



Frame-by-frame annotation of video recordings using deep neural networks

ALEXANDER M. CONWAY,¹ IAN N. DURBACH ^{1,2,†} ALISTAIR MCINNES ^{3,4} AND ROBERT N. HARRIS⁵

¹Centre for Statistics in Ecology, the Environment, and Conservation, University of Cape Town, Cape Town, South Africa

²Centre for Research into Ecological and Environmental Modelling, University of St Andrews, St Andrews, UK

³Seabird Conservation Programme, BirdLife South Africa, Johannesburg, South Africa

⁴Department of Zoology, DST/NRF Centre of Excellence at the Percy FitzPatrick Institute, Nelson Mandela University, Port Elizabeth, South Africa

⁵Sea Mammal Research Unit, University of St Andrews, St Andrews, UK

Citation: Conway, A. M., I. N. Durbach, A. McInnes, and R. N. Harris. 2021. Frame-by-frame annotation of video recordings using deep neural networks. *Ecosphere* 12(3):e03384. 10.1002/ecs2.3384

Abstract. Video data are widely collected in ecological studies, but manual annotation is a challenging and time-consuming task, and has become a bottleneck for scientific research. Classification models based on convolutional neural networks (CNNs) have proved successful in annotating images, but few applications have extended these to video classification. We demonstrate an approach that combines a standard CNN summarizing each video frame with a recurrent neural network (RNN) that models the temporal component of video. The approach is illustrated using two datasets: one collected by static video cameras detecting seal activity inside coastal salmon nets and another collected by animal-borne cameras deployed on African penguins, used to classify behavior. The combined RNN-CNN led to a relative improvement in test set classification accuracy over an image-only model of 25% for penguins (80% to 85%), and substantially improved classification precision or recall for four of six behavior classes (12–17%). Image-only and video models classified seal activity with very similar accuracy (88 and 89%), and no seal visits were missed entirely by either model. Temporal patterns related to movement provide valuable information about animal behavior, and classifiers benefit from including these explicitly. We recommend the inclusion of temporal information whenever manual inspection suggests that movement is predictive of class membership.

Key words: animal-borne video; automated detection; deep learning; image classification; neural networks; video classification.

Received 21 July 2020; accepted 5 October 2020. Corresponding Editor: Bistra Dilkina.

Copyright: © 2021 The Authors. This is an open access article under the terms of the Creative Commons Attribution License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

† **E-mail:** ian.durbach@uct.ac.za

INTRODUCTION

Technological advances in quality, size, battery life, and storage capacity have enabled video cameras to record more data at better quality on a broader variety of animals, becoming small enough to deploy on numerous animal species (Takahashi et al. 2004, Rutz and Troschianko 2013)

and on drones (Anderson and Gaston 2013, Cruzan et al. 2016), as well as in more conventional fixed locations. Footage captured using video cameras needs to be annotated for use in scientific research, a currently labor-intensive process often involving highly trained scientists manually annotating the content of videos frame by frame. Even with dedicated annotation software,

this presents a major bottleneck for scientific research based on these data, necessitating the development of computer-assisted approaches (Weinstein 2015, Schneider et al. 2019).

Video classification is a challenging modeling problem, with the challenges of image classification amplified because the same sources of natural visual variation occur not only between videos but also within videos as objects move around and change poses, scales, illuminations, and backgrounds during the course of a single video. The video camera itself can move around during recording, introducing additional variation, particularly in environments where cameras move due to wind or water movement, or because cameras are attached to animals moving around their environment. The temporal component of video also presents significant modeling challenges not only because it dramatically increases the size of video data but also because the relevant visual features required to classify a video can span several frames with no single frame containing enough information on its own. The pixels of an image representing objects are not only correlated spatially to form visual object features in a single frame but also correlated through time.

Like image classification, traditional computer-based approaches to video classification have primarily used feature engineering algorithms that create input variables based on predetermined traits. The main limitations of these approaches arise from their need to know how to represent input features in advance—this requires substantial knowledge of the study species, and hinders generalization across species and environmental contexts (Schneider et al. 2019).

Deep neural networks (DNNs, LeCun et al. 2015, Goodfellow et al. 2016) are highly flexible machine learning models that use stacked non-linear combinations of inputs, trained using gradient descent with backpropagation, that learn feature representations relevant to provided labeled data, thus no longer requiring feature engineering. DNNs have been successfully used to tackle many challenging perceptual problems involving image, video, audio, or text, where hand-designing input feature representations are nontrivial (Liu et al. 2016).

A convolutional neural network (CNN, Goodfellow et al. 2016, ch. 9) is a specialized kind of DNN architecture that takes advantage of the characteristics of image data to learn hierarchies of local features that are invariant to common translation operations like shifting, stretching, and rotation. This reduces the number of required parameters while leaving enough representational power to achieve human-level performance on image classification and other tasks involving data that have a regular grid-like topology of locally correlated hierarchical features (Schneider et al. 2019). CNNs typically involve a stacked sequence of convolutional layers—traversing the network, and the output of each of these layers can be thought of as an increasingly complex summary or “encoding” of the input image as a one-dimensional numeric vector. Beyond their use in classifying entire images, they form the basis for related tasks such as detecting objects of interest within images (Ren et al. 2016, Tian et al. 2019) and object tracking, which extends single-frame object detection to track an object of interest across multiple frames (Danelljan et al. 2017, Li et al. 2018, Yang et al. 2019).

CNNs have found numerous, and increasing, applications in ecological studies (Weinstein 2018a, Christin et al. 2019), where image classification has been used for species identification (Zhang et al. 2016, Gomez Villa et al. 2017, Weinstein 2018b), count surveys (Borowicz et al. 2018, Torney et al. 2019, Gray et al. 2019b), individual animal re-identification (Schneider et al. 2019), and morphometric measurement (Gray et al. 2019a). Applications to video classification, however, remain rare. Trinh et al. (2016) combined neural network architectures to detect birds flying into wind turbines from sequences of input frames, and Beery et al. (2020) combined an object detection model with two attention-based modules that capture short- and long-term dependencies between frames, focussing on static-camera applications such as camera traps. Otherwise, most studies have either classified frames in isolation (Siddiqui et al. 2018) or used previous frames primarily to improve the discrimination of the focal animal from background scenery, using motion-detection algorithms (Zhang et al. 2016, Weinstein 2018b).

There are three approaches to using DNNs for video classification beyond treating the problem as an image classification task by modeling frames independently. The simplest approach concatenates the vector encodings obtained from each of a sequence of input images to predict the class of the last image in the sequence; images in the input sequence are considered to be independent. The second approach uses the sequence of vector encodings produced from the sequence of input images as input to a second model—a recurrent neural network (RNN), a specialized architecture often used to process sequential data involving a temporal component (Donahue et al. 2014, Trinh et al. 2016). Finally, CNNs can be directly modified to incorporate motion information in videos by extending their convolution from two spatial dimensions (width and height) to three spatio-temporal dimensions (width, height, and time), parameters of which are jointly estimated (Tran et al. 2015).

In this paper, we have used these approaches to perform frame-by-frame annotation of two video datasets. The first was taken from a fixed underwater camera placed inside nets at a salmon trap net fishery in Scotland, for the purpose of detecting seal visits to salmon nets and ultimately reducing conflict between fisheries and seals. Here, the task was to replicate manual annotations indicating whether a seal was present in a frame, based on that and preceding frames. The second dataset was collected by animal-borne cameras deployed on African penguins in South Africa. Here, the purpose was to replicate manual annotations allocating each frame to one of six pre-defined classes covering diving and surface behavior exhibited by the birds. The first of these applications can also be addressed by object tracking methods, but the second cannot, and to the best of our knowledge, this is the first time DNNs have been applied to annotate animal-borne video. For each dataset, our primary goal was to develop classifiers that could assist manual annotation by identifying temporal regions of interest in the video, and to evaluate whether incorporating the temporal component of video brings any improvement in classification accuracy, relative to an image-only benchmark.

MATERIALS AND METHODS

Data

Seals.—An underwater video system was used to study seal behavior at a salmon trap net fishery in northeast Scotland in 2015 as part of a program of research aimed at reducing conflict between fisheries and seals. Cameras were placed inside static coastal nets to monitor seals as they moved in and out of nets to depredate salmon. There was no artificial lighting and so the cameras recorded during hours of daylight.

The labeled component of the dataset consisted of six video recordings of approximately 140 min each, converted into images at 4fps. A total of 152 instances in which a seal entered the net were observed by manual inspection, and entry and exit times for each of these recorded (Appendix S1: Fig. S1). Visits lasted between 2s and 59s, with an average duration of 13.5s. Seals were not visible in frame for the entire duration of a visit, so all images between the start and end times of a recorded visit were manually inspected and labeled as containing a seal or not. After processing, there were 4419 images containing a seal. While the vast majority of footage does not contain a seal in frame, we restricted the number of absence images to 7809, roughly twice the number of seal images, to avoid a large class imbalance (Schneider et al. 2020). Absence images were collected by randomly sampling segments of video from the remainder of the video. Images from four videos were used to train models (3826 seal, 6949 no seal), while images from each of the remaining two videos were used as validation (407 seal, 973 no seal) and test (192 seal, 111 no seal) datasets, respectively.

Penguins.—Animal-borne video recorders (AVR) were deployed on breeding African penguins attending small chicks at Stony Point, South Africa, between 2015 and 2016 (McInnes et al. 2017). The AVRs were tube-shaped, and together with the casing weighed 100 g with dimensions 104 × 26 × 28 mm. Devices were attached to the lower backs of the penguins with strips of waterproof tape during the evening preceding an anticipated foraging trip. AVRs were programmed to divide the battery life into two recording bins of approximately 30 min each, at sunset and midday to reflect potential temporal

differences in diving behavior. Recorders were retrieved when the bird returned to the colony, either on the same day that the bird was at sea and after the bird had time to provision its chicks, between 16:00 and 20:00, or the following morning if the bird could not be located the previous day.

The labeled component of the dataset consisted of 12 video recordings of approximately 30 min each, again converted into images at 4fps. These were manually classified into five diving behaviors (subsurface diving [less than 1 m]; shallow diving [1–5 m]; and the descent, bottom, and ascent phases of deep dives) and one surface behavior (searching, see Appendix S1: Fig. S2). A total of 52,722 images were obtained, with substantial imbalance between behaviors (Appendix S1: Table S1). Images from nine videos were used to train models (41,958 images, see Appendix S1: Table S1 for distribution over behaviors), while images from the remaining videos were used as validation (two videos, 7168 images) and test (one video, 3596 images) datasets, respectively.

Neural networks

We consider four broad classes of models, of increasing complexity. The first ignores the temporal aspect of video data and attempts to classify each image independently using a standard CNN-based approach. Pre-trained CNNs (VGG16, ResNet50, Inception v3, and Inception-ResNet v2) were truncated at the final convolutional layer—the output of this intermediate layer summarizes or “encodes” an image in a one-dimensional vector. Up to three dense layers were added to the truncated network, and a new output layer added for the (seal or penguin) classification task. The second model used the same approach, but classified an image by first concatenating the vector encoding obtained from the truncated layer for that image with similar vectors obtained for the previous $F-1$ images. This concatenated vector, which summarizes a set of F consecutive images rather than (as in the first model) just a single image, was then passed these to subsequent dense layers as before. The third model was the spatial-then-temporal model described in the introduction (Donahue et al. 2014). To classify a single image, it took the vector encodings from the last F images (including

the current image), as in the previous model, but instead of concatenating the encodings it passed these as input to a recurrent neural network, which combined these temporally (Fig. 1). We used two pre-trained CNNs to encode frames (ResNet50, VGG16) and three different recurrent units (long short-term memory [LSTM], SimpleRNN, and gated recurrent units [GRU]). One key step was to pre-compute the frame vector encodings from the pre-trained CNN models so that these did not have to be re-computed in each RNN model. A single training epoch for the mixed long-term recurrent convolutional network (LRCN, Donahue et al. 2014) architecture with a VGG encoder took approximately 15 min without pre-computation but only 3 seconds with pre-computed features (because most of the computation time was spent in the CNN part of LRCN). The final model jointly modeled spatial and temporal aspects using a 3-dimensional CNN that convolves simultaneously over both spatial and temporal features (Tran et al. 2015). Because convolutions occur simultaneously over space and time, the 3-D CNN cannot leverage pre-computation, and generators had to be used to stream the data from disk to avoid out-of-memory problems. Despite various attempts at optimization, a single model took approximately 3 d to converge on a single GPU, and returned substantially worse accuracy than even an image-only model. We therefore do not report on these results further.

We chose model hyperparameters using a grid search over the number of nodes in each of the three dense layers in Models 1 and 2 (32, 64, 96, ..., 512), the dropout rate (0, 0.1, 0.2, ..., 0.5), and the length of the sequence of images used in Models 2 and 3 (1, 3, 5, 7, 9, ..., 31). Following Krizhevsky et al. (2012), each model's weights for dense and recurrent layers were initialized using the Xavier initialization and each model was trained in three rounds of 20 epochs with an early stopping patience of five epochs using the Adam optimizer. The learning rate was initially set to 0.001 and reduced by a factor of 10 between training rounds, and max pooling was used. Models were evaluated based on test set accuracy (proportion of all predictions that were correct), precision (proportion of positive predictions that were correct), and recall (proportion of positive examples correctly predicted). For the

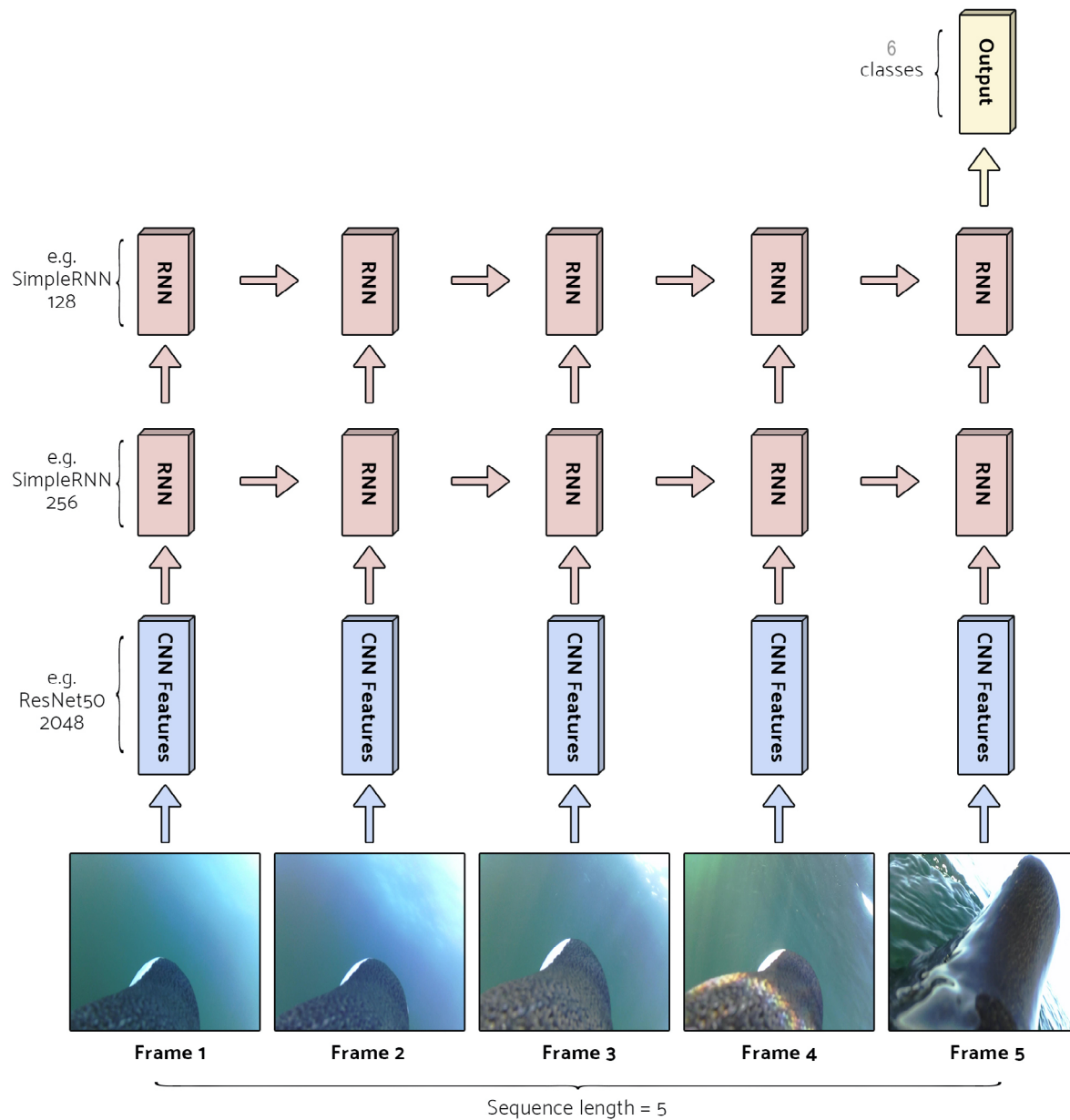


Fig. 1. A “spatial-then-temporal” neural network for frame-by-frame video classification. To predict the class of a frame (Frame 5), a pre-trained, truncated CNN (e.g., ResNet50) is used to summarize or “encode” each of a sequence of images (here, the last five frames) as one-dimensional numeric vectors. The sequence of vector encodings is then used as input in a recurrent neural network (RNN), here shown using two SimpleRNN layers. The RNN outputs predicted probabilities that the behavior in the final frame is of type $i = 1, \dots, 6$.

seals dataset, seal presence is a natural choice for the positive class. Optimal thresholds for converting predicted probabilities into binary classifications were those that maximized the F1 score, a function of precision and recall, in the

validation dataset (0.47, 0.44, 0.46, 0.79 for the models in Table 1, respectively). For the multi-class penguins dataset, images were allocated into the class with the highest predicted probability, precision and recall were obtained for each

Table 1. Classification accuracy for three best video models and best image model. Including temporal information in the form of an LRCN led to very marginal improvement in the easier seal detection task, but gave a 25% relative improvement in the ability to discriminate penguin behaviors, largely due to improved performance at the start and end of behaviors (Fig. 3). Further details on the architectures and run times of these models are given in Appendix S1: Tables S3 and S4.

Architecture	LRCN	LRCN	LRCN	IMAGE
Seal detection model				
Accuracy (% Test)	89.4	89.2	89.1	88.1
Precision (% Test)	100	98.8	99.4	99.4
Recall (% Test)	83.9	84.4	83.9	81.8
Accuracy (% Validation)	96.5	96.0	95.9	94.5
Accuracy (% Train)	95.5	95.4	95.5	95.2
Penguin behavior classifier				
Accuracy (% Test)	85.4	84.0	84.2	80.5
Precision (% Test)	85.4	84.0	84.2	80.5
Recall (% Test)	87.6	87.6	85.5	82.8
Accuracy (% Validation)	82.6	82.4	81.0	81.5
Accuracy (% Train)	90.0	88.9	94.4	88.7

class, and overall precision and recall were calculated as an average of these, weighted by sample size. Models were implemented using the TensorFlow library with Keras. Training and testing were done on a three separate Linux virtual machine instances running on Google Cloud Platform, each with eight Nvidia Tesla K80 graphics processing units (GPUs), 160 GB of RAM and 32 CPU cores. Code and analysis scripts are available online at <https://github.com/alxcnwy/Deep-Neural-Networks-for-Video-Classification>. A subset of seal and penguin video recordings, manual annotations, and results has been stored on Zenodo: <https://doi.org/10.5281/zenodo.3842040>.

RESULTS

A video component did not bring meaningful benefits in detecting seals, with both image-only and video models accurately classifying 88% and 89% of images in the test set, although both precision and recall were marginally higher in video models (Table 1). Most incorrect classifications occurred at the beginning and end of visits, as the seal was entering or exiting the field of view and where only a small part of the seal may be in

view (Fig. 2). All 152 seal visits across training, validation, and test sets were detected by either model.

Including temporal information in video data, in the form of spatial-then-temporal models, improved the accuracy of penguin behavior classifications from 80.5% (image-only benchmark) to 85.4%, a 25% relative reduction in classification error (Table 1), and improved both precision and recall. Models concatenating frame encodings occupied an intermediate position between full video and image-only models. Classification accuracy improved for most penguin behavior types (Appendix S1: Table S2), but particularly for descent and bottom dive phases (precision increasing by 17% and 14%), and for shallow and subsurface dives (recall increasing by 12% and 13%). Image-only models tended to misclassify bottom dives as descent dives, and mistook parts of the ascending and descending dive phases for shallow dives. To some extent, this reflects fuzzy boundaries between behavioral classes, but temporal information resolved some of these misclassifications (Fig. 3). Search activity, the sole surface behavior and also the most prevalent class, was almost perfectly discriminated.

Preferred LRCN models for seal detection achieved a degree of parsimony by using a relatively short sequence of frames, and in exchange used relatively complex pre-trained CNN (ResNet50) and RNN (LSTM) architectures (Appendix S1: Table S3). In contrast, equivalent preferred models for penguin behavior classification used longer sequences of frames, but simpler CNN (VGG16) and, sometimes, RNN (SimpleRNN) architectures (Appendix S1: Table S4). Both applications selected a relatively large number of nodes in the final hidden layers.

DISCUSSION

Although images are more commonly used in ecological research and are easier to work with (Swinnen et al. 2014), movement information contained in video provides richer insight into animal behavior and taking this into account can improve the identification of animals and their behaviors (Trinh et al. 2016). We found that for a relatively simple task—detecting seal activity in an image—an image-only CNN was adequate, and incorporating temporal information

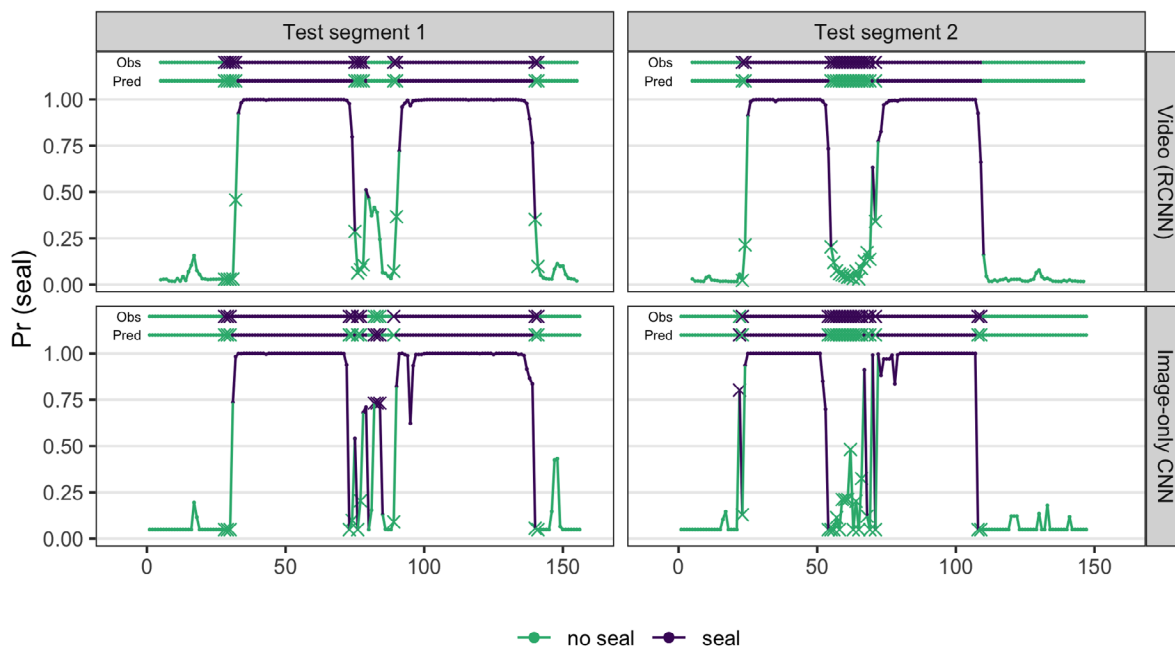


Fig. 2. Predicted probabilities of seal activity in salmon nets, with misclassifications plotted as crosses. Observed and predicted classes are plotted above the probabilities, using the same notation. Apart from one false negative (segment 2, frames 22–108), all incorrect classifications are at the beginning and end of visits, where only a small part of the seal may be in view. All visits are clearly identified.

did not meaningfully improve out-of-sample performance, even for those difficult cases in which a seal enters or exits the field of view. For a more difficult task of inferring penguin behavior from animal-borne cameras, using a video model led to substantial reduction in classification error over an image-only model and was particularly useful in disentangling certain kinds of diving behavior. In both applications, accuracy is not sufficient for full automation of the tasks, but can facilitate manual processes by partially labeling the data—identifying those classes that can be accurately discriminated and pointing the researcher to segments requiring closer inspection. Our datasets were relatively small, consisting of 6–12 h of labeled footage, and the ability of the models to generalize to new environments is unclear, but even in those classes where absolute performance was moderate, video models outperformed image-only models. Improvements are likely to be larger with larger datasets.

Practically, researchers wanting to construct a model for the frame-by-frame annotation of

video have to follow a number of steps: manually labeling a subset of the data; converting the video into images; allocating these images between training, validation, and test sets; choosing appropriate neural network architectures and estimating the parameters of those models; selecting a preferred model and using it to process the unlabeled portion of the data; and linking frame-by-frame predictions to the broader research objectives for which the classifier was developed.

Video data are manually annotated by recording the start and end times of events whose boundaries may be difficult to distinguish precisely. Poorly separated classes can reduce classification accuracy, and preprocessing steps for image classification sometimes remove ambiguous images to improve class separability. Video models, however, use a sequence of frames t , $t-1$, ..., $t-F$ to predict the class of frame t , and removing ambiguous images makes the time difference between adjacent images variable. While it is possible that removing ambiguous examples may improve accuracy more than maintaining

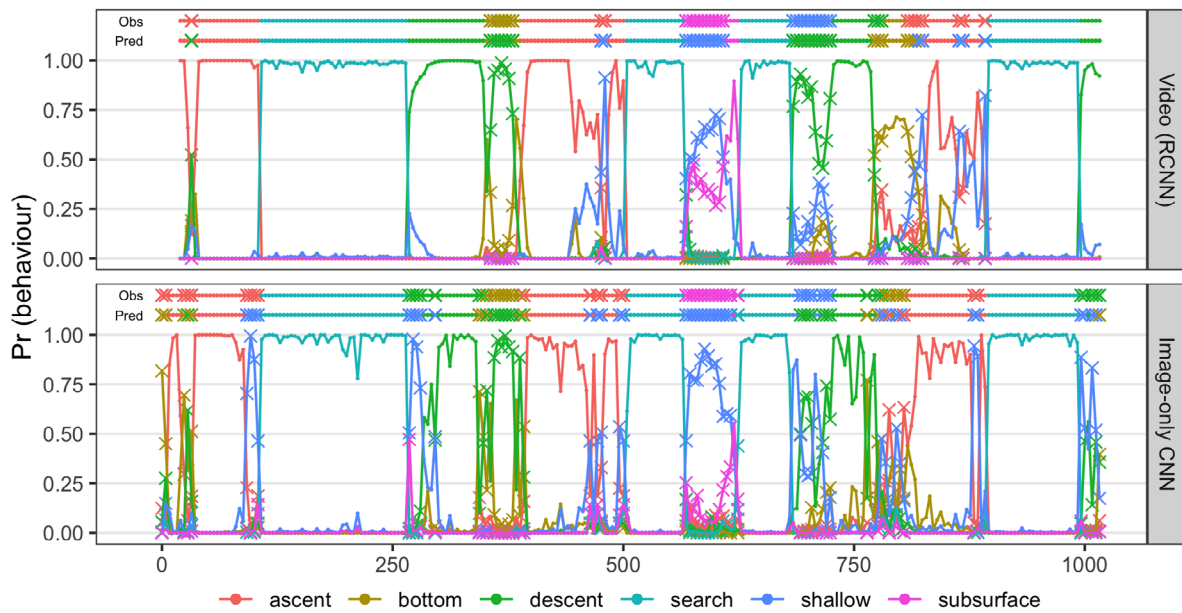


Fig. 3. Predicted probabilities for penguin behavior classes, with misclassifications plotted as crosses. Observed and predicted classes are plotted above the probabilities, using the same notation. Image-only models tend to misclassify bottom dives as descent dives (frames 350–390), and ascending and descending dive phases as shallow dives (frames 90–110 and 260–280). Video models resolve some of these errors. They also smooth transitions between behaviors (frames 260–280), better identify periods where classification uncertainty is high (frames 570–620, 750–850) and where alternate interpretations are possible (frames 570–620).

constant time difference between images, this is likely to be case-specific, and not generally recommended. Rather, the presence of ambiguous images places an effective upper limit on the accuracy that can be achieved, which may or may not impact on broader research objectives. For seal visits, for example, the detection of a seal presence is more important than identifying the exact time of entry. The first and last few frames of a visit often contain only a tiny sliver of seal or, because the times are approximate, no seal at all. These frames reduce classification accuracy but have very little bearing on the practical usefulness of the classifier.

Video data are converted to images at a user-specified frame rate, with the recording equipment setting an upper bound. A higher frame rate increases the number of images available to train models, which is always beneficial as long as there are meaningful differences between adjacent images. It is important to randomly allocate contiguous sequences of frames, that is, video sequences, to training, validation, and test

datasets, rather than randomly allocating the frames themselves. Doing the latter breaks apart sequences, losing potentially valuable information, and also means that very similar images occur in both training and test sets. We also recommend assessing whether the video in the test dataset has the same environmental conditions as video used to train the model (e.g., if a random segment of each file is used to test). If so, the ability of the model to generalize to new environments may be overestimated.

When building an LRCN, key choices are what frame rate and sequence length to use. These factors are study-specific, and the chosen frame rate need not be the same as the frame rate used to convert video to frames (Yue-Hei Ng et al. 2015). Higher frame rates allow for fine-scale changes in movement to be captured, but the same number of frames covers a shorter time interval. Increasing sequence length requires more parameters, increasing the chances of overfitting and requiring more data. Which of the two—looking back further in time or capturing fine-scale

movement—benefits classification accuracy more will be study-specific. These factors can be investigated by searching over possible frame rate/length pairs, but this quickly becomes computationally expensive. Our applications have relatively little labeled data and so we fixed the frame rate to one that would allow broad differences in behavior, observed over a few seconds, with $5 < F < 10$. Pre-trained CNNs offer a parsimonious way of summarizing images in a form that can be passed on the second-stage RNN (Donahue et al. 2014). Our best seal model combined a relatively complex CNN and RNN with a short frame sequence, whereas the best penguin model had a simple CNN and RNN, but used a longer sequence of frames. Since model complexity is primarily achieved through more parameters, this balance reflects the familiar goal of reducing validation error through model parsimony. Frame rate and sequence length can be chosen via standard hyperparameter selection practices, for example, grid search, or by first selecting a frame rate that is relatively low but still able to capture the desired transitions in behavior or class membership, and then to select the sequence length that optimizes performance on a validation dataset.

Our models allow new video footage to be classified on a frame-by-frame basis, with some expected degree of accuracy. Linking this back into research objectives is the final step in the process. The seal classifier is intended to be used as a detection system. Even with a frame-specific false-negative rate of 10%, no visits were missed entirely. An alarm system, triggered by N predicted presences in a sequence of M frames, is easily established, with N and M determined by balancing costs of false positives and negatives. Graphical displays such as Fig. 3 convey this information in an easily digested way. Higher error rates prevent the use of the penguin behavior classifier for the purpose it was intended for—replicating a human observer and calculating energy budgets—because certain classes of behavior are poorly identified. However, surface and diving behaviors were discriminated with almost no error, and deep and shallow/subsurface dives were also well differentiated. These distinctions hold practical value and also limit the amount of manual labeling that must be done.

Deep learning holds enormous promise for automating the labeling of video data, a process that looks increasingly unsustainable with manual methods. Case studies such as the ones reported here play an important role in reporting successes and failures, and developing and disseminating best practices. Classification of ecological data is difficult. Limited time and other resources, remote locations, and rare or difficult-to-detect target species, serve to decrease sample sizes at the same time that variable background environments increase the necessary sample sizes for good classification. In these contexts, full automation is perhaps, for the time being, unrealistic. Facilitating the process of manually annotating video datasets is both valuable and achievable. Video data have the great advantage that large datasets, in terms of numbers of images, are often collected relatively quickly. This offers exciting opportunities for developing and testing deep learning approaches. Our study suggests that many applications may benefit from incorporating temporal information in video, where the goal remains to predict the class to which a particular frame or image belongs. We expect these models to be widely used and developed in the near future.

ACKNOWLEDGMENTS

The penguin data collection was financially supported by Homebrew Films and the Percy Fitzpatrick Institute of African Ornithology. All fieldwork was done under permission from the South African Department of Environmental Affairs (permit nos RES 2015/38 and RES 2016/100) and Cape Nature (permit no. AAA007-00209-0056). We thank Cape Nature and the South African Department of Environmental Affairs for providing permission to carry out the study. The seal data collection was funded by the Scottish Government through the Marine Mammal Scientific Support Research Programme. ID is supported in part by funding from the National Research Foundation of South Africa (Grant ID 90782, 105782). All authors conceived the work together. RH collected and annotated seal data, and provided feedback on model usability results. AM did the same for the penguin data. AC and ID developed the modeling approach. AC implemented the models and performed analyses. AC and ID wrote the paper. All authors contributed critically to the drafts and gave final approval for publication.

LITERATURE CITED

- Anderson, K., and K. J. Gaston. 2013. Lightweight unmanned aerial vehicles will revolutionize spatial ecology. *Frontiers in Ecology and the Environment* 11:138–146.
- Beery, S., G. Wu, V. Rathod, R. Votel, and J. Huang. 2020. Context r-cnn: long term temporal context for per-camera object detection. Pages 13075–13085 in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. IEEE, Piscataway, NJ, USA.
- Borowicz, A., et al. 2018. Multi-modal survey of adélie penguin mega-colonies reveals the danger islands as a seabird hotspot. *Scientific Reports* 8:1–9.
- Christin, S., E. Hervet, and N. Lecomte. 2019. Applications for deep learning in ecology. *Methods in Ecology and Evolution* 10:1632–1644.
- Cruzan, M. B., B. G. Weinstein, M. R. Grasty, B. F. Kohn, E. C. Hendrickson, T. M. Arredondo, and P. G. Thompson. 2016. Small unmanned aerial vehicles (micro-uavs, drones) in plant ecology. *Applications in Plant Sciences* 4:1600041.
- Danelljan, M., G. Bhat, F. Shahbaz Khan, and M. Felsberg. 2017. Eco: efficient convolution operators for tracking. Pages 6638–6646 in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, Piscataway, NJ, USA.
- Donahue, J., L. A. Hendricks, M. Rohrbach, S. Venugopalan, S. Guadarrama, K. Saenko, and T. Darrell. 2014. Long-term Recurrent Convolutional Networks for Visual Recognition and Description. Pages 2625–2634 in *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. IEEE, Piscataway, NJ, USA.
- Gomez Villa, A., A. Salazar, and F. Vargas. 2017. Towards automatic wild animal monitoring: identification of animal species in camera-trap images using very deep convolutional neural networks. *Ecological Informatics* 41:24–32.
- Goodfellow, I., Y. Bengio, A. Courville, and Y. Bengio. 2016. *Deep learning*. MIT press Cambridge.
- Gray, P. C., K. C. Bierlich, S. A. Mantell, A. S. Friedlaender, J. A. Goldbogen, and D. W. Johnston. 2019a. Drones and convolutional neural networks facilitate automated and accurate cetacean species identification and photogrammetry. *Methods in Ecology and Evolution* 10:1490–1500.
- Gray, P. C., A. B. Fleishman, D. J. Klein, M. W. McKown, V. S. Bézy, K. J. Lohmann, and D. W. Johnston. 2019b. A convolutional neural network for detecting sea turtles in drone imagery. *Methods in Ecology and Evolution* 10:345–355.
- Krizhevsky, A., I. Sutskever, and G. E. Hinton. 2012. ImageNet Classification with Deep Convolutional Neural Networks. *Proceedings of the 25th International Conference on Neural Information Processing Systems*. Curran Associates Inc., Red Hook, NY, USA.
- LeCun, Y., Y. Bengio, and G. Hinton. 2015. Deep learning. *Nature* 521:436–444.
- Li, B., J. Yan, W. Wu, Z. Zhu, and X. Hu. 2018. High performance visual tracking with siamese region proposal network. Pages 8971–8980 in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, Piscataway, NJ, USA.
- Liu, W., Z. Wang, X. Liu, N. Zeng, Y. Liu, and F. E. Alsaadi. 2016. A survey of deep neural network architectures and their applications. *Neurocomputing* 234:11–26.
- McInnes, A. M., C. McGeorge, S. Ginsberg, L. Pichegru, and P. A. Pistorius. 2017. Group foraging increases foraging efficiency in a piscivorous diver, the African penguin. *Royal Society Open Science* 4:170918.
- Ren, S., K. He, R. Girshick, and J. Sun. 2016. Faster R-CNN: towards real-time object detection with region proposal networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 39:1137–1149.
- Rutz, C., and J. Troschianko. 2013. Programmable, miniature video-loggers for deployment on wild birds and other wildlife. *Methods in Ecology and Evolution* 4:114–122.
- Schneider, S., S. Greenberg, G. W. Taylor, and S. C. Kremer. 2020. Three critical factors affecting automated image species recognition performance for camera traps. *Ecology and Evolution* 10:3503–3517.
- Schneider, S., G. W. Taylor, S. Linquist, and S. C. Kremer. 2019. Past, present and future approaches using computer vision for animal re-identification from camera trap data. *Methods in Ecology and Evolution* 10:461–470.
- Siddiqui, S. A., A. Salman, M. I. Malik, F. Shafait, A. Mian, M. R. Shortis, and E. S. Harvey. 2018. Automatic fish species classification in underwater videos : exploiting pre-trained deep neural network models to compensate for limited labelled data. *ICES Journal of Marine Science* 75:374–389.
- Swinnen, K. R. R., J. Reijnen, M. Breno, and H. Leirs. 2014. A novel method to reduce time investment when processing videos from camera trap studies. *PLOS ONE* 9:e98881.
- Takahashi, A., K. Sato, Y. Naito, M. Dunn, P. Trathan, and J. Croxall. 2004. Penguin-mounted cameras glimpse underwater group behaviour. *Proceedings of the Royal Society of London. Series B: Biological Sciences* 271:S281–S282.
- Tian, Z., C. Shen, H. Chen, and T. He. 2019. FCOS: fully convolutional one-stage object detection.

- Pages 9627–9636 in Proceedings of the IEEE International Conference on Computer Vision. IEEE, Piscataway, NJ, USA.
- Torney, C. J., D. J. Lloyd-Jones, M. Chevallier, D. C. Moyer, H. T. Maliti, M. Mwita, E. M. Kohi, and G. C. Hopcraft. 2019. A comparison of deep learning and citizen science techniques for counting wildlife in aerial survey images. *Methods in Ecology and Evolution* 10:779–787.
- Tran, D., L. Bourdev, R. Fergus, L. Torresani, and M. Paluri. 2015. Learning Spatiotemporal Features with 3D Convolutional Networks. Pages 4489–4497 in Proceedings of the IEEE International Conference on Computer Vision. IEEE, Piscataway, NJ, USA.
- Trinh, T. T., R. Yoshihashi, R. Kawakami, M. Iida, and T. Naemura. 2016. Bird detection near wind turbines from high-resolution video using LSTM networks. Pages 1–4 in Proceedings of the 15th World Wind Energy Conference. Volume 5. WWEA, Tokyo.
- Weinstein, B. G. 2015. Motion meerkat: integrating motion video detection and ecological monitoring. *Methods in Ecology and Evolution* 6:357–362.
- Weinstein, B. G. 2018a. A computer vision for animal ecology. *Journal of Animal Ecology* 87:533–545.
- Weinstein, B. G. 2018b. Scene-specific convolutional neural networks for video-based biodiversity detection. *Methods in Ecology and Evolution* 9:1435–1441.
- Yang, Z., Q. Wang, L. Bertinetto, W. Hu, S. Bai, and P. H. Torr. 2019. Anchor diffusion for unsupervised video object segmentation. Pages 931–940 in Proceedings of the IEEE International Conference on Computer Vision. IEEE, Piscataway, NJ, USA.
- Yue-Hei Ng, J., M. Hausknecht, S. Vijayanarasimhan, O. Vinyals, R. Monga, and G. Toderici. 2015. Beyond short snippets: deep networks for video classification. Pages 4694–4702 in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. IEEE, Piscataway, NJ, USA.
- Zhang, Z., Z. He, G. Cao, and W. Cao. 2016. Animal Detection from Highly Cluttered Natural Scenes Using Spatiotemporal Object Region Proposals and Patch Verification. *IEEE Transactions on Multimedia* 18:2079–2092.

DATA AVAILABILITY STATEMENT

Code and analysis scripts are available online at <https://github.com/alxcnwy/Deep-Neural-Networks-for-Video-Classification>. A subset of seal and penguin video recordings, manual annotations, and results has been stored on Zenodo: <https://doi.org/10.5281/zenodo.3842040>.

SUPPORTING INFORMATION

Additional Supporting Information may be found online at: <http://onlinelibrary.wiley.com/doi/10.1002/ecs2.3384/full>