

Visual Wetland Birds Dataset: Bird Species Identification and Behavior Recognition in Videos

Javier Rodríguez-Juan¹, David Ortiz-Perez¹, Manuel Benavent-Lledo¹, David Mulero-Perez¹, Pablo Ruiz-Ponce¹, Adrian Orihuela-Torres², Jose Garcia-Rodriguez^{1,*}, and Esther Sebastián-González^{2,3}

¹Department of Computer Technology, University of Alicante, Alicante, 03690, Spain

²Department of Ecology, University of Alicante, Alicante, 03690, Spain

³'Ramón Margalef' Multidisciplinary Institute for the study of the Environment. University of Alicante, Alicante, 02690, Spain

*corresponding author(s): Jose Garcia-Rodriguez (jgarcia@dtic.ua.es)

ABSTRACT

The current biodiversity loss crisis makes animal monitoring a relevant field of study. In light of this, data collected through monitoring can provide essential insights, and information for decision-making aimed at preserving global biodiversity. Despite the importance of such data, there is a notable scarcity of datasets featuring videos of birds, and none of the existing datasets offer detailed annotations of bird behaviors in video format. In response to this gap, our study introduces the first fine-grained video dataset specifically designed for bird behavior detection and species classification. This dataset addresses the need for comprehensive bird video datasets and provides detailed data on bird actions, facilitating the development of deep learning models to recognize these, similar to the advancements made in human action recognition. The proposed dataset comprises 178 videos recorded in Spanish wetlands, capturing 13 different bird species performing 7 distinct behavior classes. In addition, we also present baseline results using state of the art models on two tasks: bird behavior recognition and species classification.

Background & Summary

Under the current scenario of global biodiversity loss, there is an urgent need for more precise and informed environmental management¹. In this sense, data derived from animal monitoring plays a crucial role in informing environmental managers for species conservation^{2,3}. Animal surveys provide important data on population sizes, distribution and trends over time, which are essential to assess the state of ecosystems and identify species at risk^{4,5}. By using systematic monitoring data on animal and bird populations, scientists can detect early warning signs of environmental changes, such as habitat loss, climate change impacts, and pollution effects⁶⁻⁸. This information helps environmental managers develop targeted conservation strategies, prioritize resource allocation, and implement timely interventions to protect vulnerable species and their habitats^{2,9,10}. Furthermore, bird surveys often serve as indicators of local ecological conditions, given birds' sensitivity to environmental changes, making them invaluable in the broader context of biodiversity conservation and ecosystem management¹¹. However, monitoring birds, as any other animal, is highly resource-consuming. Thus, automated monitoring systems that are able to reduce the investment required for accurate population data are much needed.

The first step to create algorithms that detect species automatically is to create datasets with information on the species traits to train those algorithms. For example, a common way to classify species is by their vocalizations¹². For this reason, organizations such as the Xeno-Canto Foundation¹ compiled a large-scale online database¹³ of bird sounds from more than 200,000 voice recordings and 10,000 species worldwide. This dataset was crowdsourced and today it is still growing. The huge amount of data provided by this dataset has facilitated the organization of challenges to create bird-detection algorithms using acoustic data in understudied areas, such as those led by Cornell Lab². This is the case of BirdCLEF2023¹⁴, or BirdCLEF2024¹⁵, which used acoustic recordings of eastern African and Indian birds, respectively. While these datasets contain many short recordings from a wide variety of different birds, other authors have released datasets composed of fewer but longer recordings, which imitate a real wildlife scenario. Examples of this are NIPS4BPlus¹⁶, which contains 687 recordings summing a total of 30 hours of recordings or BirdVox-full-night¹⁷, which has 6 recordings of 10 hours each.

Although audio is a common way to classify bird species and the field of bioacoustics has increased tremendously in the latest years, another possible approach to identify species automatically is using images¹⁸. One of such bird image datasets is

¹<https://xeno-canto.org/>

²<https://www.birds.cornell.edu/home/>

Birds525,³ which offers a collection of almost 90,000 images involving 525 different bird species. Another standard image dataset is CUB-200-2011¹⁹, which provides 11,788 images from 200 different bird species. This dataset not only provides bird species, but also bounding boxes and part locations for each image. There are also datasets aimed at specific world regions like NABirds²⁰, which includes almost 50,000 images from the 400 most common birds seen in North America. This dataset provides a fine-grained classification of species as its annotations differentiate between male, female and juvenile birds. These datasets can be used to create algorithms for the automatic detection of the species based on image data.

However, another important source of animal ecology information that has been much less studied because of the technological challenges of its use are videos. Video recordings may offer information not only about which species are present in a specific place, but also about their behavior. Information about animal behavior may be very relevant to inform about individual and population responses to anthropogenic impacts and has therefore been linked to conservation biology and restoration success^{21–24}. Besides its potential for animal monitoring and conservation, the number of databases on wildlife behavior are more limited. For example, the VB100 dataset²⁵, comprises 1416 clips of approximately 30 seconds. This dataset involves 100 different species from North American birds. The unique dataset comprised by annotated videos with birds behavior available in the literature is the Animal Kingdom dataset²⁶, which is not specifically aimed at birds and contains annotated videos from multiple animals. Specifically, it contains 30,000 video sequences of multi-label behaviors involving 6 different animal classes. Table 1 summarizes the main information of the datasets reviewed.

Name	Modality	Region	Samples	Species	Only birds
Ipt Xeno-canto ¹³	Audio	All	+200,000 ⁴	12115 ⁴	✓
BirdCLEF2023 ¹⁴		Eastern Africa	16,900	264	✓
BirdCLEF2024 ¹⁵		India	24,460	942	✓
NIPS4BPlus ¹⁶		Spain/France	687	61	✓
BirdVox-full-night ¹⁷		USA	6	25	✓
Birds525 ³	Image	All	89,885	525	✓
CUB-200-2011 ¹⁹		All	11,788	200	✓
NABirds ²⁰		North America	48,000	400	✓
VB100 ²⁵	Video	North America	1,416	100	✓
Animal Kingdom ²⁶		All	30,000	-	
WetlandBirds (<i>Proposed</i>)		Spain	178	13	✓

Table 1. Summary of reviewed bird datasets.

Due to the scarcity of datasets involving birds videos annotated with its behaviors, this study proposes the development of the first fine-grained behavior detection dataset for birds. Differently from Animal Kingdom, where a video is associated with the multiple behaviors happening, in our dataset, spatio-temporal behavior annotations are provided. This implies that videos are annotated per-frame, where the behavior happening and the location is annotated in each frame (*i.e.* bounding box). Moreover, the identification of the bird species appearing in the video is also provided. The proposed dataset is composed by 178 videos recorded in Spanish wetlands, more specifically in the region of Alicante (southeastern Spain). The 178 videos expand to 2765 behavior clips involving 13 different bird species. The average duration of each of the behavior clips is 19.84 seconds and the total duration of the dataset recorded is 58 minutes and 53 seconds. The annotation process involved several steps of data curation, with a technical team working alongside a group of professional ecologists.

Table 2 reflects the different species collected for the dataset, distinguishing between their common and scientific names. The number of videos and minutes recorded for each species is also included.

Seven main behaviors were identified as key activities recorded in our dataset. These represent the main activities performed by waterbirds in nature²⁷. In Figure 1, these behaviors are specified alongside the number of clips recorded per each of them and the mean duration of each behavior in frames. A clip is a piece of video where a bird is performing a specific behavior. Animals often change among actions very fast, as a response to the changing environment. Thus, to consider a collection of movements of a bird as a behavior, this had to last a minimum of 30 frames, otherwise this collection of movements was identified as a sub-movement of another main behavior, which is the one annotated for those frames.

This dataset contains not only videos with a single individual, but also videos where several bird individuals appear together. This is the case for gregarious species, which are species that concentrate in an area for the purpose of different activities. Although the individuals of gregarious species often share the same behavior at the same time, it is also common that several behaviors can be seen in the same video at the same time. Figure 2 shows some samples of videos where this happens. Videos

³<https://www.kaggle.com/datasets/gpiosenka/100-bird-species>

⁴As it is a crowdsourced project, it grows with the time

Common name	Scientific name	Videos	Recorded minutes
Yellow-legged Gull	<i>Larus michahellis</i>	13	5.08
White wagtail	<i>Motacilla alba</i>	13	4.33
Squacco Heron	<i>Ardeola ralloides</i>	15	4.94
Northern shoveler	<i>Spatula clypeata</i>	14	3.49
Mallard	<i>Anas platyrhynchos</i>	10	2.94
Little-ringed plover	<i>Charadrius dubius</i>	10	1.93
Glossy ibis	<i>Plegadis falcinellus</i>	8	3.96
Gadwall	<i>Mareca strepera</i>	13	2.59
Eurasian moorhen	<i>Gallinula chloropus</i>	18	9.18
Eurasian magpie	<i>Pica pica</i>	16	5.95
Eurasian coot	<i>Fulica atra</i>	19	4.11
Black-winged stilt	<i>Himantopus himantopus</i>	14	3.55
Black-headed gull	<i>Chroicocephalus ridibundus</i>	15	6.84

Table 2. Data for each of the annotated species.

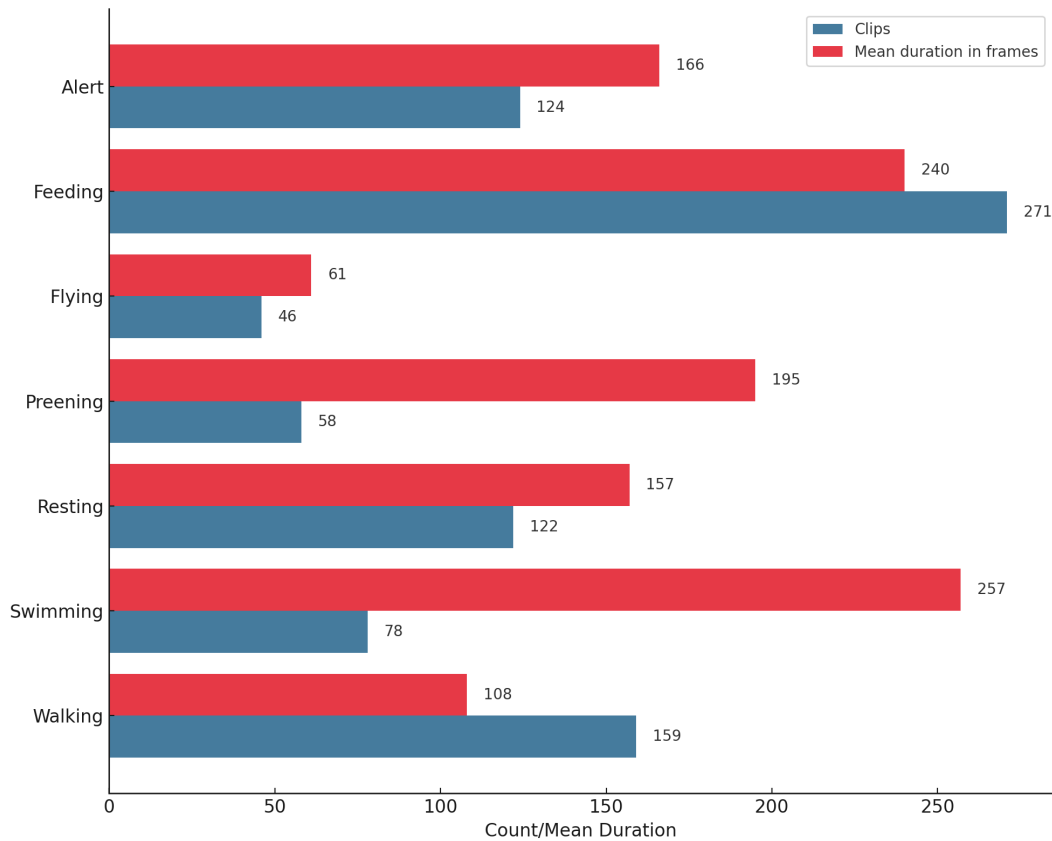


Figure 1. Comparison of total number of clips per behavior and mean duration of behavior clips.

involving different birds and/or performing different activities sequentially were cut in clips where a unique individual is performing a unique behavior in order to get the statistics shown in Figure 1.

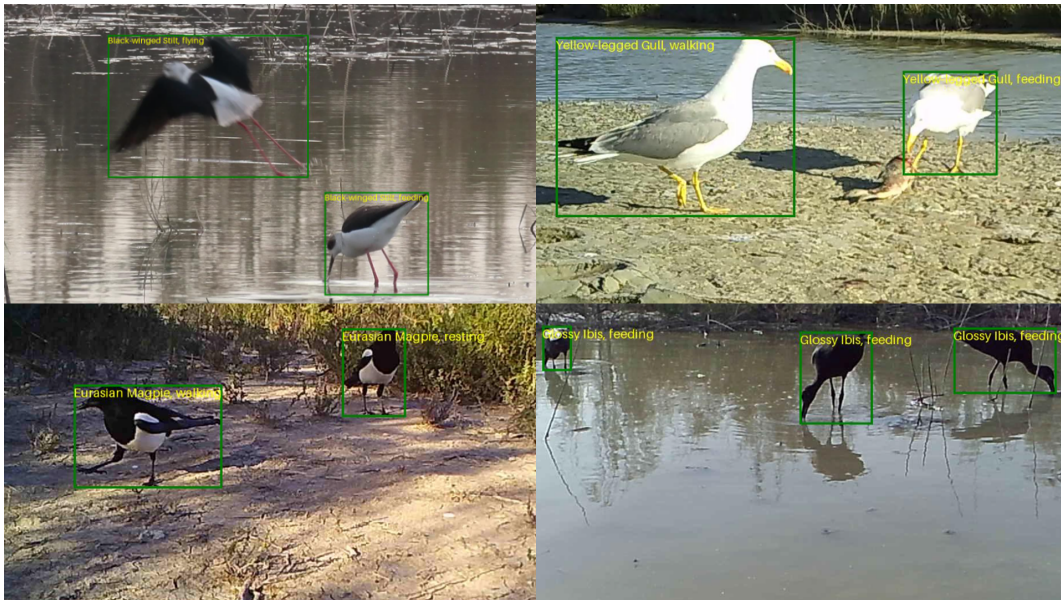


Figure 2. Frame samples where gregarious birds appear performing different behaviors.

Between the seven behaviors proposed, it should be underlined the difference between the *Alert*, *Preening* and *Resting* behaviors. These action distinctions were established by the ecology team. We considered that the bird was *Resting* when it was standing without making any movement. The bird was performing the *Alert* behavior when it was moving its head from one side to another, moving and looking around for possible dangers. Finally, we considered that the bird was *Preening* when it was standing and cleaning its body feathers with its beak. The remaining behaviors are not explained because of their obvious meaning.

As it can be seen in Figure 1, the number of clips per behavior is unbalanced between classes. This is because the recording of videos where some specific behaviors are happening is more uncommon, as happens with *Flying* or *Preening*, which represent the activities with the lowest number of clips in the dataset. These behaviors are difficult to record since they are performed with a lower frequency. In order to be able to collect more data on these less common behaviors, more hardware and human resources (*i.e.* cameras and professional ecologists) are needed to cover a wider area of the wetlands. Although the unbalanced nature of the behaviors, no balancing technique over this data was applied in the released dataset in order to maximize the number of different environments captured, ensuring in this way the variability of contexts where the birds are recorded.

Additionally, Figure 1 also shows the mean duration of the clips per behavior. It is worth noting the difference in the number of frames between *Flying*, that represents the minimum with 61 frames with respect to *Swimming*, which represents the absolute maximum with a value of 257 frames. This difference is explained in the nature of the behaviors, as swimming is naturally a slow behavior, which can be performed for a long time over the same area. However, flying is a fast behavior, and the bird quickly get outs of the camera focus, especially for videos obtained by camera traps, which cannot follow the bird while it is moving.

It is also common that birds perform two activities simultaneously. In that case, the most relevant behavior for the bird was the one annotated. In the proposed dataset, *Feeding* behavior is the one which is commonly done simultaneously to others such as *Walking* or *Swimming*. As *Feeding* behavior was considered by the ecologist experts as more relevant for birds, this behavior was always annotated in these cases.

In order to collect the videos, we deployed a set of camera traps and high quality cameras in Alicante wetlands. The camera traps were able to automatically record videos based on the motion detected in the environment. We complemented the camera trap videos with recordings from high quality cameras. In these videos, a human is controlling the focus of the camera, obtaining better views and perspectives of the birds being recorded. Species recorded, behaviors identified and the camera deployment areas were described by professional ecologist based on their expertise. In Figure 3 some video frame crops can be observed, where all the bird species developing the different behaviors available in the dataset can be seen.

After the data collection, a semi-automatic annotation method composed by an annotation tool and a deep learning model



Figure 3. Video frame crops of bird species performing the 7 behaviors composing the dataset.

was used in order to get the videos annotated. After the annotation, a cross-validation was conducted to ensure the annotation quality. This method is deeply explained in the next section.

In order to test the dataset for species and behavior identification, two baseline experimentation were carried out: one for the for the behavior detection task, which involves the correct classification of the behavior being performed by one bird during a set of frames, and a second one for the bird classification task, which involves the classification of the specie and the correct localization of the bird given input frames.

Methods

Data acquisition

The acquisition of the data was conducted within Alicante wetlands, specifically within the wetlands of *La Mata Natural Park* and *El Hondo Natural Park* (suteastern Spain). In these places, we deployed a collection of high-resolution cameras and camera traps in different areas of the wetlands. These areas were determined by the species expected to be recorded, as different species can be commonly seen in different wetland areas.

Camera traps are activated when movement is detected and thus can record for long periods of time without human intervention. The usage of automatic camera traps^{28–30} is common in the monitoring of wildlife as it provides a low-cost approach to collect video and image data from the environment. However, the focus of this camera is fix and thus the videos of the same individual are often short. Manual cameras require the presence of a human while recording and are thus more time-consuming. Also, the presence of the cameraman may affect the animal behavior. However, it permits manual changes of cameras' perspectives in order to correctly record the bird behavior. As different cameras were used videos of different resolutions were obtained: 87 videos at 1920x1080px, 75 videos at 1296x720px, 14 videos at 1280x720px, 1 video at 960x540px and 1 video at 3840x2160px.

The species selected were the most common found in the wetlands of Alicante, facilitating the recording of videos and providing valuable data to the natural parks where videos were recorded. In terms of behaviors, we identified the most representative ones of the selected species, in order to cover as much as possible the range of activities developed by the birds.

Data annotation

Accurate annotation of the captured data is a determining factor in obtaining relevant results when training deep learning models on this data. To ensure annotation accuracy, the use of annotation tools^{31,32} has been extended, as they provide a user-friendly interface that makes this process easy and accessible to non-technical staff.

There are many open-source annotation tools available on the market. CVAT⁵ is one of the most popular ones, as it provides annotation support for images and videos, including a variety of formats for exporting the data. VoTT⁶ is also popular when annotating videos, as it offers multiple annotation shapes and integration with Microsoft services to easily upload data to Azure.⁷ Other simpler annotation tools are labelme⁸ or LabelImg,⁹ which are aimed at annotating images and their capabilities are more limited. For our purpose, we decided to use CVAT because of the large number of exportable formats available, the great collaborative environment it offers, and its easy integration with semi-automatic and automatic annotation processes.

As the need for larger amounts of data to train deep learning models increased, researchers began to enhance annotation tools with automatic systems that could alleviate this task. Annotation tools integrate machine learning models³³ that can automatically infer what would otherwise be manually annotated. Common tasks performed by automated annotation tools are object detection³⁴ and semantic segmentation.³⁵ While the former predicts the bounding box and class of each object in the image, the latter predicts regions of interest associated with specific categories.

Although automated annotation systems have demonstrated strong performance, semi-automated annotation processes are ultimately used because they ensure the creation of highly accurate annotations while greatly reducing the amount of human intervention required. Semi-automated annotation studies are widely used in the medical field,^{36,37} where precision is a key factor throughout the design.

In this study, a semi-automated annotation approach was followed, based on CVAT and its possible integration with powerful computer vision models. Our approach consisted of five main steps: Species classification, bird localization, behavior classification, subject identification, data curation, and post-processing. Each of these stages is described in more detail below. Figure 4 shows this process.

⁵<https://github.com/cvat-ai/cvat>

⁶<https://github.com/microsoft/VoTT>

⁷<https://azure.microsoft.com/>

⁸<https://github.com/labelmeai/labelme>

⁹<https://github.com/HumanSignal/labelImg>

1. **Species classification:** In this first step, the ecologists labeled each video with the main bird species that appeared. The main species is that of the bird in the focus of the camera. This way, annotations of birds that are different from the main species will not be included in the video annotations.
2. **Bird localization:** Then, an object detection model is used to predict the localization of the bounding boxes of the birds that appear in each of the video frames. For ease of implementation, YOLOv7³⁸ was chosen as the object detection model because it is predefined integration into CVAT. Since the model provided by CVAT is trained on general purpose data, the class predicted by default for each bounding box is not be the bird species, but the class *bird*. To avoid manually changing all the bounding box classes, we used an option provided by CVAT to associate a user-defined class with the class detected by the model. In this way, the class *bird* was associated with the species appearing in the video. Since each video is restricted to having only one species of bird, it is possible to associate *bird* with a specific species.
3. **Behavior classification:** After automatically annotating the bounding boxes, our team of ecologists performed the second step of the annotation process, which is twofold. First, they checked and corrected erroneous bounding boxes, and second, they annotated for each bounding box the behavior performed by each bird. To annotate the behaviors, CVAT bounding boxes *tags* were used.
4. **Subject identification:** When using automatic annotation models such as YOLOv7, CVAT does not support bounding box correspondence between frames. In other words, if a video shows two birds developing different behaviors, there is no relationship between the bounding boxes of adjacent frames, so it is not possible to analyze the birds' behaviors. This is not possible because the next frame will show two new bounding boxes whose relation to the one being analyzed is not known. To solve this problem, the Euclidean distance³⁹ was used to correlate bounding boxes of adjacent frames. The euclidean distance calculates the distance between the centers of the bounding boxes of adjacent frames and then correlates the bounding boxes with the minimum distance. The center of the bounding box was calculated as follows:

$$c = \left(\frac{x_{\min} + x_{\max}}{2}, \frac{y_{\min} + y_{\max}}{2} \right) \quad (1)$$

Given the centre of the bounding boxes, the Euclidean distance was calculated as:

$$d(c_1, c_2) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2} \quad (2)$$

$$\begin{aligned} \text{where } c_1 &= (x_1, y_1) \\ \text{and } c_2 &= (x_2, y_2) \end{aligned}$$

5. **Data curation:** After the labeling of species, bounding boxes, behaviors, and subjects, an overall review of all annotations was conducted to ensure the high quality of the data. To conduct the review, videos were assigned to all ecologists equally.
6. **Post-processing:** Once the annotations were complete, their format was adapted to make them easy to use and understand. To achieve this goal, the approach used in the AVA-Kinetics dataset⁴⁰ was followed. In this approach, a CSV file was used to contain annotations containing localized behaviors of multiple subjects. To export the data into the CSV format, the data was first exported from CVAT using the CVAT Video 1.1 format. Some Python scripts were then used to extract only the relevant information from the exported data and dump it into the output CSV file.

Data Records

The dataset presented in this study is open access and accessible through Zenodo¹⁰. Within this Zenodo repository, there are four main elements:

- **Videos folder:** This folder contains the 178 videos that comprise the dataset. Videos are identified by their name, which is composed of a numeric value and the species that appears in the video. The format is the following "ID-VIDEO.SPECIES-NAME.mp4".

¹⁰<https://zenodo.org/records/14355257>

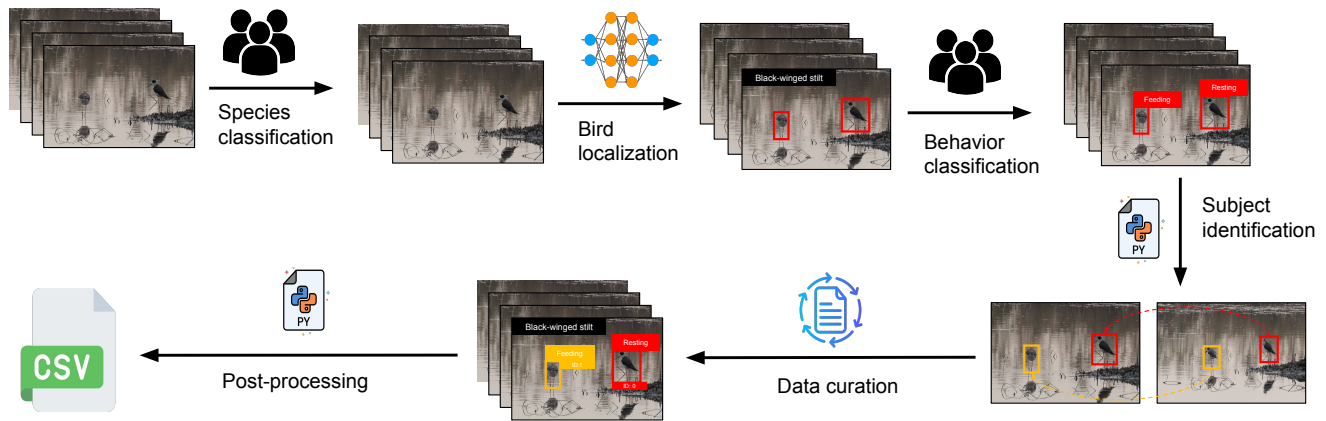


Figure 4. Visual representation of stages involved in the annotation process. Birds are first classified into species by annotators and localized using a YOLO model, then annotators recognize bird behaviors and subjects are identified using a Python script, and finally the data is curated and post-processed.

- **Bounding boxes CSV:** The *bounding_boxes.csv* file contains all the annotations of the data set. It follows a format of 10 columns, ordered as follows: Global identifier of the row, video identifier, frame identifier within the video, activity identifier, subject identifier, species appearing in the video, and the four coordinates of the bounding box (top-left x-coordinate, top-left y-coordinate, bottom-right x-coordinate, and bottom-right y-coordinate). Each of the CSV rows represents the information of one bounding box within one frame of a video.
- **Behavior identifiers CSV:** The *behavior_ID.csv* file contains a mapping of the seven behavior classes that make up the data set and their numeric identifiers.
- **Species identifiers CSV:** The file *species_ID.csv* contains a mapping between the 13 different bird species and their numerical identifier.
- **Splits JSON:** The *splits.json* file contains the videos associated with each train, validation and test split.

Technical Validation

To ensure high quality recordings and accurate annotations, the entire process was carried out by expert ecologists, as mentioned in the Data Annotation section. Firstly, video recordings were supervised by a group of experts who set up camera traps in strategic areas and also manually recorded some high-quality videos. For each video, these experts annotated the species appearing in the video. The same experts then corrected bounding box errors and annotated bird behavior, together with a number of collaborators with a background in ecology. Finally, a final stage of cross-checking of annotations was carried out by the experts and collaborators. The expertise of the annotators responsible for collecting and annotating the videos, together with the final cross-review process, ensures the quality and cleanliness of the data.

As this dataset has been conceived mainly to be used in deep learning pipelines, baseline pipelines will be given to demonstrate the applicability of the data presented in this work within deep learning workflows. As mentioned previously, the purpose of this dataset is twofold, as it provides annotation data for performing bird species detection and behavior classification tasks. Thus, one baseline per each task was developed. PyTorch¹¹ was used in both cases as coding platform.

Species classification

First, a baseline pipeline for species classification was developed. This baseline is based on a YOLOv9⁴¹ model trained over 50 epochs in the proposed dataset.

For efficient training, a downsampling of 10 is performed on the frames extracted from the videos. This can be done without affecting the performance of the model, as the difference between successive frames is minimal. The frames were extracted while maintaining the source FPS (Frames Per Second) of each video. During the training stage, a learning rate of 0.01 was used and a GeForce RTX 3090 GPU was used as the hardware platform. The test results from the baseline are shown next:

¹¹<https://pytorch.org>

Model	Precision	Recall	mAP50	mAP50-95
YOLOv9	0.835	0.759	0.801	0.556

Table 3. Results of the YOLO-based baseline developed for bird species classification.

The table 3 shows the test results for species classification in terms of precision, recall, mAP50 and mAP50-95 metrics. mAP50-95 is a common object recognition metric that refers to the mAP (mean Average Precision) computed over 10 different IoU (Intersection over Union) thresholds, specifically from 0.50 to 0.95 in increments of 0.05. The results show how YOLOv9 achieves a strong performance for the task, with a maximum precision of 0.835, while maintaining a high recall of 0.759. The mAP metrics, which measure the accuracy of bounding box localizations, are also strong, reaching 0.801 for mAP50 and 0.556 for mAP50-95, which is a good result considering the challenge of obtaining strong mAP scores when high IoU values are used.

To provide a more comprehensive understanding of the evaluation, the confusion matrix for the results is given. Figure 5 shows the confusion matrix, where it can be observed that the majority of the errors are due to the confusion of the ground truth class with the background class.

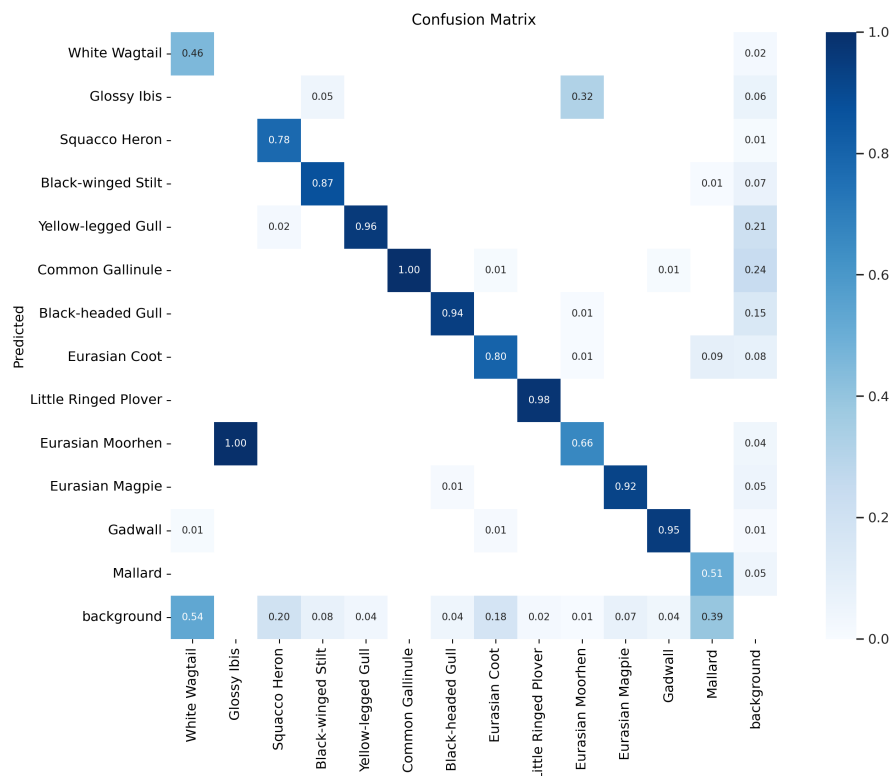


Figure 5. Confusion matrix of species classification pipeline.

Behavior detection

Secondly, the behavior detection baseline is presented. In this baseline, four different video classification models were trained end-to-end to perform the behavior classification task. The trained models were Video MViT,⁴² Video S3D,⁴³ Video SwinTransformer⁴⁴ and Video ResNet.⁴⁵ All model architectures and pretrained weights were extracted from PyTorch.

For the training stage, the input videos were downsampled with a downsample rate of 3, selecting the first frame as the representative of each set (*i.e.* only the first frame of each set of 3 is kept). Train, test and validation splits were generated from the full set of videos with a distribution 70-15-15. The splits were constructed using a stratified strategy based on the species and behaviors appearing in the videos. The distribution was computed taking into account the number of frames which constitute each video (*e.g.* 1 video with 1000 frames is equivalent to 5 videos with 200 frames). Regarding the training hyperparameters, a learning rate tuning was conducted using a uniform sampling strategy with minimum and maximum values of 0.0001 and 0.01, respectively. Similarly to the species classification baseline, training was performed on a GeForce RTX

3090 GPU. The results for each model are shown below:

Model	Learning rate	Accuracy
MViT ⁴²	0.005	0.51
S3D ⁴³	0.005	0.29
SwinTransformer ⁴⁴	0.009	0.51
Video ResNet ⁴⁵	0.003	0.56

Table 4. Results of the baseline models for behavior detection in terms of accuracy. The learning rate shown is the one that achieved the highest accuracies during hyperparameter tuning..

From the Table 4 it can be concluded that the Video ResNet model is the one which learns better the complexity of the dataset, showcasing a maximum performance of 0.56. Conversely, the model with the lowest score is the S3D model, with an accuracy of 0.29. These results show the challenge posed by the dataset under study, which presents a limited amount of data. The limited amount of data available to train complex deep learning models demonstrates the need for more resources to capture more data. Furthermore, the development of new training strategies and deep learning architectures that fit the data needs should be explored in order to improve the baseline results obtained.

Usage Notes

Since the data annotations are provided in CSV format, it is recommended to use Python libraries such as Pandas,¹² which is specifically designed to read and manage CSV data. In the GitHub repository containing the code, there are usage examples of how to load and prepare the data to be fed into deep learning models. It is recommended to read the *dataset.py* script in the *behavior_detection* directory as an example.

Code availability

The data processing and experimentation code shown in the Technical Validation section is available on GitHub¹³. The GitHub repository is organized into two main directories. The *species_classification* directory contains all the code related to the species classification, and the *behavior_detection* directory contains the experiment with the behavior detection models proposed for the dataset.

References

1. O’Riordan, T. *Environmental Science for Environmental Management* (Longman, 1995).
2. Nichols, J. D. & Williams, B. K. Monitoring for conservation. *Trends ecology & evolution* **21**, 668–673 (2006).
3. Hays, G. C. *et al.* Translating marine animal tracking data into conservation policy and management. *Trends ecology & evolution* **34**, 459–473 (2019).
4. Margules, C. & Usher, M. Criteria used in assessing wildlife conservation potential: a review. *Biol. conservation* **21**, 79–109 (1981).
5. Smallwood, K. S., Beyea, J. & Morrison, M. L. Using the best scientific data for endangered species conservation. *Environ. Manag.* **24**, 421–435 (1999).
6. Morrison, M. L. Bird populations as indicators of environmental change. In *Current Ornithology: Volume 3*, 429–451 (Springer, 1986).
7. Bonebrake, T. C., Christensen, J., Boggs, C. L. & Ehrlich, P. R. Population decline assessment, historical baselines, and conservation. *Conserv. Lett.* **3**, 371–378 (2010).
8. Carvalho, S. B., Brito, J. C., Crespo, E. J. & Possingham, H. P. From climate change predictions to actions—conserving vulnerable animal groups in hotspots at a regional scale. *Glob. Chang. Biol.* **16**, 3257–3270 (2010).
9. Joseph, L. N., Maloney, R. F. & Possingham, H. P. Optimal allocation of resources among threatened species: a project prioritization protocol. *Conserv. biology* **23**, 328–338 (2009).

¹²<https://pandas.pydata.org/>

¹³<https://github.com/3dperceptionlab/Visual-WetlandBirds>

10. Nuttall, M. N. *et al.* Long-term monitoring of wildlife populations for protected area management in southeast asia. *Conserv. Sci. Pract.* **4**, e614 (2022).
11. Fraixedas, S. *et al.* A state-of-the-art review on birds as indicators of biodiversity: Advances, challenges, and future directions. *Ecol. Indic.* **118**, 106728 (2020).
12. Stastny, J., Munk, M. & Juranek, L. Automatic bird species recognition based on birds vocalization. *EURASIP J. on Audio, Speech, Music. Process.* **2018**, 1–7 (2018).
13. Vellinga, W.-P. & Planqué, R. The xeno-canto collection and its relation to sound recognition and classification. In *CLEF (Working Notes)* (2015).
14. Kahl, S. *et al.* Overview of birdclef 2023: Automated bird species identification in eastern africa. In *CLEF (Working Notes)*, 1934–1942 (2023).
15. Miyaguchi, A., Cheung, A., Gustineli, M. & Kim, A. Transfer learning with pseudo multi-label birdcall classification for ds@gt birdclef 2024 (2024). [2407.06291](https://arxiv.org/abs/2407.06291).
16. Morfi, V., Bas, Y., Pamuła, H., Glotin, H. & Stowell, D. Nips4bplus: a richly annotated birdsong audio dataset. *PeerJ Comput. Sci.* **5**, e223 (2019).
17. Lostanlen, V., Salamon, J., Farnsworth, A., Kelling, S. & Bello, J. P. Birdvox-full-night: A dataset and benchmark for avian flight call detection. In *2018 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, 266–270 (IEEE, 2018).
18. Huang, Y.-P. & Basanta, H. Bird image retrieval and recognition using a deep learning platform. *IEEE Access* **7**, 66980–66989, [10.1109/ACCESS.2019.2918274](https://doi.org/10.1109/ACCESS.2019.2918274) (2019).
19. Wah, C., Branson, S., Welinder, P., Perona, P. & Belongie, S. The caltech-ucsd birds-200-2011 dataset. Tech. Rep. CNS-TR-2011-001, California Institute of Technology (2011).
20. Van Horn, G. *et al.* Building a bird recognition app and large scale dataset with citizen scientists: The fine print in fine-grained dataset collection. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 595–604, [10.1109/CVPR.2015.7298658](https://doi.org/10.1109/CVPR.2015.7298658) (2015).
21. Alados, C. L., Escos, J. M. & Emlen, J. Fractal structure of sequential behaviour patterns: an indicator of stress. *Animal Behav.* **51**, 437–443 (1996).
22. Lindell, C. A. The value of animal behavior in evaluations of restoration success. *Restor. Ecol.* **16**, 197–203 (2008).
23. Berger-Tal, O. *et al.* A systematic survey of the integration of animal behavior into conservation. *Conserv. Biol.* **30**, 744–753 (2016).
24. Goldenberg, S., Douglas-Hamilton, I., Daballen, D. & Wittemyer, G. Challenges of using behavior to monitor anthropogenic impacts on wildlife: a case study on illegal killing of african elephants. *Animal Conserv.* **20**, 215–224 (2017).
25. Harvey, S. Deepbird: A deep learning pipeline for wildlife camera data analysis (2019).
26. Ng, X. L. *et al.* Animal kingdom: A large and diverse dataset for animal behavior understanding (2022). [2204.08129](https://arxiv.org/abs/2204.08129).
27. Rose, P. *et al.* Evaluation of the time-activity budgets of captive ducks (anatidae) compared to wild counterparts. *Appl. Animal Behav. Sci.* **251**, 105626 (2022).
28. Fontúrbel, F. E., Orellana, J. I., Rodríguez-Gómez, G. B., Tabilo, C. A. & Castaño-Villa, G. J. Habitat disturbance can alter forest understory bird activity patterns: A regional-scale assessment with camera-traps. *For. Ecol. Manag.* **479**, 118618 (2021).
29. Fontúrbel, F. E. *et al.* Sampling understory birds in different habitat types using point counts and camera traps. *Ecol. Indic.* **119**, 106863 (2020).
30. Murphy, A. J. *et al.* Using camera traps to examine distribution and occupancy trends of ground-dwelling rainforest birds in north-eastern madagascar. *Bird Conserv. Int.* **28**, 567–580 (2018).
31. Arandjelovic, M. *et al.* Highly precise community science annotations of video camera-trapped fauna in challenging environments. *Remote. Sens. Ecol. Conserv.* (2024).
32. Aljabri, M., AlAmir, M., AlGhamdi, M., Abdel-Mottaleb, M. & Collado-Mesa, F. Towards a better understanding of annotation tools for medical imaging: a survey. *Multimed. tools applications* **81**, 25877–25911 (2022).
33. Guillermo, M. *et al.* Implementation of automated annotation through mask rcnn object detection model in cvat using aws ec2 instance. In *2020 IEEE REGION 10 CONFERENCE (TENCON)*, 708–713, [10.1109/TENCON50793.2020.9293906](https://doi.org/10.1109/TENCON50793.2020.9293906) (2020).

34. Kiyokawa, T., Tomochika, K., Takamatsu, J. & Ogasawara, T. Fully automated annotation with noise-masked visual markers for deep-learning-based object detection. *IEEE Robotics Autom. Lett.* **4**, 1972–1977, [10.1109/LRA.2019.2899153](https://doi.org/10.1109/LRA.2019.2899153) (2019).
35. Pavoni, G. *et al.* Taglab: Ai-assisted annotation for the fast and accurate semantic segmentation of coral reef orthoimages. *J. field robotics* **39**, 246–262 (2022).
36. Krenzer, A. *et al.* Fast machine learning annotation in the medical domain: a semi-automated video annotation tool for gastroenterologists. *BioMedical Eng. OnLine* **21**, 33 (2022).
37. Li, H. *et al.* A semi-automated annotation algorithm based on weakly supervised learning for medical images. *Biocybern. Biomed. Eng.* **40**, 787–802 (2020).
38. Wang, C.-Y., Bochkovskiy, A. & Liao, H.-Y. M. Yolov7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors (2022). [2207.02696](https://arxiv.org/abs/2207.02696).
39. Sadli, R., Afkir, M., Hadid, A., Rivenq, A. & Taleb-Ahmed, A. Aggregated euclidean distances for a fast and robust real-time 3d-mot. *IEEE Sensors J.* **21**, 21872–21884, [10.1109/JSEN.2021.3104390](https://doi.org/10.1109/JSEN.2021.3104390) (2021).
40. Li, A. *et al.* The ava-kinetics localized human actions video dataset (2020). [2005.00214](https://arxiv.org/abs/2005.00214).
41. Wang, C.-Y. & Liao, H.-Y. M. YOLOv9: Learning what you want to learn using programmable gradient information (2024).
42. Li, Y. *et al.* Mvitv2: Improved multiscale vision transformers for classification and detection (2022). [2112.01526](https://arxiv.org/abs/2112.01526).
43. Xie, S., Sun, C., Huang, J., Tu, Z. & Murphy, K. Rethinking spatiotemporal feature learning: Speed-accuracy trade-offs in video classification (2018). [1712.04851](https://arxiv.org/abs/1712.04851).
44. Liu, Z. *et al.* Video swin transformer (2021). [2106.13230](https://arxiv.org/abs/2106.13230).
45. Tran, D. *et al.* A closer look at spatiotemporal convolutions for action recognition (2018). [1711.11248](https://arxiv.org/abs/1711.11248).

Acknowledgements

We would like to thank "A way of making Europe" European Regional Development Fund (ERDF) and MCIN/AEI/10.13039/501100011033 for supporting this work under the "CHAN-TWIN" project (grant TED2021-130890B-C21 and HORIZON-MSCA-2021-SE-0 action number: 101086387, REMARKABLE, Rural Environmental Monitoring via ultra wide-Area networkS And distriButed federated Learning. This work is part of the HELEADE project (TSI-100121-2024-24), funded by Spanish Ministry of Digital Processing and by the European Union NextGeneration EU. This work has also been supported by three Spanish national and two regional grants for PhD studies, FPU21/00414, FPU22/04200, FPU23/00532, CIACIF/2021/430 and CIACIF/2022/175.

Author contributions statement

J.R.J. conceived the manuscript, J.R.J., D.O.P., M.B.L., D.M.P. and P.R.P. conducted the experiments and analysed the results, A.O.T. and E.S.G. managed the data collection and annotation, J.G.R. led the project management. All authors contributed to the annotation of the data and reviewed the manuscript.

Competing interests

The authors declare no competing interests.