

From Forest to Zoo: Great Ape Behavior Recognition with ChimpBehave

Michael Fuchs, Emilie Genty, Adrian Bangerter, Klaus Zuberbühler, Paul Cotofrei

University of Neuchâtel, Switzerland

{michael.fuchs, emilie.genty, adrian.bangerter, klaus.zuberbuehler, paul.cotofrei}@unine.ch

Abstract

This paper addresses the significant challenge of recognizing behaviors in non-human primates, specifically focusing on chimpanzees. Automated behavior recognition is crucial for both conservation efforts and the advancement of behavioral research. However, it is significantly hindered by the labor-intensive process of manual video annotation. Despite the availability of large-scale animal behavior datasets, the effective application of machine learning models across varied environmental settings poses a critical challenge, primarily due to the variability in data collection contexts and the specificity of annotations.

In this paper, we introduce ChimpBehave, a novel dataset featuring over 2 hours of video (approximately 193,000 video frames) of zoo-housed chimpanzees, meticulously annotated with bounding boxes and behavior labels for action recognition. ChimpBehave uniquely aligns its behavior classes with existing datasets, allowing for the study of domain adaptation and cross-dataset generalization methods between different visual settings. Furthermore, we benchmark our dataset using a state-of-the-art CNN-based action recognition model, providing the first baseline results for both within and cross-dataset settings. The dataset, models, and code can be accessed at: <https://github.com/MitchFuchs/ChimpBehave>

1. Introduction

The development of machine learning tools to recognize animal behaviors from videos plays a critical role in ecology and ethology. Automated systems for recognizing chimpanzee behaviors could offer a broad spectrum of applications, from enhancing conservation efforts to providing valuable insights into the behavior of great apes. Furthermore, non-invasive technologies developed for their well-being can significantly benefit chimpanzees, an endangered species, in both wild and captive settings. For example, these systems could monitor population dynamics in natural

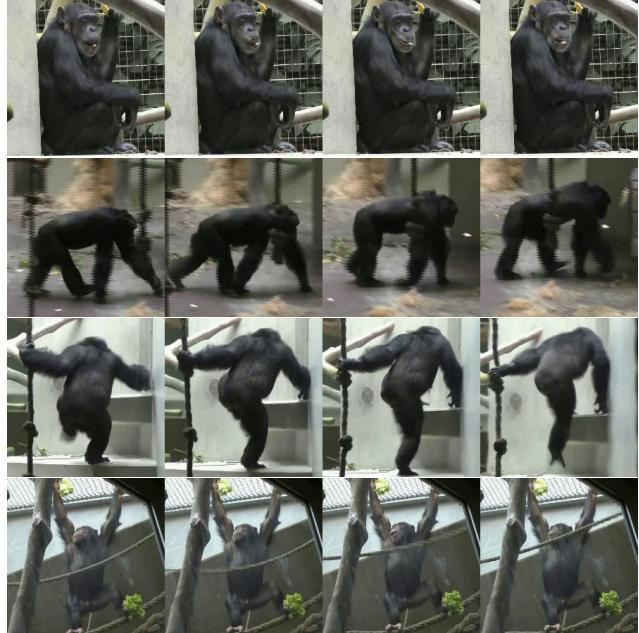


Figure 1. Walking, hanging, sitting, or climbing up? Identifying which chimpanzee behavior is depicted in these images is trivial for most humans. For algorithms, however, this is not always the case, especially when exposed to videos from previously unseen environments. **TL;DR:** We propose a new dataset and methods to investigate the classification of chimpanzee behaviors across different visual settings.

habitats or timely signal behavioral abnormalities in unwell individuals to caretakers in zoos.

As one of humans' closest living relatives, chimpanzees have been the subject of extensive scientific research in fields such as ecology, comparative cognition, neuroscience, and evolutionary biology. This research often relies on videos, whose manual annotation can be time-consuming and labor-intensive. The advancement of algorithms in animal behavior classification can therefore significantly benefit researchers by speeding up the labeling

process and/or reducing its overall cost.

Large animal datasets have recently been created to adapt human-centered action recognition models for animal behavior classification (see e.g., [5, 17]). Although comprehensive, these datasets lack the fine-grained annotations needed to capture the complex behaviors of great apes. To address this, more focused datasets like ChimpACT [13] and PanAf [4] have been created, targeting species-specific behaviors in different environments - from zoo settings to wild forests - and under distinct filming conditions (see Sec. 2.1 for details). While valuable, these datasets have two main issues. First, they typically have their own unique set of annotated behaviors. Second, they are often captured under very different conditions (e.g., fixed vs. moving cameras; zoo vs. forest environments). As a result, their joint exploitation for developing action recognition systems that work in diverse conditions is quite challenging.

To address these challenges, we make the following contributions:

- We introduce ChimpBehave, a dataset for great ape behavior recognition, which features over 2 hours of video (approximately 193,000 video frames) annotated with fine-grained behaviors and bounding boxes. Its label classes are specifically aligned with existing datasets to facilitate the study of domain adaptation techniques and cross-dataset generalization experiments.
- We establish a first comparative baseline using X3D [7], a state-of-the-art CNN-based model, in different scenarios in both within-dataset and cross-dataset settings. We use a rigorous evaluation protocol including a stratified 5-fold cross-validation procedure to validate our results.



Figure 2. Tracking example after interpolation for missing frames (red bounding box) and corrected ID. As the individual passes behind the pole, its track was lost and the ID was swapped with another.

2. Related work

2.1. Non-human primate datasets

The growing number of animal datasets designed for computer vision tasks has significantly included non-human primates (NHP), reflecting their significance across various biological and ecological studies. They span a wide array, from encompassing multiple animal orders [5, 12, 17, 21] to focusing specifically on primates [22], apes [4, 6, 13], monkeys [20], and particularly macaques [2, 10, 15]. These datasets showcase notable diversity in annotations and tasks, including species identification [4–6, 21, 22], animal detection and tracking [4, 12, 13, 21], pose estimation [2, 6, 10, 12, 13, 15, 17, 20–22], and behavior recognition [2, 4, 5, 12, 13, 15, 17].

Notably, ChimpACT [13] and PanAf [4] stand out for their focus on great ape behavior, each adopting distinct approaches and highlighting the challenges of applying computer vision across diverse settings. These datasets exhibit significant differences in both their visual environments and the scope of behaviors annotated. ChimpACT documents the life of a young chimpanzee within a zoo environment, characterized by man-made backgrounds, dynamic, hand-held camera work, and a longitudinal focus on a single individual. In contrast, PanAf offers insights into the lives of chimpanzees and gorillas in their natural habitats, featuring static cameras placed in African forests and capturing a wide array of ape populations. These varying contexts highlight the challenges that computer vision applications can encounter across diverse settings. Additionally, their behavior annotations diverge: ChimpACT captures broader locomotive behaviors and social interactions, whereas PanAf specifies actions like 'climbing up', 'running', and 'standing'. This disparity complicates cross-datasets analyses due to mismatches in annotated behaviors, challenging direct model comparison. ChimpBehave aims to bridge these gaps by combining a visual and filming setup reminiscent of ChimpACT with a behavioral annotation scheme aligned with PanAf's detailed action categories.

2.2. Behavior recognition for non-human primates

Studies on automated behavior recognition in non-human primates have primarily concentrated on macaques [2, 11, 14], monkeys [12] and apes [1, 3, 4, 9, 13, 19].

Key advancements in this domain leverage action recognition techniques, which are critical for classifying behaviors from video sequences. These techniques are categorized into three main approaches: video-based, skeleton-based and multi-modal. Video-based action recognition analyzes the visual features directly from the raw video data, capturing movements and interactions within the frame's pixel data [3, 4, 12–14, 19]. In contrast, skeleton-based

approaches focus on the movement of key body points or joints, constructing a skeletal representation of the subject to discern specific actions or behaviors [2, 9]. Additionally, some studies enhance behavior classification by also incorporating multimodal signals, such as audio cues in [1], providing a more comprehensive analysis by combining visual movement patterns with relevant sounds and/or vocalizations.

3. The ChimpBehave dataset

In this section, we present the dataset we have built and its main features.

Dataset Description: ChimpBehave consists of 1,121 annotated video segments, derived from 50 focal sampling videos recorded in 2016 at the Basel Zoo indoor enclosure (see Fig. 1 and other examples in Fig. 5 in Supplementary Materials). Each video follows one of nine chimpanzees to ensure focused observation of the individual, while still capturing its surroundings and other conspecifics. The filming conditions are naturalistic and include camera motion, zooming, and shaking, similar to [13]. All recordings were made with a hand-held camera in 1920x1080 pixel resolution at 25 fps.

Behavior Annotation: An expert primatologist (EG) meticulously labeled the dataset’s behavior annotations in all video segments using ELAN (<https://archive.mpi.nl/tla/elan>), leading to a total of approximately 193,000 annotated video frames. The annotator focused on seven behavioral classes which represent common primate behaviors, namely ‘sitting’, ‘standing’, ‘walking’, ‘running’, ‘hanging’, ‘climbing down’, and ‘climbing up’. These classes were intentionally selected to match those annotated in [4], in order to facilitate cross-dataset analysis.

Bounding Box Annotations: A single annotator (MF) labeled approximately 12,000 video frames from 131 video segments with bounding boxes at intervals of every tenth frame, with annotations for the remaining frames interpolated. All annotations were made in Label Studio (<https://labelstud.io/>) on video segments covering a variety of scenes. For each chimpanzee, a minimum of 500 frames were annotated, and for each of the 50 original focal videos, at least two segments or 100 frames were included. These labels were later used to fine-tune QDTrack [18] on the MMaction2 platform [16], a state-of-the-art Multiple Object Tracking (MOT) model pretrained on the ChimpACT dataset [13]. This particular model was selected due to its demonstrated adequacy for this task as highlighted in [13]. Fine-tuning was conducted for 2 epochs using 107 video segments, leaving out 24 video segments for evaluation purposes. The final model achieves the following scores on common MOT metrics at test time: Recall: 0.7640, Precision: 0.9740, HOTA: 0.6630, mAP: 0.8400.

Using this model, we then predicted tracking bounding boxes for all video segments, manually refined the tracks to fix any ID swaps, and interpolated the predictions where the model’s predictions were missing (Fig. 2). Each track was individually reviewed, and only valid sequences were kept as part of the final dataset. The code for MOT fine-tuning and data conversion between MMaction2 and Label Studio will be made available on our code repository. **Data**

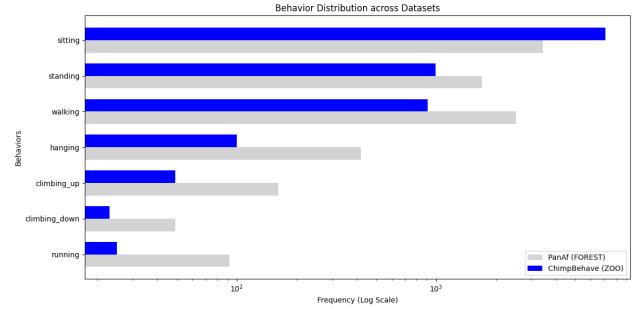


Figure 3. Behavior frequency distribution in the ChimpBehave and PanAf datasets, used in our experiments. Frequencies are plotted on a logarithmic scale to highlight the long-tailed characteristic of the data

Preprocessing: For our experiments, we extracted 9,187 unique *miniclip*s (a sequence of 20 consecutive frames, without overlap) from all video segments. All remaining frames in segments were discarded. The coordinates of the individual’s bounding boxes within the miniclip are used to calculate the global minimum and maximum coordinates of a cropping area, which is adjusted to ensure a minimum size of 224x224 pixels. This process ensures that each action recognition input includes 20 frames of size 224x224, associated with one behavioral class.

4. Method and experiments

Our goal is to study cross-data generalization performance of a standard action recognition model. To this end, we first describe the second dataset we have used, before presenting the tested behavior recognition model and our experimental protocol.

4.1. PanAf dataset

The Pan African Programme ‘The Cultured Chimpanzee’ [8] aimed to enhance the understanding of evolutionary-ecological factors influencing chimpanzee behavioral diversity. In its efforts, it collected numerous hours of footage from camera traps placed in the forests of Central Africa. From this collection, 500 videos of chimpanzees or gorillas, each 15 seconds long (180,000 frames at 24 FPS, resolution 720x404), were annotated with bounding box coordinates for ape detection and behaviors for action recogni-

tion [4, 19] (see image examples in Fig 6 in Supplementary Materials). The nine labeled behaviors include ‘walking’, ‘standing’, ‘sitting’, ‘running’, ‘hanging’, ‘climbing up’, ‘climbing down’, ‘sitting on back’, and ‘camera interaction’.

Data Preprocessing: We include in our experiments all videos that are fully annotated and extract from them 8,404 miniclip in a similar fashion as detailed in Sec. 3 from all behavioral classes except ‘camera interaction’ and ‘sitting on back’, which are absent in ChimpBehave.

4.2. Behavior Recognition Model

X3D [7], a state-of-the-art CNN-based model for action recognition, was chosen to conduct our comparative baseline. This model was selected due to its high performance in the PanAf benchmark as presented in [4]. Notably, X3D’s architecture expands incrementally from a small 2D image classification model along multiple network axes, including space, time, depth, and width. This network expansion process is designed to find an optimal balance between model complexity and performance, resulting in efficient models without compromising accuracy.

4.3. Experimental protocol

Evaluation Metrics: We evaluate each trained model using commonly used action recognition metrics: Top-1 Accuracy, Mean Class Accuracy, F1 (Micro), and Mean Average Precision (mAP). The Top-1 Accuracy metric measures the proportion of correctly classified miniclip, making it sensitive to imbalanced class distributions. F1 (Micro) is particularly useful for imbalanced datasets, as it aggregates the contributions of all classes to compute overall precision, recall, and F1 score, giving equal weight to each miniclip. Mean Class Accuracy evaluates the average accuracy across all classes, ensuring each class is equally represented. Mean Average Precision (mAP) provides a comprehensive measure by averaging the precision across different recall levels for each class, capturing the model’s overall ability to identify relevant miniclip. For deeper insights into model predictions, we render all confusion matrices in Supplementary Materials (see Fig. 8 to Fig. 11).

Cross-Validation Procedure: To validate all experimental results, we follow a standard stratified 5-fold cross-validation procedure. In this approach, each of the five models is trained on 4 folds (representing 80% of the dataset) and validated on the remaining 20%. This ensures that all dataset miniclip are used exactly once for validation. To ensure the exact same class distribution across folds, we compute the largest multiple of 5 in the total number of miniclip for each class and discard the rest. The miniclip for each fold were selected in the order of their appearance in the database, which was sorted by video name and then by frame numbering, to group them as much as

possible at the video level. This step was taken to ensure maximum generalization across videos, similar to the train/validation/test partitioning followed in [4].

Confidence Intervals: To assess statistical significance, we calculated 95% confidence intervals for each metric using the five evaluations from the 5-fold cross-validation procedure, applying a t-distribution ($\alpha = 0.025, \nu = 4$).

Cross-Dataset Evaluation: When referring to ‘cross-dataset’ in our evaluations, we imply that the model was trained on 80% of the miniclip from one of the datasets (as described above) and later tested on all miniclip of the second dataset.

Comparative Scenarios: As both datasets suffer from class distribution imbalance, we investigate two distinct scenarios referred to as **A** for **ALL** and **B** for **BALANCED**:

- **Scenario A** includes all miniclip from all seven classes.
- **Scenario B** only includes miniclip from the most frequent classes, namely ‘sitting’, ‘standing’, and ‘walking’. Additionally, in Scenario B, the number of miniclip kept for experiments in each class is set to the same value across classes and is determined by the number of miniclip in the least frequent class of the three. For the two other classes, a set of miniclip is randomly selected from all its miniclip.

Experimental details. Each model training and evaluation was conducted on the MMAction2 platform [16] using default hyperparameters. Each X3D model was trained on the HPC cluster of the University of Neuchâtel, on 4x NVIDIA RTX 2080 Ti (4x 11GB) for 50 epochs, using *SGD* optimization, with an initial learning rate of 0.01. The final epoch was selected based on its Top1 accuracy on the validation set.

5. Results and discussion

To emphasize the uniqueness of the visual features of our two datasets and to simplify discussions, in this Section we will refer to the ChimpBehave dataset as ZOO and the PanAf dataset as FOREST.

Within dataset results. Our results are shown in Fig. 4.

When comparing scores between ZOO and FOREST, the model performs statistically significantly better in both scenarios for ZOO with respect to Top-1 Accuracy and F1 Score (micro). Both of these metrics set equal importance on all miniclip. A different pattern emerges when examining mean class accuracy and mAP, two metrics whose final scores are averaged by the number of classes. For these two metrics, we observe no significant difference between ZOO and FOREST in scenario A. This can be partially explained by the long-tailed characteristic of both datasets, where the model systematically fails to classify classes such as ‘climbing up’, ‘climbing down’, and ‘running’ across all folds, as seen in Fig. 8 in Supplementary Materials. However, the model performs significantly better in ZOO across all eval-

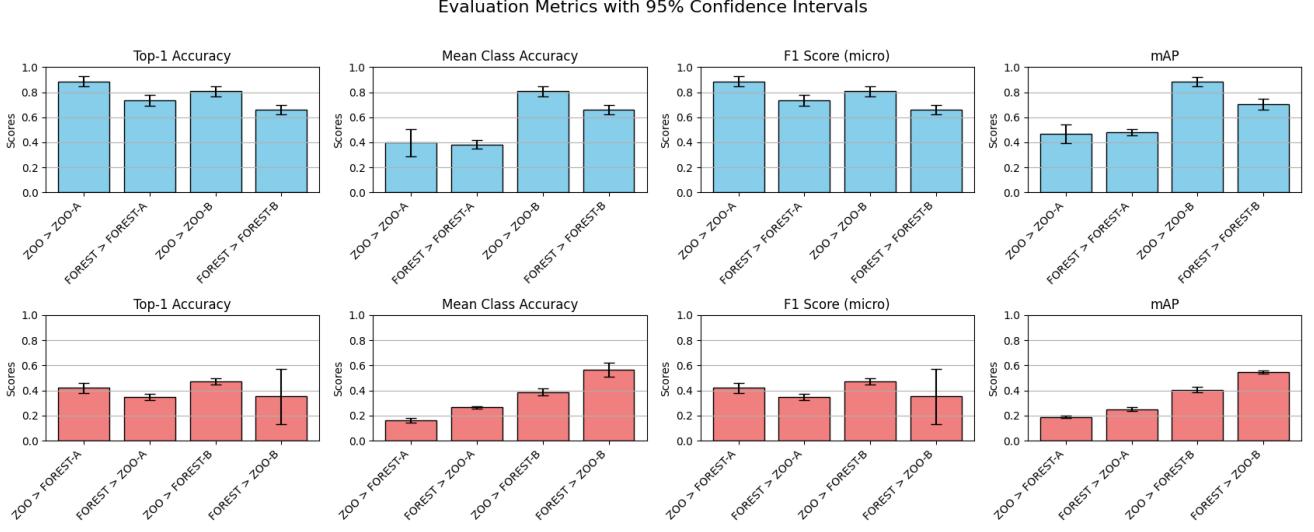


Figure 4. Evaluation metrics on both datasets across both scenarios. The top row (blue) represents the scores within-dataset, while the bottom row (red) represents the scores in cross-dataset settings.

uation metrics, indicating that it learns to recognize visual features in ZOO more easily.

We can identify three major factors that may contribute to this difference. First, when visually comparing sequences of images from both datasets in Fig. 7 in Supplementary Materials, one can observe higher contrast between the individual and its background in ZOO. Second, the overall image resolution is markedly higher in ZOO (1920x1080) compared to FOREST (720x404). Third, the relative size of individuals can be much smaller in FOREST, as the camera is fixed, whereas focal sampling videos in ZOO follow each individual.

Evaluating generalization: cross-dataset results. The model’s generalization capabilities show contrasting patterns, depending on the evaluation metric. For mean class accuracy and mean average precision, the ZOO→FOREST models obtained statistically significant lower performances than the FOREST→ZOO models in both scenarios. Conversely, for Top-1 Accuracy and F1-score, the ZOO→FOREST models achieved statistically significant better performances than the FOREST→ZOO models, but only in scenario A. Additionally, the performances of all models in the cross-dataset context (independent of scenario or evaluation metric) are approximately 50% of the same performances in the within-dataset context.

6. Limitations and future work

As highlighted in the previous sections, ChimpBehave suffers from class distribution imbalance, which we aim to address in future work. In this regard, spatio-temporal action detection models, such as those proposed in [13],

could facilitate future data annotation efforts. On another note, while pose estimation is currently a much-researched topic in animal behavior recognition, especially for NHP, our current results only provide a baseline for video-based approaches. In the future, our work could be extended to include a skeleton-based action recognition baseline using methods similar to the one proposed in [9].

7. Conclusion

In this paper, we introduced ChimpBehave, a new dataset which we hope will be useful for future research in automated behavior classification of non-human primates. By specifically aligning the labeling scheme with existing datasets, we were able to demonstrate how to perform cross-dataset evaluations and provide an initial baseline for future research. We hope this will prove helpful for the conservation, study, understanding, and well-being of chimpanzees and great apes in general.

8. Acknowledgement

We extend our gratitude to the Basel Zoo, its staff, and its director, Adrian Baumeyer, for granting us the opportunity to conduct our data collection within their facilities. Furthermore, we wish to express our appreciation to NCCR Evolving Language, Swiss National Science Foundation Agreement #51NF40_180888 and Grant No. CR31I3_166331 awarded to A.B. and K.Z. for their financial support in data collection and annotation. Special thanks are owed to the members of the *SIG Ape Gestures*, including Daphné Bavelier, Richard Hahnloser, Nianlong Gu, and Remo Nitschke, for their contributions.

Thank you to Jean-Marc Odobez, at Idiap Research Institute, for his dedicated help.

9. Ethical Statement

We received ethical agreement for this study from the Commission d’Ethique de la Recherche of the University of Neuchâtel (agreement number: 01-FS-2017) and the Kantonales Veterinäramt BS at Basel Zoo.

References

- [1] Max Bain, Arsha Nagrani, Daniel Schofield, Sophie Berdugo, Joana Bessa, Jake Owen, Kimberley J. Hockings, Tetsuro Matsuzawa, Misato Hayashi, Dora Biro, Susana Carvalho, and Andrew Zisserman. Automated audiovisual behavior recognition in wild primates. *Science Advances*, 7(46):eabi4883, 2021. [2](#), [3](#)
- [2] Praneet C Bala, Benjamin R Eisenreich, Seng Bum Michael Yoo, Benjamin Y Hayden, Hyun Soo Park, and Jan Zimmermann. Automated markerless pose estimation in freely moving macaques with openmonkeystudio. *Nature communications*, 11(1):1–12, 2020. [2](#), [3](#)
- [3] Otto Brookes, Majid Mirmehdi, Hjalmar Kühl, and Tilo Burghardt. Triple-stream deep metric learning of great ape behavioural actions. *arXiv preprint arXiv:2301.02642*, 2023. [2](#)
- [4] 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. *International Journal of Computer Vision*, pages 1–17, 2024. [2](#), [3](#), [4](#)
- [5] 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 (CVPR)*, pages 13052–13061, 2023. [2](#)
- [6] Nisarg Desai, Praneet Bala, Rebecca Richardson, Jessica Raper, Jan Zimmermann, and Benjamin Hayden. Openapepose: a database of annotated ape photographs for pose estimation. *arXiv preprint arXiv:2212.00741*, 2022. [2](#)
- [7] Christoph Feichtenhofer. X3d: Expanding architectures for efficient video recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 203–213, 2020. [2](#), [4](#)
- [8] Max Planck Institute for Evolutionary Anthropology. Pan african programme: The cultured chimpanzee. [3](#)
- [9] Michael Fuchs, Emilie Genty, Klaus Zuberbühler, and Paul Cotofrei. Asbar: an animal skeleton-based action recognition framework. recognizing great ape behaviors in the wild using pose estimation with domain adaptation. *bioRxiv*, pages 2023–09, 2023. [2](#), [3](#), [5](#)
- [10] Rollyn Labuguen, Jumpei Matsumoto, Salvador Blanco Negrete, Hiroshi Nishimaru, Hisao Nishijo, Masahiko Takada, Yasuhiro Go, Ken-ichi Inoue, and Tomohiro Shiba. Macaquepose: A novel “in the wild” macaque monkey pose dataset for markerless motion capture. *Frontiers in behavioral neuroscience*, 14:581154, 2021. [2](#)
- [11] Chuxi Li, Zifan Xiao, Yerong Li, Zhinan Chen, Xun Ji, Yiqun Liu, Shufei Feng, Zhen Zhang, Kaiming Zhang, Jian-feng Feng, et al. Deep learning-based activity recognition and fine motor identification using 2d skeletons of cynomolgus monkeys. *Zoological Research*, 44(5):967, 2023. [2](#)
- [12] Dan Liu, Jin Hou, Shaoli Huang, Jing Liu, Yuxin He, Bochuan Zheng, Jifeng Ning, and Jingdong Zhang. Loternal: A long time-span dataset for endangered animal behavior understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 20064–20075, 2023. [2](#)
- [13] Xiaoxuan Ma, Stephan Kaufhold, Jiajun Su, Wentao Zhu, Jack Terwilliger, Andres Meza, Yixin Zhu, Federico Rossano, and Yizhou Wang. Chimpact: A longitudinal dataset for understanding chimpanzee behaviors. In *Advances in Neural Information Processing Systems*, pages 27501–27531. Curran Associates, Inc., 2023. [2](#), [3](#), [5](#)
- [14] Markus Marks, Qiuhan Jin, Oliver Sturman, Lukas von Ziegler, Sepp Kollmorgen, Wolfgang von der Behrens, Valerio Mante, Johannes Bohacek, and Mehmet Fatih Yanik. Deep-learning-based identification, tracking, pose estimation and behaviour classification of interacting primates and mice in complex environments. *Nature machine intelligence*, 4(4):331–340, 2022. [2](#)
- [15] Lucas M. Martini, Anna Bognár, Rufin Vogels, and Martin A. Giese. Macaction: Realistic 3d macaque body animation based on multi-camera markerless motion capture. *bioRxiv*, 2024. [2](#)
- [16] MMAAction2 Contributors. OpenMMLab’s Next Generation Video Understanding Toolbox and Benchmark, 2020. [3](#), [4](#)
- [17] 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. [2](#)
- [18] Jiangmiao Pang, Linlu Qiu, Xia Li, Haofeng Chen, Qi Li, Trevor Darrell, and Fisher Yu. Quasi-dense similarity learning for multiple object tracking. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 164–173, 2021. [3](#)
- [19] Faizaan Sakib and Tilo Burghardt. Visual recognition of great ape behaviours in the wild. 2021. International Conference on Pattern Recognition (ICPR) Workshop on Visual Observation and Analysis of Vertebrate And Insect Behavior , VAIB ; Conference date: 10-01-2021 Through 15-01-2021. [2](#)
- [20] Zheng Sun, Xiangyu Zhu, Zhen Lei, and Xibo Ma. Caged monkey dataset: A new benchmark for caged monkey pose estimation. In *Pattern Recognition and Computer Vision*, pages 694–706, Cham, 2022. Springer Nature Switzerland. [2](#)
- [21] Yuxiang Yang, Yingqi Deng, Yufei Xu, and Jing Zhang. Aptv2: Benchmarking animal pose estimation and tracking with a large-scale dataset and beyond, 2023. [2](#)

- [22] Yuan Yao, Praneet Bala, Abhiraj Mohan, Eliza Bliss-Moreau, Kristine Coleman, Sienna M Freeman, Christopher J Machado, Jessica Raper, Jan Zimmermann, Benjamin Y Hayden, et al. Openmonkeychallenge: Dataset and benchmark challenges for pose estimation of non-human primates. *International Journal of Computer Vision*, 131(1): 243–258, 2023. [2](#)

From Forest to Zoo: Great Ape Behavior Recognition with ChimpBehave

Supplementary Material

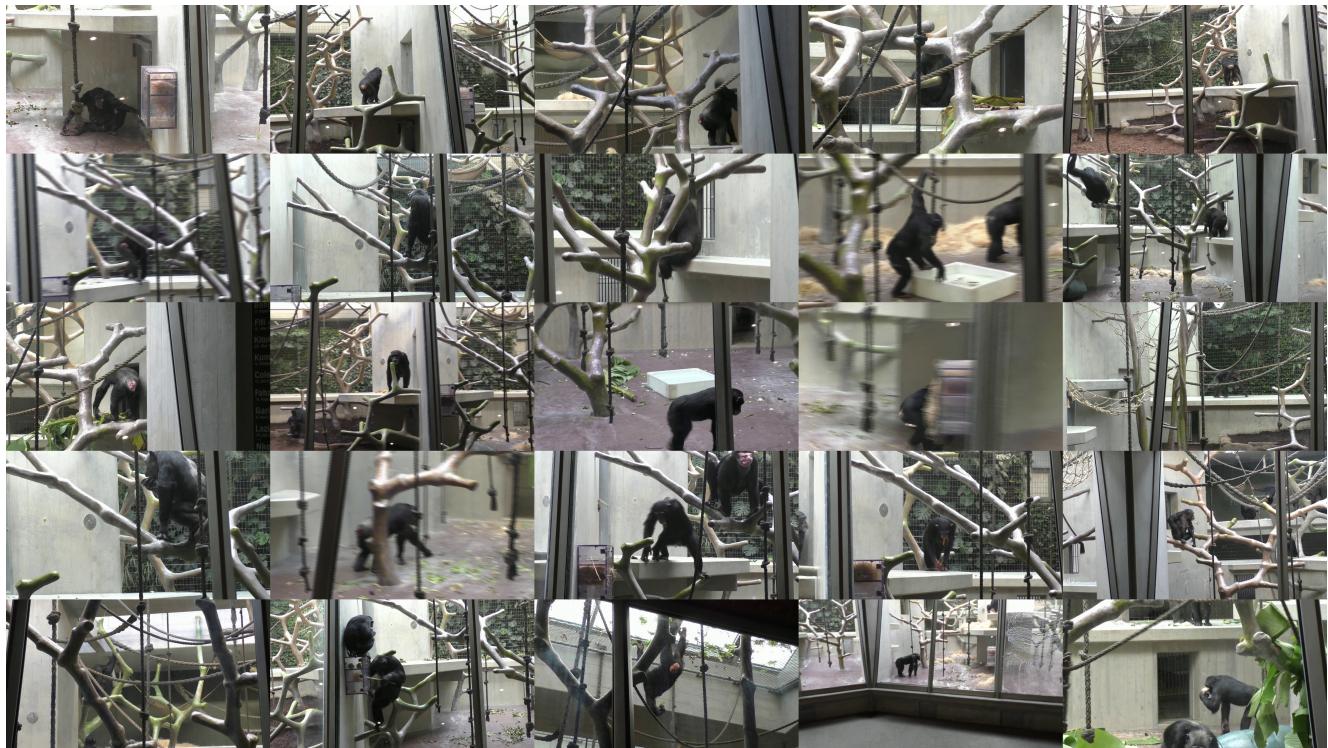


Figure 5. Image examples of ChimpBehave

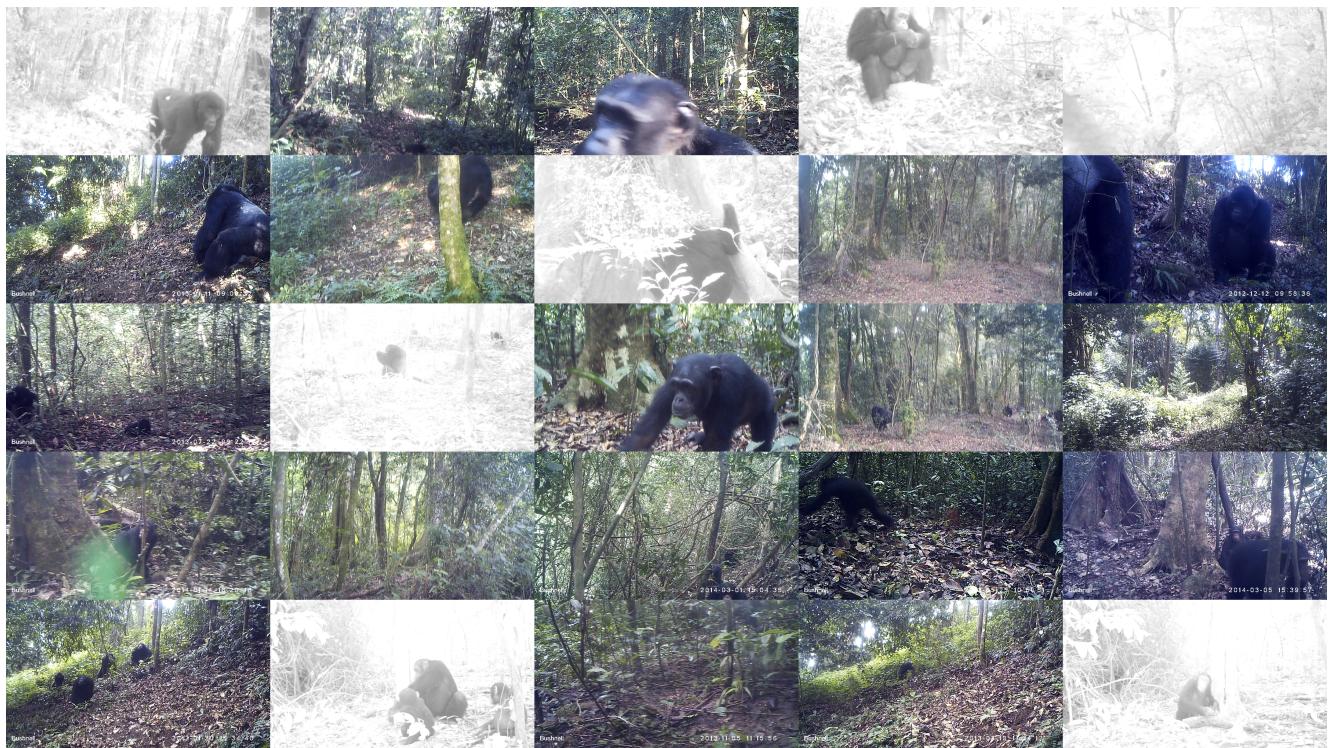


Figure 6. Image examples of PanAf

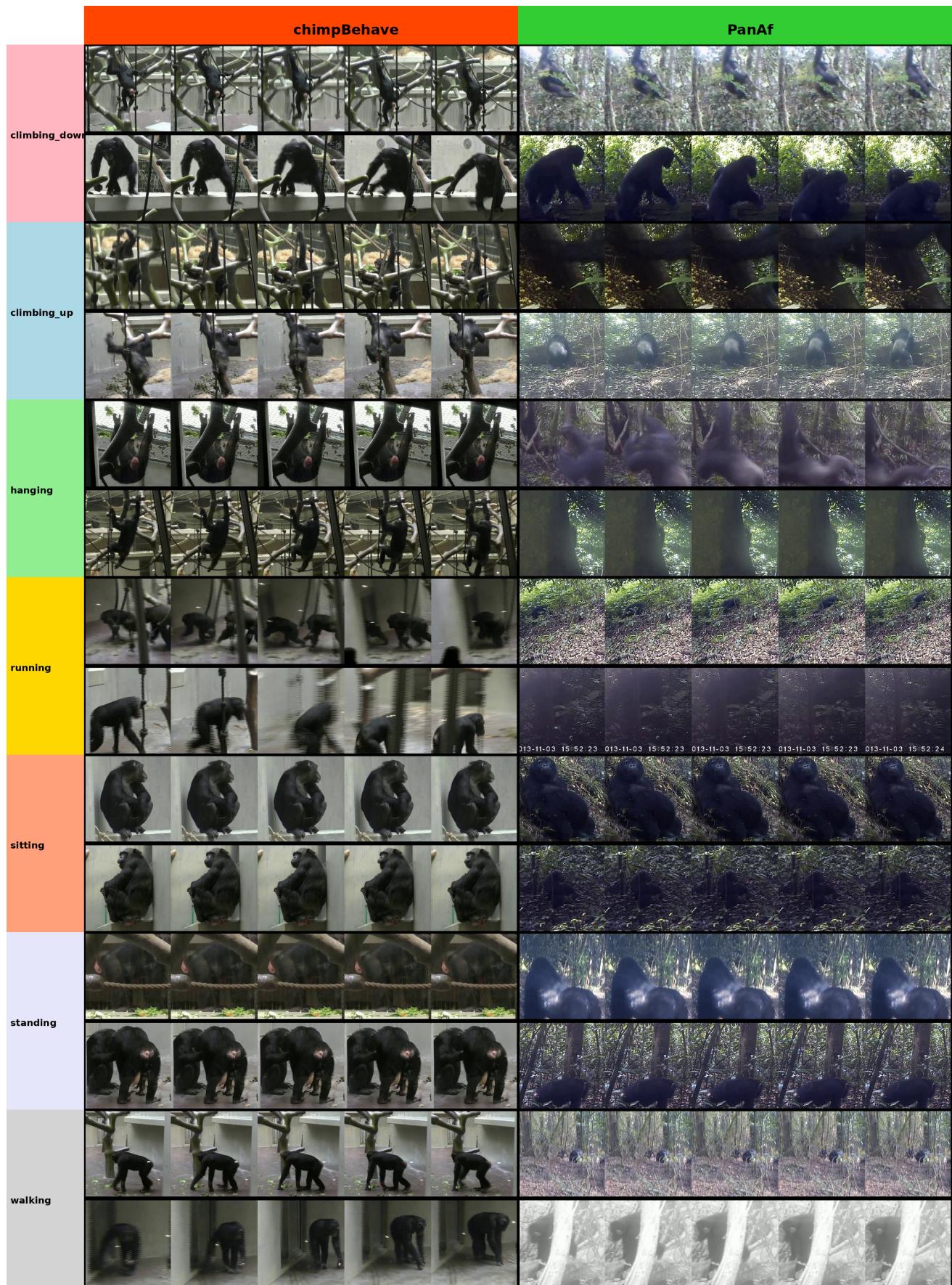


Figure 7. Examples of miniclip between datasets and behavior classes. Note that we sampled 4 out of 20 frames from each miniclip for visualization purposes.

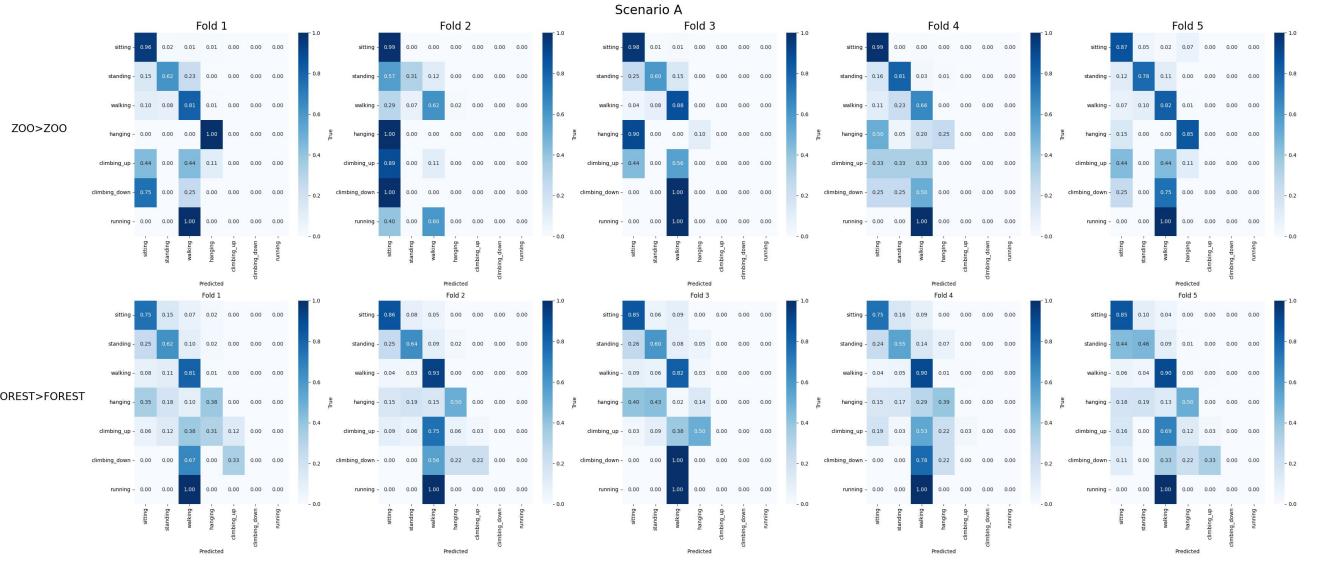


Figure 8. Confusion matrices for each fold in cross-validation: Scenario A - Within-dataset

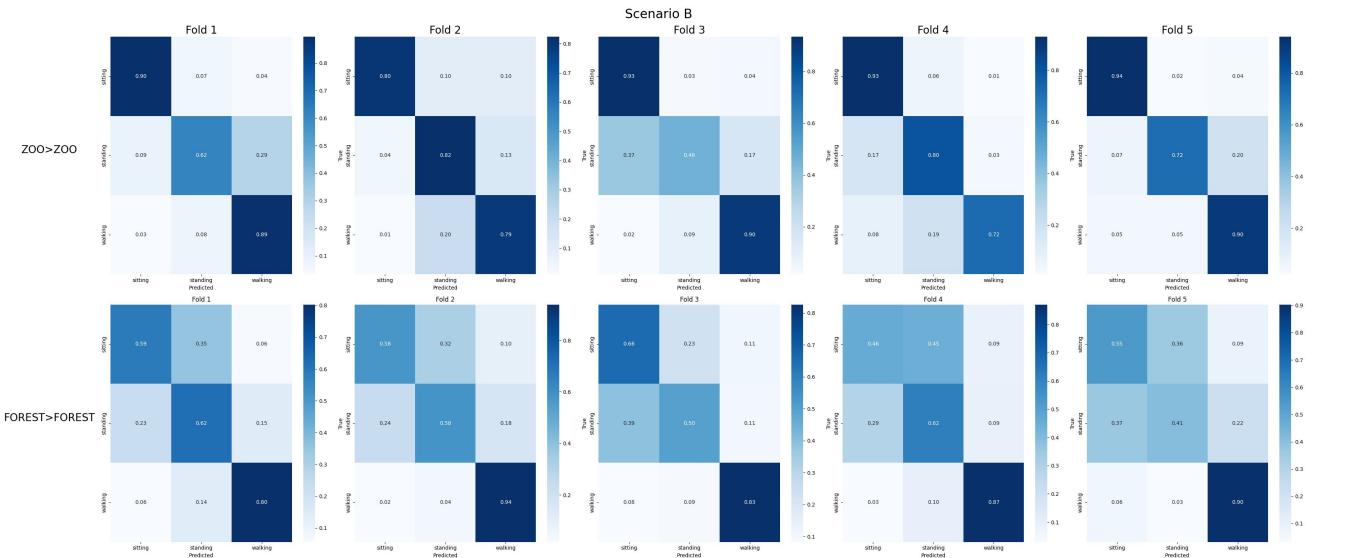


Figure 9. Confusion matrices for each fold in cross-validation: Scenario B - Within-dataset

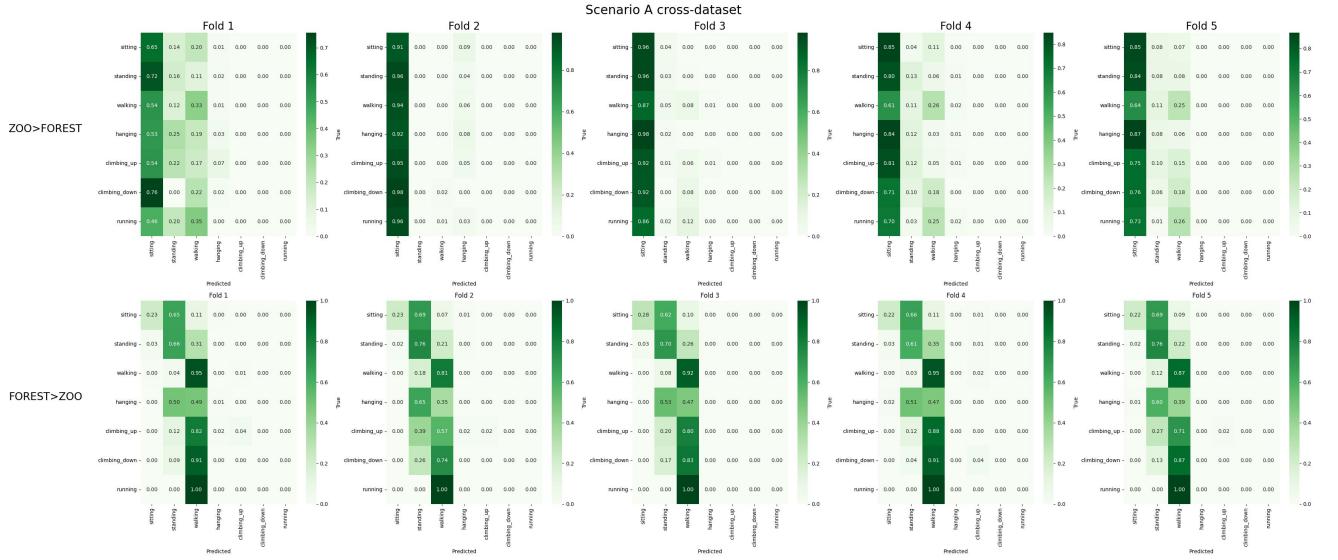


Figure 10. Confusion matrices for each fold in cross-validation: Scenario A - Cross-dataset

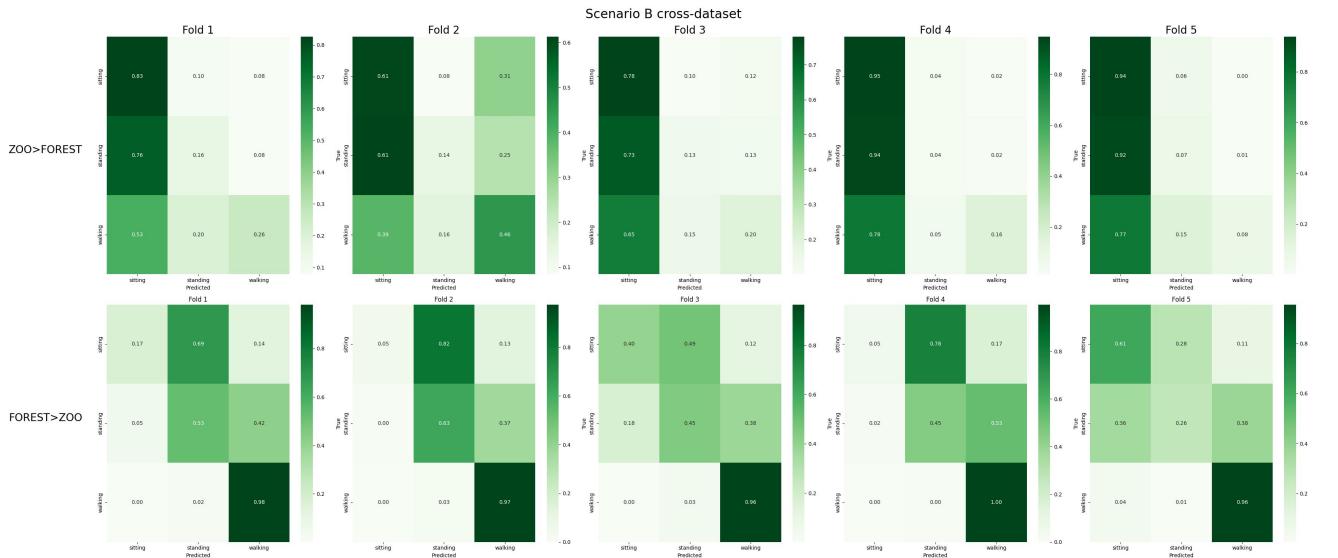


Figure 11. Confusion matrices for each fold in cross-validation: Scenario B - Cross-dataset