 9, 1–9. <https://doi.org/10.1371/journal.pone.0088655>
- Gray, P. C., Fleishman, A. B., Klein, D. J., McKown, M. W., Bézy, V. S., Lohmann, K. J., & Johnston, D. W. (2019). Data from: A convolutional neural network for detecting sea turtles in drone imagery. *Dryad Digital Repository*. <https://doi.org/10.5061/dryad.585t4>
- Gupta, P., & Verma, G. K. (2018). Wild animal detection using deep convolutional neural network. In B. Chaudhuri, M. Kankanalli, & B. Raman (Eds.), *Proceedings of 2nd international conference on computer vision & image processing* (pp. 327–338). Singapore: Springer. https://doi.org/10.1007/978-981-10-7898-9_27
- Hodgson, J. C., Baylis, S. M., Mott, R., Herrod, A., & Clarke, R. H. (2016). Precision wildlife monitoring using unmanned aerial vehicles. *Scientific Reports*, 6, 22574. <https://doi.org/10.1038/srep22574>
- Hodgson, A., Kelly, N., & Peel, D. (2013). Unmanned aerial vehicles (UAVs) for surveying Marine Fauna: A dugong case study. *PLoS ONE*, 8, e79556. <https://doi.org/10.1371/journal.pone.0079556>
- Hodgson, J. C., Mott, R., Baylis, S. M., Pham, T. T., Wotherspoon, S., Kilpatrick, A. D., ... Koh, L. P. (2018). Drones count wildlife more accurately and precisely than humans. *Methods in Ecology and Evolution*, 9, 1160–1167. <https://doi.org/10.1111/2041-210X.12974>
- Huang, J., Rathod, V., Sun, C., Zhu, M., Korattikara, A., Fathi, A., ... Murphy, K. (2017). Speed/accuracy trade-offs for modern convolutional object detectors. In *Proceedings—30th IEEE conference on computer vision and pattern recognition, CVPR 2017, 2017-Janua* (pp. 3296–3305). Honolulu, HI: IEEE. <https://doi.org/10.1109/cvpr.2017.351>
- James, M. C., Ottensmeyer, C. A., & Myers, R. A. (2005). Identification of high-use habitat and threats to leatherback sea turtles in northern waters: New directions for conservation. *Ecology Letters*, 8, 195–201. <https://doi.org/10.1111/j.1461-0248.2004.00710.x>
- Johnston, D. W. (2019). Unoccupied aircraft systems in marine science and conservation. *Annual Review of Marine Science*, 11. <https://doi.org/10.1146/annurev-marine-010318-095323>
- Johnston, D. W., Dale, J., Murray, K. T., Josephson, E., Newton, E., & Wood, S. (2017). Comparing occupied and unoccupied aircraft surveys of wildlife populations: Assessing the gray seal (*Halichoerus grypus*) breeding colony on Muskeget Island, USA. *Journal of Unmanned Vehicle Systems*, 5, 178–191. <https://doi.org/10.1139/juvs-2017-0012>
- Krebs, C. (1978). *Ecology: The experimental analysis of distribution and abundance*. New York, NY: Harper and Row. <https://doi.org/10.2307/3545>

- Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). ImageNet classification with deep convolutional neural networks. *Advances in Neural Information Processing Systems*, 1–9. <https://doi.org/10.1016/j.protocy.2014.09.007>
- Lecun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, 521, 436–444. <https://doi.org/10.1038/nature14539>
- Lin, T. Y., Goyal, P., Girshick, R., He, K., & Dollar, P. (2017). Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision* (pp. 2999–3007). <https://doi.org/10.1109/iccv.2017.324>
- Lin, T. Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., ... Zitnick, C. L. (2014). *Microsoft COCO: Common objects in context*. Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics) 8693 LNCS (pp. 740–755). https://doi.org/10.1007/978-3-319-10602-1_48
- Lynch, H. J., & Schwaller, M. R. (2014). Mapping the abundance and distribution of Adélie penguins using landsat-7: First steps towards an integrated multi-sensor pipeline for tracking populations at the continental scale. *PLoS ONE*, 9, 113301. <https://doi.org/10.1371/journal.pone.0113301>
- Moxley, J. H., Bogomolni, A., Hammill, M. O., Moore, K. M. T., Polito, M. J., Sette, L., ... Johnston, D. W. (2017). Google haul out: Earth observation imagery and digital aerial surveys in coastal wildlife management and abundance estimation. *BioScience*, 67, 760–768. <https://doi.org/10.1093/biosci/bix059>
- Norouzzadeh, M. S., Nguyen, A., Kosmala, M., Swanson, A., Palmer, M., Packer, C., & Clune, J. (2017). Automatically identifying, counting, and describing wild animals in camera-trap images with deep learning. *Proceedings of the National Academy of Sciences of the United States of America*, 115, E5716–E5725. <https://doi.org/10.1073/pnas.171936711>
- Oquab, M., Bottou, L., Laptev, I., & Sivic, J. (2014). Learning and transferring mid-level image representations using convolutional neural networks. *Proceedings of IEEE conference on computer vision and pattern recognition*. Columbus, OH: IEEE. <https://doi.org/10.1109/cvpr.2014.222>
- Razavian, A. S., Azizpour, H., Sullivan, J., & Carlsson, S. (2014). CNN features off-the-shelf: An astounding baseline for recognition. *IEEE conference on computer vision and pattern recognition workshops* (pp. 512–519). Columbus, OH: IEEE. <https://doi.org/10.1109/cvprw.2014.131>
- Rees, A., Avens, L., Ballorain, K., Bevan, E., Broderick, A., Carthy, R., ... Godley, B. (2018). The potential of unmanned aerial systems for sea turtle research and conservation: A review and future directions. *Endangered Species Research*, 35, 81–100. <https://doi.org/10.3354/esr00877>
- Ren, S., He, K., Girshick, R., & Sun, J. (2017). Faster R-CNN: Towards real-time object detection with region proposal networks. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 39, 1137–1149. <https://doi.org/10.1109/TPAMI.2016.2577031>
- Schneider, S., Taylor, G. W., & Kremer, S. C. (2018). *Deep learning object detection methods for ecological camera trap data*. In Computer vision and pattern recognition. Retrieved from <http://arxiv.org/abs/1803.10842>
- Seymour, A. C., Dale, J., Hammill, M., Halpin, P. N., & Johnston, D. W. (2017). Automated detection and enumeration of marine wildlife using unmanned aircraft systems (UAS) and thermal imagery. *Scientific Reports*, 7, 45127. <https://doi.org/10.1038/srep45127>
- Shrader-Frechette, K. (1994). *Ethics of Scientific Research* (1st ed.). Lanham, MD: Rowman & Littlefield Publishers.
- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014). Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15, 1929–1958. <https://doi.org/10.1214/12-AOS1000>
- Sykora-Bodie, S. T., Bezy, V., Johnston, D. W., Newton, E., & Lohmann, K. J. (2017). Quantifying nearshore sea turtle densities: Applications of unmanned aerial systems for population assessments. *Scientific Reports*, 7, 17690. <https://doi.org/10.1038/s41598-017-17719-x>
- Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., ... Arbo, A. (2015). Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 1–9). <https://doi.org/10.1109/cvpr.2015.7298594>
- Van Horn, G., Branson, S., Farrell, R., Haber, S., Barry, J., Ipeirotis, P., ... Belongie, S. (2015). Building a bird recognition app and large scale dataset with citizen scientists: The fine print in fine-grained dataset collection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 07–12 June (pp. 595–604). <https://doi.org/10.1109/cvpr.2015.7298658>
- Weinstein, B. G. (2017). A computer vision for animal ecology. *Journal of Animal Ecology*. <https://doi.org/10.1111/1365-2656.12780>
- Weinstein, B. G. (2018). Scene-specific convolutional neural networks for video-based biodiversity detection. *Methods in Ecology and Evolution*, 9, 1435–1441. <https://doi.org/10.1111/2041-210X.13011>
- Yosinski, J., Clune, J., Bengio, Y., & Lipson, H. (2014). How transferable are features in deep neural networks? In *Advances in neural information processing systems* (pp. 1–9). <https://doi.org/10.1109/ijcnn.2016.7727519>
- Yousif, H., He, Z., & Kays, R. (2018). Object segmentation in the deep neural network feature domain from highly cluttered natural scenes. In *Proceedings—International Conference on Image Processing, ICIP* (pp. 3095–3099). <https://doi.org/10.1109/icip.2017.8296852>

How to cite this article: Gray PC, Fleishman AB, Klein DJ, et al.

A convolutional neural network for detecting sea turtles in drone imagery. *Methods Ecol Evol.* 2019;10:345–355.

<https://doi.org/10.1111/2041-210X.13132>