

Mouse Lockbox Dataset: Behavior Recognition for Mice Solving Lockboxes

Patrik Reiske ^{*† a b}Marcus N. Boon ^{† a b}Niek Andresen ^{a b c}Sole Traverso ^{a b}Katharina Hohlbaum ^{a d}Lars Lewejohann ^{a c d}Christa Thöne-Reineke ^{a c}Olaf Hellwich ^{‡ a b}Henning Sprekeler ^{‡ a b}

Abstract

Machine learning and computer vision methods have a major impact on the study of natural animal behavior, as they enable the (semi-)automatic analysis of vast amounts of video data. Mice are the standard mammalian model system in most research fields, but the datasets available today to refine such methods focus either on simple or social behaviors. In this work, we present a video dataset of individual mice solving complex mechanical puzzles, so-called lockboxes. The more than 110 hours of total playtime show their behavior recorded from three different perspectives. As a benchmark for frame-level action classification methods, we provide human-annotated labels for all videos of two different mice, that equal 13% of our dataset. Our keypoint (pose) tracking-based action classification framework illustrates the challenges of automated labeling of fine-grained behaviors, such as the manipulation of objects. We hope that our work will help accelerate the advancement of automated action and behavior classification in the computational neuroscience community. Our dataset is publicly available at <https://doi.org/10.14279/depositonce-23850>

1. Introduction

Ethology, the study of non-human behavior, [29] is one of the cornerstones of understanding complex biological systems. In recent years, with the integration of machine learning into the field, computational ethology [1] emerged as a powerful new paradigm offering new pathways for advancing both fields and beyond. For instance, it has significantly influenced neuroscience, enabling the development of com-

putational frameworks that bridge neural mechanisms with observations of behaviors [6, 14, 21, 30]. In robotics, animal behavior datasets allow researchers to learn artificial agents to navigate and interact autonomously in natural environments. The hypothesized learning models used in this process can then be tested by comparing the performance of the learned agents against that of natural agents [2].

The available datasets of freely moving animals [4, 5, 7–9, 13, 15, 16, 18, 19, 24–28, 32] provide the foundation for the development of automated behavioral analysis tools, e.g., [12, 17, 31]. However, all of these datasets and their descending methods focus on trivial and social behaviors, but neglect the structure imposed by well-defined tasks that provoke complex behaviors. This absence limits their applicability for studying goal-directed actions, problem-solving, and other behaviors critical to understanding cognitive processes in natural and artificial intelligence.

In this work, we provide the first large-scale labeled, single-agent, multi-perspective video dataset of mice showing complex behavior as they solve mechanical puzzles, so-called lockboxes. Every lockbox is baited with a food reward and consists of a single or a combination of four different mechanisms. As a benchmark, we provide labels for 13% of our data, including mechanism state, mouse-to-mechanism proximity, and both mouse-mechanism and mouse-reward actions. This amounts to about 15 hours and 25 minutes of total playtime. In doing so, we increase the longest total playtime, i.e., the number of perspectives multiplied by the real time recorded, available through any mouse dataset from 88 hours [5] by more than 33% to 117 hours and 52 minutes.

To provide high-quality label data, each labeled video is annotated by two skilled human raters who have been instructed prior to annotating. The consistency between raters is assessed by their inter-rater reliability [22], a well-established and objective measure of agreement. We regard such rigorous and transparent annotation protocols as essential for creating a dataset that allows us to reliably assess the performance of future machine learning methods.

We use our state-of-the-art keypoint-based method [3]

^{*}Corresponding author; patrik.reiske@tu-berlin.de

[†]Authors with equal contribution as first authors

[‡]Authors with equal contribution as last authors

^aCluster of Excellence “Science of Intelligence,” Berlin, Germany

^bTechnische Universität Berlin, Berlin, Germany

^cFreie Universität Berlin, Berlin, Germany

^dGerman Centre for the Protection of Laboratory Animals, German Federal Institute for Risk Assessment, Berlin, Germany

DATASET	AGENT	LABELS	DURATION
CRIM13	Mice	13 social	$2 \times 44\text{h} \approx 88\text{h}$
PAIR-R24M	Rats	14 social	$24 \times 9\text{h} \approx 220\text{h}$
MARS	Mice	3 social	$2 \times 14\text{h} \approx 28\text{h}$
CalMS21	Mice	3 social	$1 \times 70\text{h} \approx 70\text{h}$
Ours	Mice	20 task-specific	$3 \times 40\text{h} \approx 120\text{h}$

Table 1. Overview of video datasets showing rodents. Durations, i.e., the total playtime calculated as the number of perspectives multiplied by the real time recorded, are rounded. Our 20 behaviors reflect five labeled interactions on four lockbox mechanisms.

as an initial benchmark for our dataset, and compare our human-human agreement against its human-machine agreement. In the absence of established benchmark methods for the interaction of natural agents with their environment, this will allow others to assess the performance of their methods.

We hope that our dataset will serve two purposes. First, that it will promote the advancement and adoption of more diverse machine learning methods in computational ethology. Our dataset provides an interesting challenge to the machine learning community, since the classification of the labels we provide require both large-scale pose and fine-level visual information. And second, that analyses of our dataset by the community will advance our understanding of how natural agents learn to solve complex problems.

2. Related Work

We limit the following overview of available datasets to those that show rodents and provide behavior labels, as their largely similar visual appearance and motor apparatus has proven to allow for domain transfer settings [7]. Tab. 1 summarizes some of their distinguishing properties.

CRIM13 [5] has, to this date, been the largest mouse dataset with a total playtime of 88 hours, i.e., 44 hours of real time recorded from two (top-down and side) perspectives. It shows mice in a resident-intruder context, and provides 13 human-annotated (social) behavior labels—approach, attack, coitus, chase, circle, drink, eat, clean, human, sniff, up, walk, and other—for each of the 237 pairs of 10 minute long videos. There, human raters reach an agreement of 70% while the method proposed alongside the dataset reaches 61.2% human-machine agreement.

PAIR-R24M [19] is the longest rodent (rat) dataset with a total playtime of 220 hours, i.e., 9 hours of real-time recorded from 24 perspectives. It provides 14 human-annotated (social) behavior labels—amble, crouch, explore, head tilt, idle, investigate, locomotion, rear down, rear up,

small movement, sniff, groom, as well as close to, explore, and chase—for the entire dataset.

MARS [27] has a total playtime of 28 hours, i.e., 14 hours of real-time recorded from two (top-down and front) perspectives. It provides three human-annotated social behavior labels—attack, investigation, and mount—for then videos with a total playtime of 3 hours.

CalMS21 [28] is a 70 hour long mouse dataset recorded from a single (top-down) perspective. It provides three human-annotated social behavior labels—attack, investigate, and mount—for 10 hours worth of video data.

3. Dataset

This section describes our dataset in detail. Sec. 3.1 specifies the mouse breed, the arena including the lockboxes, the camera setup, the schedule at which mice were presented with the lockboxes, and the preprocessing of the recorded videos. Sec. 3.2 describes the annotation of behavior labels including our ethogram. Secs. 3.3.1 and 3.3.2 provide statistics on both playtimes and labels. And in Sec. 3.4 we report the limitations of our dataset.

3.1. Data Collection and Preprocessing

The video data was initially recorded for [11], where more detailed information on this section’s contents can be found. Twelve female C57BL/6J mice obtained from Charles River Laboratories (Sulzfeld, Germany) were recorded in a free-standing Makrolon type III cage, that was connected to a home cage of the same type by a tube. The mice were housed in groups of 4 animals in an artificial 12/12-hour light/dark cycle. The lockbox trials took place in the light phases, and only one animal could enter the arena at a time. The arena was closed with a cutout top grid to allow for unobstructed view on the lockbox. Three Basler acA1920-40um cameras (LM25HC7 lens, $f = 25\text{mm}$, $k = 1.4$; Kowa, Nagoya, Japan) were used to record grayscale videos at the maximum resolution of $1936 \times 1216\text{px}$ at the common [19, 27, 28] 30Hz frame rate. Additionally, two infrared lights (Synergy 21 IR-Strahler 60W, ALLNET GmbH Computersysteme, Germering, Germany) illuminated the cage. The advantages of infrared lights is that they enhance the video quality of used infrared-sensitive cameras while not being aversive to the animals.

Fig. 1a depicts the described setup. All cameras were connected to a computer and controlled by software to synchronize frame capturing. The mice were presented with five different lockboxes: a combined lockbox (Fig. 1b) consisting of four interlocked mechanisms, and four simpler lockboxes (Fig. 1c) presenting these mechanisms individually. A hidden food reward (oatmeal flake) was used as a

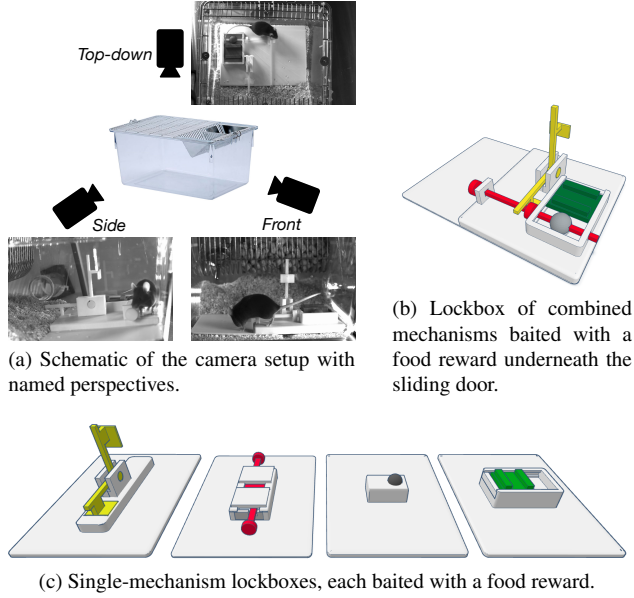


Figure 1. Camera setup used for recording the videos, as well as lockboxes and their mechanisms: lever (yellow), stick (red), ball (gray), and sliding door (green). Each lockbox is baited with a food reward underneath the (last) mechanism. Sec. 5 of our supplementary material provides figures of the unlocked lockboxes.

bait. The mice were not subjected to food or water deprivation, but had *ad libitum* access to food pellets (LASvendi, LAS QCDiet, Rod 16, autoclavable) as well as tap water. However, the food reward was exclusively provided within the lockboxes. To familiarize the mice with the food reward, they were habituated over three consecutive days prior to the lockbox trials by placing eight oat flakes at the location where the lockbox would be placed in the arena. The freely behaving mice were presented with the combined lockbox for a total of 6 and with the single-mechanism lockboxes 11 trials. The mice were first presented with the combined lockbox trial followed by 11 trials of a randomized order of each of the single-mechanism lockboxes, followed by another 5 combined lockbox trials. The videos end shortly after the reward is reached, or if a trial reached the maximum duration of 30 minutes for combined and 15 minutes for single-mechanism lockboxes.

We manually cut the videos to remove disturbances, e.g., the experimenter’s hands switching lockboxes. Any videos where the lockbox mechanisms were not captured in their entirety were filtered out.

3.2. Label Annotation

We provide human annotations for mechanism states, mouse-to-mechanism proximity, as well as both mouse-mechanism and mouse-reward action labels. To prevent information leakage between labeled and unlabeled data splits, we labeled all videos of two specific mice (mouse

LABEL	DEFINITION
Proximity	The mouse’s snout is within a distance of 1cm to a specific mechanism.
Touch	The mouse touches a specific mechanism with one or both of its front paws.
Bite	The mouse bites into a specific mechanism.
Unlock	The state of a specific mechanism changes to unlocked. This may make the reward accessible or enabling the next mechanism to be unlocked. State changes may occur without the mouse manipulating a mechanism directly.
Lock	The state of a specific mechanism changes to locked. This may make the reward inaccessible or preventing the next mechanism from being unlocked. State changes may occur without the direct manipulation of a mechanism.
Reach reward	The mouse is in first contact with the reward with any of its body parts.

Table 2. Ethogram used for label annotation. Sec. 6 of our supplementary material provides example frames for the different labels.

numbers 291 and 324) that have a combined total playtime of 15 hours and 25 minutes in 270 videos spanning 90 trials. This equals about 13% of our dataset.

Tab. 2 defines the ethogram used by our nine skilled human raters. We used these labels as they express trivial truths in order to minimize anthropomorphic biases, that would otherwise distort both experiment evaluations and drawn conclusions. These biases are especially apparent when using more high-level labels, such as exploring and deliberately manipulating lockbox mechanisms, that strongly depend on subjective human interpretation. Using more explicit labels not only leads to higher label quality but also lowers the risk of computer vision and machine learning models learning said biases before reintroducing them as noise to any analysis based on their outputs.

For annotating the labels, we merged every video triplet (top-down, side, and front perspective) into a combined video.¹ All labels have been annotated by randomized pairs of raters with a temporal accuracy of ± 100 milliseconds, i.e., ± 3 frames, using BORIS [10]. To annotate any of the videos, it took our raters about 6.2 to 11.5 times longer than their playtime. This aligns with the factor of about 5 to 10 that is usually reported throughout the available literature.

¹Merging the video triplets into combined videos was necessary as BORIS version 8.27 suffers from a software issue that occurs more frequently when using it with multiple videos opened at once, and that causes the software to crash only minutes into using it. The published dataset does not include the merged videos.

	Lever	Stick	Ball	Door	Any
Proximity	15.73	19.05	13.41	18.97	55.39
Touch	7.06	4.07	7.00	9.32	25.50
Bite	1.81	1.50	3.41	1.42	8.12

Table 3. Action labels in the labeled videos in percent.

We account our slightly higher efforts to the multitude of mouse body parts and lockbox mechanisms that needed to be observed at the same time.

3.3. Dataset Statistics

This section gives an overview over various data and label statistics. Sec. 3.3.1 provides insight on playtime statistics, and Sec. 3.3.2 on label statistics. We provide further statistics and analyses in our previous work [3, 11].

3.3.1. Playtime Statistics

Our dataset was recorded from 3 perspectives (top-down, side, and front) and has a total playtime of 117 hours and 52 minutes, i.e., 39 hours and 17 minutes of recorded real time. It consists of a total of 1629 videos, i.e., 543 trials, with a mean playtime of 4 minutes and 21 seconds. Any combination of individual mice and lockboxes accounts for 0.3–7.6% of the data. Each mouse accounts for 5.3–15.3% of the data, while the single-mechanisms lockboxes account for 9.7–14.2% and the combined lockbox for 52% of the data.

3.3.2. Label Statistics

We provide human-annotated labels for a total playtime of 15 hours and 25 minutes, i.e., 5 hours and 8 minutes of recorded real time. Sec. 3.2 provides details on both the labels and their annotation. Tab. 3 shows the resulting label distribution. The uneven label distribution is rooted in the mice behaving freely, and reflects their naturally occurring preference for different actions and mechanisms.

Fig. 2 shows the inter-rater reliability, i.e., Cohen’s kappa coefficients, [22] among human raters. On average, our human raters annotate most proximity and touch labels with moderate to strong agreement, except for the stick mechanism. However, they annotate bite labels only with minimal to weak agreement. We account this to biting being particularly hard to annotate as it rarely is directly visible in the videos. Fig. 2 also shows the inter-rater reliability between the human raters and our benchmark method [3], on which we provide details in Sec. 7 of our supplementary material. In short, the benchmark method uses 2-dimensional keypoints, extracted with DeepLabCut [20, 23], reconstructs the scene in 3D, and refines the tracks using (extended) Kalman filtering. These refined 3-dimensional tracks are then used to extract frame-per-frame

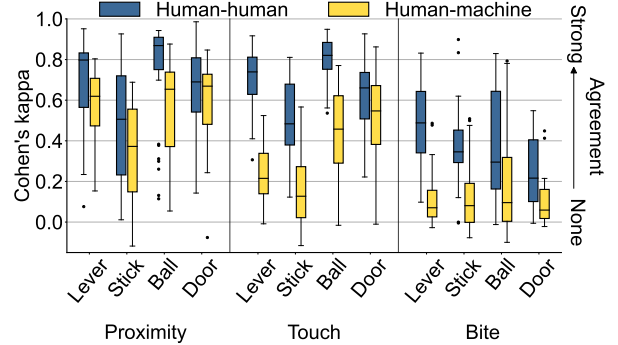


Figure 2. Human-human and human-machine inter-rater reliability, i.e., Cohen’s kappa coefficients, for different action labels.

action labels based on bounding boxes around the mechanisms of interest. While it almost reaches human-human reliability for proximity labels, its reliability is outperformed by humans for touch and bite labels. We account the decreased performance in touch and bite labels to the necessity for increased spatial accuracy in tracking the mechanisms and mouse body parts.

In addition to the inter-rater reliability, we report the F_1 scores for all action labels (Tab. 4) where we use tolerance of ± 3 frames that matches the temporal accuracy our human annotators were instructed to adhere to. As we have mentioned before, annotating actions is a nontrivial task even for humans. Therefore, we report on both the F_1 scores comparing humans against each other, and individual humans against our pipeline, as well as the union of human-annotated labels against our pipeline. We consider the inter-human F_1 scores to be the upper performance limit for any benchmark method.

3.4. Limitations

Our dataset has three limitations. First, since the video recording was pseudo-synchronized through software, the frames of different cameras have been captured with a temporal desynchronization. We sampled the average asynchronicity to be 1.39 frames with a standard deviation of 1.50 frames. Since this is lower than the accuracy we annotated our labels with, we do not expect it to cause any major issues. Second, not all videos share the same exact positioning of the cameras as they have been recorded in several trial periods over the course of months, during which the setup had to be rearranged. And third, due to insufficient lighting conditions and severe camera dislocation, some trials had to be discarded from the dataset resulting in an imbalanced number of videos per mouse.

4. Conclusion

In this work, we presented the—to the best of our knowledge—first available single-agent, multi-perspective

	Proximity				Touch				Bite			
	Lever	Stick	Ball	Door	Lever	Stick	Ball	Door	Lever	Stick	Ball	Door
Human 1 × human 2	0.85	0.70	0.88	0.83	0.76	0.63	0.88	0.72	0.62	0.32	0.58	0.17
Human 1 × machine	0.76	0.69	0.73	0.77	0.38	0.25	0.60	0.60	0.15	0.13	0.23	0.09
Human 2 × machine	0.74	0.54	0.73	0.77	0.37	0.28	0.60	0.65	0.12	0.14	0.24	0.13
∪ humans × machine	0.77	0.64	0.74	0.80	0.40	0.29	0.63	0.67	0.19	0.19	0.30	0.18

Table 4. F_1 scores (± 3 frames tolerance) comparing human annotators and our benchmark method. Our benchmark method is compared against both individual humans and the union of their annotated labels. The scores comparing humans (top row) can be used as a reference.

video dataset of mice showing complex behavior as they learn to solve mechanical puzzle mechanisms. These so-called lockboxes consist of either one of four mechanisms or their combination, and are baited with a food reward. In total, we provide videos with a total playtime, i.e., the number of perspectives multiplied by the real time recorded, of 117 hours and 52 minutes.

As a benchmark, we provide human-annotated behavior labels for 13% of our dataset, and report the inter-rater reliability among humans as well as between humans and our benchmark method. We find that the benchmark method almost reaches human-level performance for proximity labels, but not for touching and biting labels.

We hope that our dataset will contribute to this advancement by challenging and inspiring others. Our dataset is publicly available at DOI <https://doi.org/10.14279/depositonce-23850>

Acknowledgments

We thank our encouraged lab assistants Clara Bekemeier, Sophia Meier, Jule Detmers, and Andreas Pauli for their support with cleaning the raw video data and annotating the labels. Their dedication and hard work were essential to composing the presented dataset.

This project was funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany’s Excellence Strategy—EXC 2002/1 “Science of Intelligence”—project number 390523135.

Ethics Statement

This research does not involve human subjects, sensitive data, harmful insights, nor methodologies or applications that may raise ethical concerns.

As reported in [11], animal research was conducted in compliance with the German Animal Welfare Act and Directive 2010/63/EU on the protection of animals used for scientific purposes. The experimental procedures and maintenance of the animals were preregistered in the Animal Study Registry (DOI [10.17590/asr.0000237](https://doi.org/10.17590/asr.0000237)) and approved by the Berlin State Authority, Landesamt für Gesundheit und Soziales (permit number G0249/19).

The authors declare that they have no conflicts of interest. No sponsorships influenced this research.

References

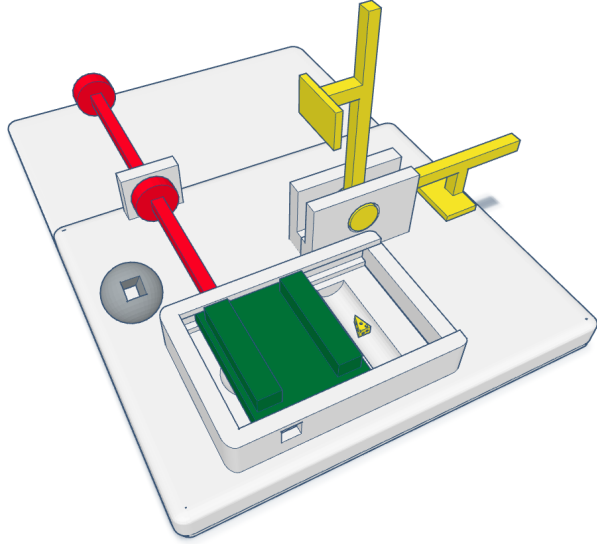
- [1] David J. Anderson and Pietro Perona. Toward a Science of Computational Ethology. *Neuron*, 84(1):18–31, 2014. 1
- [2] Manuel Baum, Lukas Schattenhofer, Theresa Rössler, Antonio Osuna-Mascaró, Alice Auersperg, Alex Kacelnik, and Oliver Brock. Yoking-Based Identification of Learning Behavior in Artificial and Biological Agents. In *From Animals to Animats*, pages 67–78, 2022. 1
- [3] Marcus N. Boon, Niek Andresen, Soledad Traverso, Sophia Meier, Friedrich Schuessler, Olaf Hellwich, Lars Lewejohann, Christa Thöne-Reineke, Henning Sprekeler, and Katharina Hohlbaum. Mechanical problem solving in mice. *bioRxiv*, 2024. 1, 4
- [4] Otto Brookes, Majid Mirmehdi, Colleen Stephens, Samuel Angedakin, Katherine Corogenes, Dervla Dowd, Paula Dieguez, Thurston C. Hicks, Sorrel Jones, Kevin Lee, Vera Leinert, Juan Lapuente, Maureen S. McCarthy, Amelia Meier, Mizuki Murai, Emmanuelle Normand, Virginie Vergnes, Erin G. Wessling, Roman M. Wittig, Kevin Langergraber, Nuria Maldonado, Xinyu Yang, Klaus Zuberbühler, Christophe Boesch, Mimi Arandjelovic, Hjalmar Kühl, and Tilo Burghardt. PanAf20K: A Large Video Dataset for Wild Ape Detection and Behaviour Recognition. *International Journal of Computer Vision*, 132(8):3086–3102, 2024. 1
- [5] Xavier P. Burgos-Artizzu, Piotr Dollár, Dayu Lin, David J. Anderson, and Pietro Perona. Social behavior recognition in continuous video. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1322–1329, 2012. 1, 2
- [6] Sandeep Robert Datta, David J. Anderson, Kristin Branson, Pietro Perona, and Andrew Leifer. Computational Neuroethology: A Call to Action. *Neuron*, 104(1):11–24, 2019. 1
- [7] Timothy W. Dunn, Jesse D. Marshall, Kyle S. Severson, Diego E. Aldarondo, David G. C. Hildebrand, Selmaan N. Chettih, William L. Wang, Amanda J. Gellis, David E. Carlson, Dmitriy Aronov, Winrich A. Freiwald, Fan Wang, and Bence P. Ölveczky. Geometric deep learning enables 3D kinematic profiling across species and environments. *Nature Methods*, 18:564–573, 2021. 1, 2
- [8] Isla Duporge, Maksim Kholiavchenko, Roi Harel, Scott Wolf, Dan Rubenstein, Meg Crofoot, Tanya Berger-Wolf,

- Stephen Lee, Julie Barreau, Jenna Kline, Michelle Ramirez, and Charles Stewart. BaboonLand Dataset: Tracking Primates in the Wild and Automating Behaviour Recognition from Drone Videos, 2024.
- [9] Eyrun Eyjolfssdottir, Steve Branson, Xavier P. Burgos-Artizzu, Eric D. Hoopfer, Jonathan Schor, David J. Anderson, and Pietro Perona. Fly v. Fly Dataset, 2021. doi: 10.22002/D1.1893. 1
- [10] Olivier Friard and Marco Gamba. BORIS: a free, versatile open-source event-logging software for video/audio coding and live observations. *Methods in Ecology and Evolution*, 7: 1325001330, 2016. 3
- [11] Katharina Hohlbaum, Niek Andresen, Paul Mieske, Pia Kahnau, Benjamin Lang, Kai Diedrich, Rupert Palme, Lars Mundhenk, Henning Sprekeler, Olaf Hellwich, Christa Thöne-Reineke, and Lars Lewejohann. Lockbox enrichment facilitates manipulative and cognitive activities for mice. *Open Research Europe*, 4(108), 2024. 2, 4, 5
- [12] Alexander I. Hsu and Eric A. Yttri. B-SOiD, an open-source unsupervised algorithm for identification and fast prediction of behaviors. *Nature Communications*, 12:5188, 2021. 1
- [13] Bo Hu, Bryan Seybold, Shan Yang, Avneesh Sud, Karla Barron, Yi Liu, Pauly Cha, Marcelo Cosino, Ellie Karlsson, Janessa Kite, Ganesh Kolumam, Joseph Preciado, Chunlian Solorio, José Zavala-and Zhang, Xiaomeng Zhang, Martin Voorbach, Ann E. Tovcimak, J. Graham Ruby, and David A. Ross. 3D mouse pose from single-view video and a new dataset. *Scientific Reports*, 13:13554, 2023. 1
- [14] Ann Kennedy. The what, how, and why of naturalistic behavior. *Current Opinion in Neurobiology*, 74:102549, 2022. 1
- [15] Maksim Kholiavchenko, Jenna Kline, Michelle Ramirez, Sam Stevens, Alec Sheets, Reshma Babu, Namrata Banerji, Elizabeth Campolongo, Matthew Thompson, Nina Van Tiel, Jackson Miliko, Eduardo Bessa, Isla Duporge, Tanya Berger-Wolf, Daniel Rubenstein, and Charles Stewart. KABR: In-Situ Dataset for Kenyan Animal Behavior Recognition from Drone Videos. In *2024 IEEE/CVF Winter Conference on Applications of Computer Vision Workshops (WACVW)*, pages 31–40, 2024. 1
- [16] Ci Li, Ylva Mellbin, Johanna Krogager, Senya Polikovsky, Martin Holmberg, Nima Ghorbani, Michael J. Black, Hedvig Kjellström, Silvia Zuffi, and Elin Hernlund. The Poses for Equine Research Dataset (PFERD). *Science Data*, 11:497, 2024. 1
- [17] Kevin Luxem, Petra Mocellin, Falko Fuhrmann, Johannes Kürsch, Stephanie R. Miller, Jorge J. Palop, Stefan Remy, and Pavol Bauer. Identifying behavioral structure from deep variational embeddings of animal motion. *Communications Biology*, 5:1267, 2022. 1
- [18] 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*, 2023. 1
- [19] Jesse D. Marshall, Ugne Klibaite, Amanda Gellis, Diego E. Aldarondo, Bence P. Ölveczky, and Timothy W. Dunn. The PAIR-R24M Dataset for Multi-animal 3D Pose Estimation. In *Advances in Neural Information Processing Systems*, 2021. 1, 2
- [20] 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. *Nature Neuroscience*, 21:1281–1289, 2018. 4, 1
- [21] Michael H. McCullough and Geoffrey J. Goodhill. Unsupervised quantification of naturalistic animal behaviors for gaining insight into the brain. *Current Opinion in Neurobiology*, 70:89–100, 2021. 1
- [22] Mary L. McHugh. Interrater reliability: the kappa statistic. *Biochemia Medica*, 22(3):276–282, 2012. 1, 4
- [23] Tanmay Nath, Alexander Mathis, An Chi Chen, Amir Patel, Matthias Bethge, and Mackenzie W Mathis. Using deeplabcut for 3d markerless pose estimation across species and behaviors. *Nature Protocols*, 2019. 4, 1
- [24] 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 *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 19001–19012, 2022. 1
- [25] Malte Pedersen, Joakim Bruslund Haurum, Stefan Hein Bengtson, and Thomas B. Moeslund. 3D-ZeF: A 3D Zebrafish Tracking Benchmark Dataset. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2423–2433, 2020.
- [26] Mitchell Rogers, Gaël Gendron, David Arturo Soriano Valdez, Mihailo Azhar, Yang Chen, Shahrokh Heidari, Caleb Perelini, Padriac O’Leary, Kobe Knowles, Izak Tait, Simon Eyre, Michael Witbrock, and Patrice Delmas. Meerkat Behaviour Recognition Dataset, 2023.
- [27] Cristina Segalin, Jalani Williams, Tomomi Karigo, May Hui, Moriel Zelikowsky, Jennifer J. Sun, Pietro Perona, David J. Anderson, and Ann Kennedy. The Mouse Action Recognition System (MARS) software pipeline for automated analysis of social behaviors in mice. *eLife*, 10:e63720, 2021. 2
- [28] Jennifer J. Sun, Tomomi Karigo, Dipam Chakraborty, Sharada P. Mohanty, Benjamin Wild, Quan Sun, Chen Chen, David J. Anderson, Pietro Perona, Yisong Yue, and Ann Kennedy. The Multi-Agent Behavior Dataset: Mouse Dyadic Social Interactions. In *Advances in Neural Information Processing Systems*, 2021. 1, 2
- [29] Nikolaas Tinbergen. On aims and methods of ethology. *Zeitschrift für Tierpsychologie*, 20(1):410–433, 1961. 1
- [30] Lukas von Ziegler, Oliver Sturman, and Johannes Bohacek. Big behavior: challenges and opportunities in a new era of deep behavior profiling. *Neuropsychopharmacology*, 46:33–44, 2021. 1
- [31] Caleb Weinreb, Jonah E. Pearl, Sherry Lin, Mohammed Abdal Monium Osman, Libby Zhang, Sidharth Annapragada, Eli Conlin, Red Hoffmann, Sofia Makowska, Winthrop F. Gillis, Maya Jay, Shaokai Ye, Alexander Mathis, Mackenzie W. Mathis, Talmo Pereira, Scott W. Linderman, and Sandeep Robert Datta. Keypoint-MoSeq: parsing behavior by linking point tracking to pose dynamics. *Nature Methods*, 21:1329–1339, 2024. 1

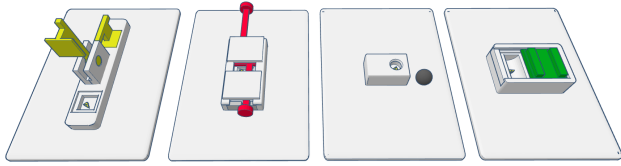
- [32] Ali Zia, Renuka Sharma, Reza Arablouei, Greg Bishop-Hurley, Jody McNally, Neil Bagnall, Vivien Rolland, Brano Kusy, Lars Petersson, and Aaron Ingham. CVB: A Video Dataset of Cattle Visual Behaviors, 2023. [1](#)

Mouse Lockbox Dataset: Behavior Recognition for Mice Solving Lockboxes

Supplementary Material



(a) Unlocked lockbox of combined mechanisms baited with a symbolized food reward underneath the sliding door.



(b) Unlocked single-mechanism lockboxes baited with a symbolized food reward underneath each mechanisms.

Figure 3. Unlocked lockboxes and their mechanisms: lever (yellow), stick (red), ball (gray), and sliding door (green). This depiction contains symbolized food baits.

5. Lockboxes with Unlocked Mechanisms

Fig. 3 shows the opened lockboxes with symbolized food baits; see Figs. 1b and 1c for reference.

6. Example Frames for Labels

Figure 4 shows a selection of examples for our different label classes.

7. Benchmark Method

Our benchmark experiments are based on the pose-tracking approach used by Boon et al. [3]. The method consists of three steps: the use of DeepLabCut (DLC) for 2-dimensional pose tracking, 3-dimensional reconstruction and the refinement of keypoint data using (Extended)

Kalman filtering, and the detection of action labels. A high-level description of steps is given below.

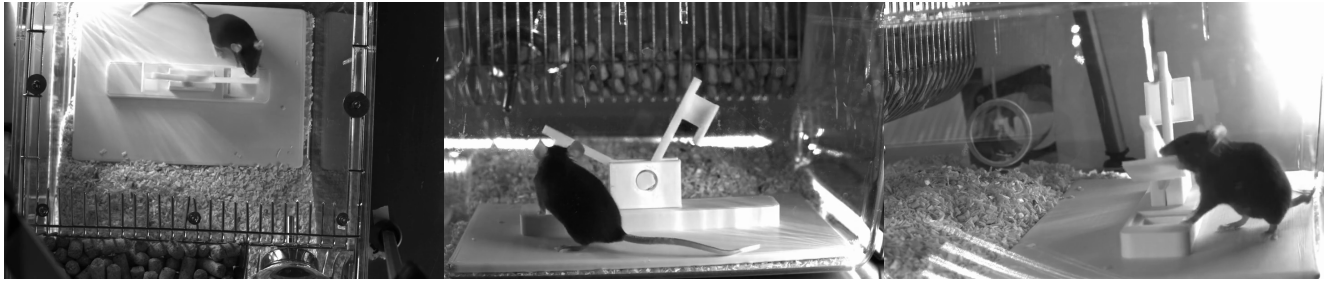
First, 2-dimensional poses of the mice and lockbox mechanisms are extracted from the videos on a frame-level by learning DLC models under supervision. We learn one DLC model to locate keypoints of mice, and two that locate keypoints of lockbox mechanisms—one for the single-mechanism lockboxes, and one for the lockbox combining them—using default parameters [20, 23]. Next, the scene is reconstructed by utilizing the known 3-dimensional locations of the lockbox mechanisms given by their CAD models. We linearly map the known 3-dimensional locations onto the corresponding triplets of 2-dimensional keypoints using the RANSAC algorithm and construct a triangulation matrix for each trial. For trials in which the lockbox does not have a well-defined third dimension (i.e., single stick, ball, and door), the lever trial of that mouse (and day) provides reference triangulation data instead. Potential rotations due to the lockbox not being placed in the same orientation as the lever reference are accounted for by rotating the 3-dimensional reference coordinates from the CAD model in accordance to the observed rotation in the xy-plane from the top-down view camera.

The triangulation matrices for each trial are used as observation matrices for (Extended) Kalman filters to refine the observed triplets of 2-dimensional keypoints into a common 3-dimensional space. The head and the tail of the mouse are inferred using a skeletal model, while the keypoints of the mechanisms and the paws of the mouse are inferred as single keypoints.

Finally, the interactions of the mice with the lockbox mechanisms are detected based on the 3-dimensional poses of the mouse and predefined bounding boxes spanned by the 3-dimensional keypoint locations. For the proximity labels, the snout of mouse is used to detect the actions: each frame in which the snout of the mouse is inside of a bounding box defined around each lockbox mechanism, the corresponding action label (e.g., proximity lever) is detected. Biting is detected using the mouth of the mouse, which location is computed from the rigid body model of the mouse head. The touch labels are detected using the locations of the front paws. Note that the bite and touch labels have different predefined bounding boxes than the proximity labels, as these actions have a finer level of granularity than proximity labels. As a last processing step, each action label is filtered by a weighted moving average filter with a Gaussian shape (with a window size of 30 frames), to remove faulty single-frame detections from the data.



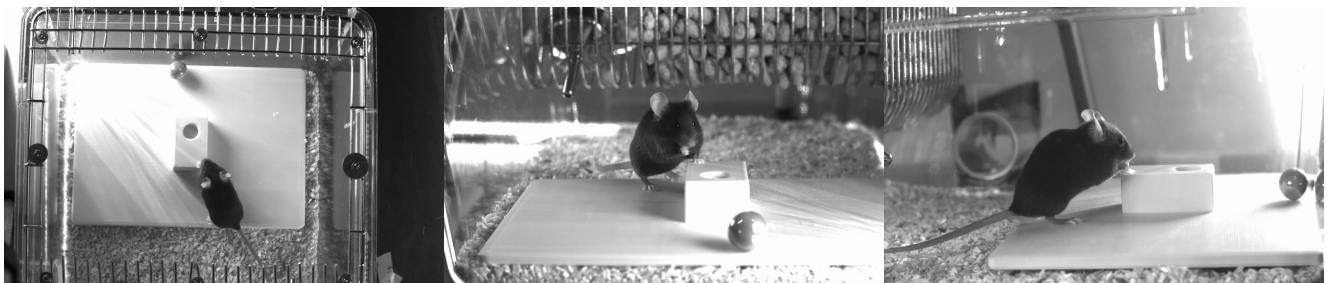
(a) Frame example with mouse in proximity to lever and touching the sliding door while all mechanisms are locked.



(b) Frame example with mouse in proximity to and biting the lever while the mechanism is unlocked.



(c) Frame example with mouse in proximity to the stick while the mechanism is locked.



(d) Frame example with no action label active while the ball mechanism is unlocked.



(e) Frame example with mouse in proximity to the sliding door while the mechanism is unlocked.

Figure 4. Example frames from labeled videos showing mice performing different actions.