# Towards Application-Specific Evaluation of Vision Models: Case Studies in Ecology and Biology

Alex Hoi Hang Chan\*1,2,3, Otto Brookes\*4,5, Urs Waldmann<sup>1,6</sup>, Hemal Naik<sup>1,2,7</sup>, Iain D. Couzin<sup>1,2,3</sup> Majid Mirmehdi<sup>4</sup>, Noël Adiko Houa<sup>5</sup>, Emmanuelle Normand<sup>5</sup>, Christophe Boesch<sup>5</sup>, Lukas Boesch<sup>5</sup> Mimi Arandjelovic<sup>9</sup>, Hjalmar Kühl<sup>8</sup>, Tilo Burghardt<sup>†4</sup>, Fumihiro Kano<sup>†1,2,3</sup>

<sup>1</sup>Centre for the Advanced Study of Collective Behaviour, University of Konstanz, Germany
 <sup>2</sup>Department of Collective Behavior, Max Planck Institute of Animal Behavior, Germany
 <sup>3</sup>Department of Biology, University of Konstanz, Germany
 <sup>4</sup>School of Computer Science, University of Bristol, United Kingdom
 <sup>5</sup>Wild Chimpanzee Foundation, Germany

<sup>6</sup>Department of Computer and Information Science, University of Konstanz, Konstanz, Germany
<sup>7</sup>Department of Ecology of Animal Societies, Max Planck Institute of Animal Behavior, Konstanz, Germany
<sup>8</sup>Senckenberg Museum of Natural History Goerlitz, Goerlitz, Germany
<sup>9</sup>Max Planck Institute for Evolutionary Anthropology, Germany.
\*,†contributed equally.

hoi-hang.chan@uni-konstanz.de, otto.brookes@bristol.ac.uk

#### **Abstract**

Computer vision methods have demonstrated considerable potential to streamline ecological and biological workflows, with a growing number of datasets and models becoming available to the research community. However, these resources focus predominantly on evaluation using machine learning metrics, with relatively little emphasis on how their application impacts downstream analysis. We argue that models should be evaluated using applicationspecific metrics that directly represent model performance in the context of its final use case. To support this argument, we present two disparate case studies: (1) estimating chimpanzee abundance and density with camera trap distance sampling when using a video-based behaviour classifier and (2) estimating head rotation in pigeons using a 3D posture estimator. We show that even models with strong machine learning performance (e.g., 87% mAP) can yield data that leads to substantial discrepancies in abundance estimates compared to expert-derived data. Similarly, the highest-performing models for posture estimation do not produce the most accurate inferences of gaze direction in pigeons. Motivated by these findings, we call for researchers to integrate application-specific metrics in ecological/biological datasets, allowing for models to be benchmarked in the context of their downstream application and to facilitate better integration of models into application workflows.

Computer vision methods have demonstrated considerable potential to streamline ecological and biological workflows, with a growing number of datasets and models becoming available to the research community. However, these resources focus predominantly on evaluation using machine learning metrics, with relatively little emphasis on how their application impacts downstream analysis. We argue that models should be evaluated using applicationspecific metrics that directly represent model performance in the context of its final use case. To support this argument, we present two disparate case studies: (1) estimating chimpanzee abundance and density with camera trap distance sampling when using a video-based behaviour classifier and (2) estimating head rotation in pigeons using a 3D posture estimator. We show that even models with strong machine learning performance (e.g., 87% mAP) can yield data that leads to substantial discrepancies in abundance estimates compared to expert-derived data. Similarly, the highest-performing models for posture estimation do not produce the most accurate inferences of gaze direction in pigeons. Motivated by these findings, we call for researchers to integrate application-specific metrics in ecological/biological datasets, allowing for models to be benchmarked in the context of their downstream application and to facilitate better integration of models into application workflows. searchers that stands to benefit both fields.

#### 1. Introduction

Computer vision (CV) has emerged as a powerful tool in ecology and biology, with the potential to dramatically reduce the effort required for data extraction [50]. These advantages have important implications for many disciplines, including biodiversity monitoring [47], conservation [50], behavioural ecology [13], neurobiology [38] and more.

In recent years, there has been significant development of datasets and algorithms [4, 5, 9, 18, 28, 32, 36, 37, 40–42, 46, 52] for animal studies. While these algorithms are rigorously evaluated using machine learning (ML) metrics (e.g., accuracy, mean average precision (mAP), root mean squared error (RMSE)), their impact on ecological and biological research remains largely underexplored. Only a small number of studies have examined how CV solutions translate to real-world applications and their effectiveness in deriving meaningful measurements in downstream analysis [8, 11, 20, 41, 44, 47, 53]. However, the emergence of application-driven ML learning [49] and its adoption by major conferences [12, 24, 43] is beginning to bridge this gap and it is a timely opportunity to evaluate the practical usability of ML models in their intended applications.

In this paper, we argue models should be evaluated using application-specific metrics that reflect effectiveness in their intended ecological or biological use cases, in addition to standard ML metrics. To support this argument, we present two case studies with both methods of evaluation. First, we evaluate a behaviour classification model for removing video segments containing behaviours known to bias population abundance estimates calculated using camera trap distance sampling (CTDS). We demonstrate that; (a) highperforming ML models (e.g., 87% mAP) can output data that, when used in CTDS, produce abundance estimates that differ significantly from those obtained through manual expert annotation and; (b) the way model predictions interact with statistical approaches, such as CTDS, is difficult to predict. Second, we employ 3D posture estimation algorithms to infer attention in pigeons by estimating head rotation. We find that; (c) there is a mismatch between the accuracy of keypoint estimation with rotation estimation, and; (d) mismatches can result in misleading conclusions about which model is the most appropriate for the final application.

To better align models with their intended use case, we recommend that application-specific metrics are included in existing/future datasets and reported on when proposing new models, as complementary metrics to existing ML benchmarks. Building on recent calls for closer interdisciplinary collaboration [47, 50], we highlight application-specific metrics as a practical and accessible starting point for joint efforts between ML and ecology/biology re-

#### 2. Related Work

As early as 2012, Wagstaff [51] highlighted over-reliance on benchmark datasets, lack of method interpretatability, and limited capacity for real-world application in the field of machine learning. More recent commentaries have highlighted limitations with the current benchmarking system [35], and need for alternative evaluation procedures [48]. Recently, Rolnick et al. [49] formulated a paradigm to separate ML research into methods-driven ML and application-driven ML. The latter is intended to overcome the aforementioned limitations by promoting closer collaboration with end users, conservative use of standardized datasets, and evaluation metrics tailored to specific domains. In biological and ecological studies, there is increasing interest in evaluating ML models in an application-specific context [8, 11, 20, 23, 29, 41, 44, 47, 53].

Ecology. Whytock et al. [53] trained the ResNet50 architecture to classify 26 Central African mammal and bird species with a top-1 accuracy of 77.63% and showed that predictions were equivalent to expert labelling when considering species richness, occupancy and activity patterns. Pantazis et al. [45] assess the impact of model architecture, label noise, and dataset size on the same ecological metrics, showing that model architecture is of minimal importance although noise and dataset size is significant. Henrich et al. [20] perform semi-automatic extraction of distances between animal and camera with a deep learning model and predict population densities for 10 species using CTDS. They show densities predicted with manual and semi-automatic processes do not differ significantly.

Biology. Chistensen et al. [11] highlight that ML metrics (F1-score) of a behavioral classifier from accelerometer data can be misleading, and even low accuracy can lead to a reliable inference of behavioural states in wild birds. Benchmarking results from the BuckTales dataset [41] suggest that commonly used MOT metric, MOTA, is not suitable to evaluate animal tracking algorithms, as it does not consider ID switches, which is highly relevant for biological analysis. Recent studies have also evaluated ML-based pipelines by their ability to predict simple biological hypotheses in animal behaviour [8, 11], or detect differences in medical treatment groups, e.g [23, 29]. Specifically for animal posture estimation, to the best of our knowledge, there is limited work evaluating the relationship between the accuracy of keypoint estimation models and downstream applications, despite the widespread application of the method.

## 3. Case study 1: Abundance & Density Estimation of Chimpanzees

Overview. CTDS is a common approach to estimating the abundance of a species through camera trapping [16, 17, 21, 22]. It measures distances from cameras to detected animals and models detection probability as a function of distance to estimate population abundance and density. At its core, it is a detection function, which describes how the probability of detecting an animal decreases with increasing distance from the camera. The data required by CTDS is difficult to obtain, necessitating the annotation of intravideo animal identities, species, distances from the camera at regular time intervals, and behaviour in camera trap videos.

In this case study, we automate one stage of the data acquisition pipeline: behaviour recognition. Camera reactivity behaviour is known to bias abundance estimates by causing chimps to remain in (i.e., attraction) or out (i.e., avoidance) of the cameras view, which can lead to overestimation and underestimation of population abundances, respectively [17,21]. Identifying and removing video clips where camera reactivity is observed currently represents the most effective method to debias estimates [17]. Below, we describe the training and evaluation of a behaviour recognition model, and evaluate its ability to identify clips in which chimpanzees exhibit camera reactivity and the effect of using it on abundance estimates.

#### 3.1. Experimental Setup

Model Implementation & Training. We train UniformerV2 [34] on a modified version of the PanAf20k dataset [6] to classify camera reactivity (i.e., presence or absence of reaction) of chimpanzees in camera trap videos. To achieve this, we convert the labels of the dataset into a binary format: videos annotated with camera reactivity are assigned a positive label, while those without are assigned a negative label. We follow the training protocol detailed in [6], except we use a class balanced focal loss [14] to mitigate the effect of class imbalance (a 70:30 class imbalance between non-reactive and reactive videos exists). Macroaveraged mean average precision (mAP) is used for model evaluation.

Abundance Estimation & Model Evaluation. To evaluate impact on downstream abundance estimates a different dataset comprising 413 videos from the Taï National Park, Côte d'Ivoire, is used. Each video is manually annotated with all the information required to generate abundance estimates using CTDS [21]. The calculations are performed following the approach outlined by Howe et al. [22]. The trained model is applied over 2 second segments of all videos (i.e., the same time intervals used during manual annotation), and clips classified as containing camera reactivity are removed prior to CTDS. Note that inference

is performed in an out-of-distribution (OOD) setting since the camera trap locations here are unseen during training.

#### 3.2. Results

First, the performance of the behaviour recognition model is reported, as evaluated using typical ML metrics (mAP). Then, the abundance estimates produced using the model filtered video footage are compared with those calculated using expert annotation and removal. Results are reported for two different detection functions: half-normal (Hn1) and hazard rate (Hr1). Hn1 assumes a smooth decline in detection, while Hr1 allows for a flat detection zone followed by a steep drop. Note that all other data required for distance sampling (i.e., intra-video individual identities, species, distance from camera to individual etc.) are produced expertly. Thus, only the effect of automated behaviour classification on CTDS is examined.

**Result 1: Behaviour Recognition Performance.** The model achieves an mAP score of 87.82%, although there is a noticeable discrepancy in class-wise performance. While it effectively detects the absence of camera reactivity with 95.31% AP, performance on detecting presence is lower at 73.95%. This discrepancy highlights the significance of class imbalance *even* when imbalances are not extreme.

Result 2: Abundance & Density Estimation. If camera reactivity clips are not removed prior to CTDS then population abundance and densities are significantly overestimated (Tab. 1; 2575 vs. 1680 and 2529 vs. 1954 for Hr1 and Hn1, respectively). Across both detection functions, estimates are highest when no removal of camera reactivity clips is performed (row 1 and 4 vs. rest). Tab. 1 also shows that using the behaviour classifier to remove camera reactivity clips still results in over estimation of abundance/density, across both detection functions. However, the degree of abundance overestimation is far greater in Hr1 (+20.77%) than Hn1 (+14.38%). Density is also overestimated but there is relatively little difference in the increase between detection functions. This difference shows that even when the same model predictions are used to remove clips, the resulting abundance estimates can vary depending on the choice of detection function. This highlights that it is not just the number of correctly classified segments that matters — but which segments are correctly identified and removed.

#### 4. Case study 2: Gaze Estimation in Pigeons

**Overview.** Vision is a primary sensory channel for many species, and the direction of attention—where an animal directs its eyes or head (gaze)—provides crucial insights into how animals acquire information about their environments. In this case study, we evaluate markerless 3D posture estimation methods by comparing their absolute Euclidean position errors with angular head orientation errors, the latter being a more meaningful measure for assessing the accuracy of inferred gaze direction. Similar to hu-

	<b>Detection Function</b>	Removal	Camera Reaction	Abundance			Density (ind./km2)		
				Estimate; SE	95%	CIs	Estimate; SE	95%	CIs
1		None	Yes	2575; 921	1122	4981	0.47; 0.16	0.24	0.80
2	Hr1	Manual	No	1680; 755	578	3235	0.34; 0.14	0.13	0.64
3		Auto	No	2029; 767	703	3658	0.36; 0.17	0.17	0.68
4		None	Yes	2529; 892	1313	4332	0.49; 0.18	0.21	0.83
5	Hn1	Manual	No	1954; 943	693	4015	0.36; 0.16	0.12	0.71
6		Auto	No	2235; 909	745	3806	0.38; 0.15	0.16	0.63

Table 1. Chimpanzee density & abundance estimations calculated with CTDS under two different detection functions (Hr1 and Hn1), with camera reactivity clips removed from videos using both manual (i.e., human expert) and automated (i.e., behaviour classifier) methods. The automated method leads to significant overestimates of abundance when compared with manual removal.

ML Metrics					Application Specific Metrics					
Method	RMSE (mm)	Median (mm)	PCK05 (%)	PCK10 (%)	RMSE Angles (°)	Median Angles (°)	Median Yaw (°)	Median Pitch (°)	Median Roll (°)	
3D-KP-RCNN	22.8	5.46	79.9	91.7	18.2	4.23	4.22	3.42	5.60	
3D-DLC*	21.1	5.28	82.4	93.1	13.8	3.34	2.37	3.11	4.23	
3D-ViTPose*	21.2	5.09	83.5	92.9	15.7	4.17	5.38	3.21	4.17	
LToHP	15.7	4.12	89.3	96.1	9.3	3.61	3.63	2.82	4.34	

Table 2. 3D posture estimation benchmarks presented in 3D-MuPPET [52], compared with angular error of the head. 3D estimates are obtained by first detecting keypoints in 2D from multiple views, then triangulated into 3D. Only head keypoints are used to compute error metrics. 3D-KP-RCNN: Keypoint RCNN model [19], 3D-DLC\*: first detecting pigeons using YOLOv8l [27], then keypoints using DeepLabCut [37], 3D-ViTPose\*: first detecting pigeons using YOLOv8l, then keypoints using ViTPose [55]. LToHP: Learnable triangulation of human postures framework [25], used as a baseline model for comparison. Bold denotes best model for given metric.

mans, many bird species have foveas—specialized regions in the eye where visual acuity is highest [2]. Recent studies have shown that by measuring a bird's 3D head rotation, researchers can infer its gaze direction using both marker-based [15,26,30,39] and markerless motion capture systems [10,40]. When combined with simple behavioral experiments that estimate the location and extent of the bird's visual field during object-directed attention, head orientation can serve as a reliable proxy for gaze direction [7,30,31].

#### 4.1. Experimental setup

Dataset & Evaluation Metrics. Here, we use the 3D-POP dataset [40], a large scale 2D to 3D posture dataset in pigeon flocks, to reproduce benchmarks introduced in the 3D-MuPPET framework [52] for 3D posture estimation. In 3D-MuPPET, the authors provide benchmark results. We refer to original publication for implementation details [52]. in terms of RMSE and percentage correct keypoints (PCK). RMSE is calculated by taking the root mean squared of the Euclidean distances between all predicted and ground truth keypoints, while PCK is the percentage of predicted keypoints that fall within 5% or 10% of the maximum length of any two ground truth 3D keypoints [52]. In the current study, we only report these metrics for four head keypoints - not the nine keypoints reported in [52], so the comparison is fair with the alternative metric we propose below.

We propose an alternative metric that is more aligned to the biological use case, simply the absolute rotational error of the head. We employ the same method shown in [10] to calculate the rotation error of the head of each pigeon in each frame, in terms of the yaw, pitch and roll axis relative to the forward axis (axis from the mid-point between the eyes towards the beak). The acceptable rotational error in head orientation depends on the spatial arrangement of gaze targets, such as conspecifics or objects of interest (e.g. food) within the observation arena. In most experimental setups, these targets are typically separated by more than  $10^{\circ}$ . Moreover, pigeons generally move their eyes within a narrow range—typically no more than  $5^{\circ}$ —when attending to distant objects [30,54]. Therefore, a head orientation error within  $5^{\circ}$  from 3D posture tracking is considered an acceptable margin for reliably estimating gaze direction.

#### 4.2. Results

We refer to Table 2 for the comparison between ML and application-specific metrics. Firstly, we find that based on the median head rotation error, all of the models compared are below 5°, confirming that the 3D-MuPPET framework is appropriate for the gaze estimation task. In this context, any of the architectures would have been within the acceptable accuracy, so the end user can consider other factors like inference speed or trajectory tracking accuracy when deciding which architecture to use (see [52]).

Next, we show that the model with the best performance in terms of Euclidean distance errors and PCK—LToHP—does not correspond to the best performance when evaluated using head rotation angle, where 3D-DLC\* performs better (Table 2). This discrepancy appears to be driven by 3D-DLC\*'s lower accuracy in estimating yaw (i.e., horizontal head rotation), which is important if gaze targets are arranged on the same eye-level as the pigeon, which is typically the case. However, we also find that LToHP yields the lowest RMSE in head rotation angle,

suggesting it may produce fewer outliers. In the original 3D-MuPPET publication [52], the authors recommend 3D-ViTPose\* as the most accurate model (excluding LToHP). Yet, based on our evaluation using an application-specific metric, 3D-DLC\* appears to perform best. This further highlights the discrepancy between standard ML metrics and application-specific evaluation.

#### 5. Discussion & Conclusion

We presented two case studies that evaluate computer vision models using ML and application-specific metrics. In both cases we report misalignment between metrics, demonstrating that current ML-focused benchmarking is not always optimized to facilitate effective downstream application of models, and complementary mechanisms to evaluate them in this context are needed. Case Study 1 shows that model-derived abundance estimates of chimpanzees diverges notably from those produced using manual annotations, despite high ML performance. Similarly, Case Study 2 shows that the best model for 3D posture estimation is not the best model for estimating head rotation. This highlights the potential for end-users to be misled by impressive ML benchmarks or invest substantial resources (e.g., data collection, annotation, compute [50]) optimising model performance that fails to translate in practice.

Further analysis of *why* ML and application-specific metrics can be mismatched is needed, though many plausible explanations exist. For example, in Case Study 1 the model is applied in an OOD setting (see Sec. 3.1) which may degrade performance and ultimately contribute to the mismatch between metrics – a key challenge in ML [1,3,33]. In case study 2, rotation estimation is likely more sensitive than position estimation, as rotational errors can be compounded in unpredictable ways by small errors in point estimates. These challenges are often found in real-world deployments, and these aspects deserve consideration when developing and benchmarking SOTA algorithms.

Finally, while we present case studies where ML and application-specific metrics can be mismatched, we acknowledge that in certain cases, these metrics can also be aligned. However, this does not downplay the importance of introducing complementary application-specific metrics in current benchmarking pipelines, to better bridge novel algorithms to downstream applications.

**Future Directions** To work toward solutions of the presented challenges, we recommend introducing application-specific metrics in existing and new CV datasets with biological/ecological relevance, as complementary benchmarks to traditional ML metrics. While datasets can be used to solve many different problems, researchers could introduce complementary metrics for core problems relating to the proposed application domain of the dataset. This can hopefully encourage novel models to benchmark comple-

mentary metrics that better represent the efficacy of a model for downstream applications. Finally, similar to [47,50], we encourage more interdisciplinary collaboration between CV researchers and ecologists/biologists, to better align computer vision solutions to real-life problems, which will lead to more holistic solutions as the field continues to grow.

### Acknowledgments

This work is funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany's Excellence Strategy—EXC 2117—422037984, and DFG project number 15990824. U.W. acknowledges funding from the Connected Minds Program, supported by Canada First Research Excellence Fund, Grant #CFREF-2022-00010. This work was supported by the UKRI CDT in Interactive AI (grant EP/S022937/1).

#### References

- [1] Sara Beery, Grant Van Horn, and Pietro Perona. Recognition in terra incognita. In <u>Proceedings of the European conference on computer vision (ECCV)</u>, pages 456–473, 2018. 5
- [2] Andreas Bringmann. Structure and function of the bird fovea. <u>Anatomia, histologia, embryologia</u>, 48(3):177–200, 2019. 4
- [3] Otto Brookes, Maksim Kukushkin, Majid Mirmehdi, Colleen Stephens, Paula Dieguez, Thurston C. Hicks, Sorrel Jones, Kevin Lee, Maureen S. McCarthy, Amelia Meier, Emmanuelle Normand, Erin G. Wessling, Roman M. Wittig, Kevin Langergraber, Klaus Zuberbühler, Lukas Boesch, Thomas Schmid, Mimi Arandjelovic, Hjalmar Kühl, and Tilo Burghardt. Panaffgbg: Understanding the impact of backgrounds in wildlife behaviour recognition. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2025. 5
- [4] Otto Brookes, Majid Mirmehdi, Hjalmar Kuhl, and Tilo Burghardt. Chimpvlm: Ethogram-enhanced chimpanzee behaviour recognition. 2024. 2
- [5] Otto Brookes, Majid Mirmehdi, Hjalmar S. Kühl, and Tilo Burghardt. Triple-stream deep metric learning of great ape behavioural actions. In <u>Proceedings of the 18th International Joint Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications</u>, pages 294–302, 2023. 2
- [6] Otto Brookes, Majid Mirmehdi, Colleen Stephens, Samuel Angedakin, Katherine Corogenes, Dervla Dowd, Paula Dieguez, Thurston C Hicks, Sorrel Jones, Kevin Lee, et al. Panaf20k: a large video dataset for wild ape detection and behaviour recognition. <u>International Journal of Computer Vision</u>, 132(8):3086–3102, 2024. 3
- [7] Shannon R Butler, Jennifer J Templeton, and Esteban Fernández-Juricic. How do birds look at their world? a novel avian visual fixation strategy. <u>Behavioral ecology and sociobiology</u>, 72:1–11, 2018. 4

- [8] Alex Hoi Hang Chan, Jingqi Liu, Terry Burke, William D Pearse, and Julia Schroeder. Comparison of manual, machine learning, and hybrid methods for video annotation to extract parental care data. <u>Journal of Avian Biology</u>, 2024(3-4):e03167, 2024. 2
- [9] Jun Chen, Ming Hu, Darren J Coker, Michael L Berumen, Blair Costelloe, Sara Beery, Anna Rohrbach, and Mohamed Elhoseiny. Mammalnet: A large-scale video benchmark for mammal recognition and behavior understanding. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 13052–13061, 2023.
- [10] Michael Chimento, Alex Hoi Hang Chan, Lucy M Aplin, and Fumihiro Kano. Peering into the world of wild passerines with 3d-socs: synchronized video capture and posture estimation. bioRxiv, pages 2024–06, 2024. 4
- [11] Charlotte Christensen, André C Ferreira, Wismer Cherono, Maria Maximiadi, Brendah Nyaguthii, Mina Ogino, Daniel Herrera, and Damien R Farine. Moving towards more holistic validation of machine learning-based approaches in ecology and evolution. bioRxiv, 2024. 2
- [12] Conference on Computer Vision and Pattern Recognition. Conference on Computer Vision and Pattern Recognition (CVPR), 2025. Accessed: 2025-03-17.
- [13] Iain D Couzin and Conor Heins. Emerging technologies for behavioral research in changing environments. <u>Trends</u> in Ecology & Evolution, 2023. 2
- [14] Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. Class-balanced loss based on effective number of samples. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 9268–9277, 2019. 3
- [15] Mathilde Delacoux and Fumihiro Kano. Fine-scale tracking reveals visual field use for predator detection and escape in collective foraging of pigeon flocks. <u>Elife</u>, 13:RP95549, 2024. 4
- [16] Zackary J Delisle, Elizabeth A Flaherty, Mackenzie R Nobbe, Cole M Wzientek, and Robert K Swihart. Nextgeneration camera trapping: systematic review of historic trends suggests keys to expanded research applications in ecology and conservation. <u>Frontiers in Ecology and</u> Evolution, 9:617996, 2021. 3
- [17] Zackary J Delisle, Maik Henrich, Pablo Palencia, and Robert K Swihart. Reducing bias in density estimates for unmarked populations that exhibit reactive behaviour towards camera traps. <u>Methods in Ecology and Evolution</u>, 14(12):3100–3111, 2023. 3
- [18] Timm Haucke, Hjalmar S Kühl, Jacqueline Hoyer, and Volker Steinhage. Overcoming the distance estimation bottleneck in estimating animal abundance with camera traps. Ecological Informatics, 68:101536, 2022. 2
- [19] Kaiming He, Georgia Gkioxari, Piotr Dollar, and Ross Girshick. Mask r-cnn. In ICCV, 2017. 4
- [20] Maik Henrich, Mercedes Burgueño, Jacqueline Hoyer, Timm Haucke, Volker Steinhage, Hjalmar S Kühl, and Marco Heurich. A semi-automated camera trap distance sampling approach for population density estimation. <u>Remote Sensing in Ecology and Conservation</u>, 10(2):156– 171, 2024. 2

- [21] Noël Adiko Houa, Noémie Cappelle, Eloi Anderson Bitty, Emmanuelle Normand, Yves Aka Kablan, and Christophe Boesch. Animal reactivity to camera traps and its effects on abundance estimate using distance sampling in the taï national park, côte d'ivoire. PeerJ, 10:e13510, 2022. 3
- [22] Eric J Howe, Stephen T Buckland, Marie-Lyne Després-Einspenner, and Hjalmar S Kühl. Distance sampling with camera traps. <u>Methods in Ecology and Evolution</u>, 8(11):1558–1565, 2017. 3
- [23] Yujia Hu, Carrie R Ferrario, Alexander D Maitland, Rita B Ionides, Anjesh Ghimire, Brendon Watson, Kenichi Iwasaki, Hope White, Yitao Xi, Jie Zhou, et al. Labgym: Quantification of user-defined animal behaviors using learning-based holistic assessment. <u>Cell Reports Methods</u>, 3(3), 2023. 2
- [24] International Conference on Machine Learning. International Conference on Machine Learning (ICML), 2025. Accessed: 2025-03-17. 2
- [25] Karim Iskakov, Egor Burkov, Victor Lempitsky, and Yury Malkov. Learnable triangulation of human pose. In <u>ICCV</u>, 2019. 4
- [26] Akihiro Itahara and Fumihiro Kano. "corvid tracking studio": A custom-built motion capture system to track head movements of corvids. <u>Japanese Journal of Animal</u> Psychology, 72(1):1–16, 2022. 4
- [27] Glenn Jocher, Ayush Chaurasia, and Jing Qiu. YOLO by Ultralytics, Jan. 2023. 4
- [28] Peter Johanns, Timm Haucke, and Volker Steinhage. Automated distance estimation for wildlife camera trapping. Ecological Informatics, 70:101734, 2022. 2
- [29] Takaaki Kaneko, Jumpei Matsumoto, Wanyi Lu, Xincheng Zhao, Louie Richard Ueno-Nigh, Takao Oishi, Kei Kimura, Yukiko Otsuka, Andi Zheng, Kensuke Ikenaka, et al. Deciphering social traits and pathophysiological conditions from natural behaviors in common marmosets. <u>Current Biology</u>, 34(13):2854–2867, 2024. 2
- [30] Fumihiro Kano, Hemal Naik, Göksel Keskin, Iain D Couzin, and Máté Nagy. Head-tracking of freely-behaving pigeons in a motion-capture system reveals the selective use of visual field regions. <u>Scientific Reports</u>, 12(1):19113, 2022. 4
- [31] Fumihiro Kano, James Walker, Takao Sasaki, and Dora Biro. Head-mounted sensors reveal visual attention of free-flying homing pigeons. <u>Journal of experimental biology</u>, 221(17), 2018. 4
- [32] Maksim Kholiavchenko, Jenna Kline, Maksim Kukushkin, Otto Brookes, Sam Stevens, Isla Duporge, Alec Sheets, Reshma R Babu, Namrata Banerji, Elizabeth Campolongo, et al. Deep dive into kabr: a dataset for understanding ungulate behavior from in-situ drone video. <u>Multimedia Tools</u> and Applications, pages 1–20, 2024. 2
- [33] Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanas Phillips, Irena Gao, et al. Wilds: A benchmark of in-thewild distribution shifts. In <u>International conference on</u> machine learning, pages 5637–5664. PMLR, 2021. 5

- [34] Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Limin Wang, and Yu Qiao. Uniformerv2: Spatiotemporal learning by arming image vits with video uniformer. <u>arXiv</u> preprint arXiv:2211.09552, 2022. 3
- [35] Thomas Liao, Rohan Taori, Inioluwa Deborah Raji, and Ludwig Schmidt. Are we learning yet? a meta review of evaluation failures across machine learning. In <u>Thirty-fifth Conference on Neural Information Processing</u> <u>Systems Datasets and Benchmarks Track (Round 2)</u>, 2021.
- [36] Dan Liu, Jin Hou, Shaoli Huang, Jing Liu, Yuxin He, Bochuan Zheng, Jifeng Ning, and Jingdong Zhang. Loteanimal: A long time-span dataset for endangered animal behavior understanding. In <u>Proceedings of the IEEE/CVF</u> <u>International Conference on Computer Vision</u>, pages 20064– 20075, 2023. 2
- [37] Alexander Mathis, Pranav Mamidanna, Kevin M. Cury, Taiga Abe, Venkatesh N. Murthy, Mackenzie Weygandt Mathis, and Matthias Bethge. Deeplabcut: markerless pose estimation of user-defined body parts with deep learning. Nat. Neurosci., 21:1281–1289, 2018. 2, 4
- [38] Mackenzie Weygandt Mathis and Alexander Mathis. Deep learning tools for the measurement of animal behavior in neuroscience. <u>Current opinion in neurobiology</u>, 60:1–11, 2020. 2
- [39] Máté Nagy, Hemal Naik, Fumihiro Kano, Nora V Carlson, Jens C Koblitz, Martin Wikelski, and Iain D Couzin. Smart-barn: Scalable multimodal arena for real-time tracking behavior of animals in large numbers. <u>Science Advances</u>, 9(35):eadf8068, 2023. 4
- [40] Hemal Naik, Alex Hoi Hang Chan, Junran Yang, Mathilde Delacoux, Iain D Couzin, Fumihiro Kano, and Máté Nagy. 3d-pop—an automated annotation approach to facilitate markerless 2d-3d tracking of freely moving birds with markerbased motion capture. <a href="mailto:arXiv:2303.13174"><u>arXiv preprint arXiv:2303.13174</u></a>, 2023. 2, 4
- [41] Hemal Naik, Junran Yang, Dipin Das, Margaret Crofoot, Akanksha Rathore, and Vivek Hari Sridhar. Bucktales: A multi-uav dataset for multi-object tracking and re-identification of wild antelopes. <u>Advances in Neural</u> Information Processing Systems, 37:81992–82009, 2024. 2
- [42] Xun Long Ng, Kian Eng Ong, Qichen Zheng, Yun Ni, Si Yong Yeo, and Jun Liu. Animal kingdom: A large and diverse dataset for animal behavior understanding. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 19023–19034, 2022.
- [43] International Conference on Computer Vision. International conference on computer vision (iccv), 2025. Accessed: 2025-03-17. 2
- [44] Omiros Pantazis, Peggy Bevan, Holly Pringle, Guilherme Braga Ferreira, Daniel J. Ingram, Emily Madsen, Liam Thomas, Dol Raj Thanet, Thakur Silwal, Santosh Rayamajhi, Gabriel Brostow, Oisin Mac Aodha, and Kate E. Jones. Deep learning-based ecological analysis of camera trap images is impacted by training data quality and size, 2024. 2
- [45] Omiros Pantazis, Peggy Bevan, Holly Pringle, Guilherme Braga Ferreira, Daniel J Ingram, Emily Madsen, Liam

- Thomas, Dol Raj Thanet, Thakur Silwal, Santosh Rayamajhi, et al. Deep learning-based ecological analysis of camera trap images is impacted by training data quality and size. arXiv preprint arXiv:2408.14348, 2024. 2
- [46] Talmo D. Pereira, Nathaniel Tabris, Arie Matsliah, David M. Turner, Junyu Li, Shruthi Ravindranath, Eleni S. Papadoyannis, Edna Normand, David S. Deutsch, Z. Yan Wang, Grace C. McKenzie-Smith, Catalin C. Mitelut, Marielisa Diez Castro, John D'Uva, Mikhail Kislin, Dan H. Sanes, Sarah D. Kocher, Samuel S.-H. Wang, Annegret L. Falkner, Joshua W. Shaevitz, and Mala Murthy. Sleap: A deep learning system for multi-animal pose tracking. Nat. Methods, 19:486–495, 2022.
- [47] Laura J Pollock, Justin Kitzes, Sara Beery, Kaitlyn M Gaynor, Marta A Jarzyna, Oisin Mac Aodha, Bernd Meyer, David Rolnick, Graham W Taylor, Devis Tuia, et al. Harnessing artificial intelligence to fill global shortfalls in biodiversity knowledge. <a href="Nature Reviews Biodiversity">Nature Reviews Biodiversity</a>, pages 1–17, 2025. 2, 5
- [48] Inioluwa Deborah Raji, Emily M Bender, Amandalynne Paullada, Emily Denton, and Alex Hanna. Ai and the everything in the whole wide world benchmark. arxiv. <a href="mailto:arXiv">arXiv</a> preprint arXiv:2111.15366, 2021. 2
- [49] David Rolnick, Alan Aspuru-Guzik, Sara Beery, Bistra Dilkina, Priya L. Donti, Marzyeh Ghassemi, Hannah Kerner, Claire Monteleoni, Esther Rolf, Milind Tambe, and Adam White. Position: Application-Driven Innovation in Machine Learning. In Proceedings of the 41st International Conference on Machine Learning, pages 42707–42718. PMLR, July 2024. ISSN: 2640-3498.
- [50] Devis Tuia, Benjamin Kellenberger, Sara Beery, Blair R Costelloe, Silvia Zuffi, Benjamin Risse, Alexander Mathis, Mackenzie W Mathis, Frank Van Langevelde, Tilo Burghardt, et al. Perspectives in machine learning for wildlife conservation. <u>Nature communications</u>, 13(1):792, 2022. 2, 5
- [51] Kiri Wagstaff. Machine learning that matters. <u>arXiv preprint</u> arXiv:1206.4656, 2012. 2
- [52] Urs Waldmann, Alex Hoi Hang Chan, Hemal Naik, Máté Nagy, Iain D Couzin, Oliver Deussen, Bastian Goldluecke, and Fumihiro Kano. 3d-muppet: 3d multi-pigeon pose estimation and tracking. <u>International Journal of Computer</u> <u>Vision</u>, 132(10):4235–4252, 2024. 2, 4, 5
- [53] Robin C Whytock, Jędrzej Świeżewski, Joeri A Zwerts, Tadeusz Bara-Słupski, Aurélie Flore Koumba Pambo, Marek Rogala, Laila Bahaa-el din, Kelly Boekee, Stephanie Brittain, Anabelle W Cardoso, et al. Robust ecological analysis of camera trap data labelled by a machine learning model. <u>Methods in Ecology and Evolution</u>, 12(6):1080–1092, 2021.
- [54] Andreas Wohlschläger, Ralf Jäger, and Juan D Delius. Head and eye movements in unrestrained pigeons (columba livia). Journal of Comparative Psychology, 107(3):313, 1993. 4
- [55] Yufei Xu, Jing Zhang, Qiming Zhang, and Dacheng Tao. Vitpose: Simple vision transformer baselines for human pose estimation. Advances in neural information processing systems, 35:38571–38584, 2022. 4