

Deep-learning based identification, tracking, pose estimation, and behavior classification of interacting primates and mice in complex environments

Markus Marks^{1,3}, Jin Qiuhan⁴, Oliver Sturman^{2,3}, Lukas von Ziegler^{2,3}, Sepp Kollmorgen^{1,3},
Wolfgar von der Behrens^{1,3}, Valerio Mante^{1,3}, Johannes Bohacek^{2,3}, Mehmet Fath Yank^{1,3,*}

¹Institute of Neuroinformatics ETH Zürich and University of Zürich, Switzerland

²Laboratory of Molecular and Behavioral Neuroscience, Institute for Neuroscience,
Department of Health Sciences and Technology, ETH Zurich, Switzerland

³Neuroscience Center Zurich, ETH Zürich and University of Zürich, Switzerland,

⁴Laboratory for Neuro- & Psychophysiology, Department of Neurosciences, KU Leuven,
Belgium

*correspondence to: yanik@ethz.ch

Abstract

The quantification of behaviors of interest from video data is commonly used to study brain function, the effects of pharmacological interventions, and genetic alterations. Existing approaches lack the capability to analyze the behavior of groups of animals in complex environments. We present a novel deep learning architecture for classifying individual and social animal behavior, even in complex environments directly from raw video frames, while requiring no intervention after initial human supervision. Our behavioral classifier is embedded in a pipeline (SIPEC) that performs segmentation, identification, pose-estimation, and classification of complex behavior, outperforming the state of the art. SIPEC successfully recognizes multiple behaviors of freely moving individual mice as well as socially interacting non-human primates in 3D, using data only from simple mono-vision cameras in home-cage setups.

Introduction

While the analysis of animal behavior is crucial for systems neuroscience¹ and preclinical assessment of therapies, it remains a highly laborious and error-prone process. Over the last few years, there has been a surge in machine learning tools for behavioral analysis, including segmentation, identification, and pose estimation^{2–11}. Although this has been an impressive feat for the field, a key element, the direct recognition of behavior itself, has been rarely addressed. Unsupervised analysis of behavior^{12–17} can be a powerful tool to capture the diversity of the underlying behavioral patterns, but the results of these methods do not align with human annotations and therefore require subsequent inspection¹⁵. There have been advances also in the supervised analysis of mouse behavior, using classifiers on top of pose-estimation generated features^{18–21} or manually defined features such as ellipses^{22–25}. Sturman et. al.²⁰ demonstrated that the classification of mouse behaviors using features generated from pose-estimation

39 algorithms can outperform the behavioral classification performance of commercial systems.
40 Yet, such pose-estimation-based behavior classification remains a labor-intensive and error-
41 prone process as we show below. Moreover, pose estimation in primates is difficult to achieve
42 with current methods²⁶.

43
44 Here, we demonstrate a complementary approach for researchers who automatically seek to
45 identify behaviors of interest. Our approach relies on the initial annotation of exemplar
46 behaviors, i.e. snippets of video footage. These video snippets are subsequently used to train a
47 Deep Neural Network (DNN) to subsequently recognize such particular behaviors in arbitrarily
48 long videos and complex environments. To achieve this, we designed a novel DNN
49 architecture, called SIPEC:BehaveNet, which uses raw videoframes as input and significantly
50 outperforms a pose-estimation-based approach tested on a well-annotated mouse dataset and
51 reaches human-level performances for counting grouped behavioral events. In addition to this
52 behavioral classification network, we developed the first all-inclusive pipeline, called SIPEC,
53 with modules for segmentation (SIPEC:SegNet), identification (SIPEC:IdNet), behavioral
54 classification (SIPEC:BehaveNet), and pose estimation (SIPEC:PoseNet) of multiple and
55 interacting animals in complex environments. This pipeline utilizes four DNNs operating
56 directly on videos, developed and optimized for analyzing animal behavior and providing state-
57 of-the-art performance. We use this pipeline to classify, for the first time, social interactions in
58 home-caged primates from raw video frames and without needing to use any pose estimation.

59
60 SIPEC:SegNet is a Mask R-CNN architecture²⁷, optimized to robustly segment animals despite
61 occlusions, multiple scales, and rapid movement, and enables tracking of animal identities
62 within a session. SIPEC:IdNet has a DenseNet²⁸ backbone, that yields visual features, that are
63 integrated over time through a gated-recurrent-unit network (GRU)^{29,30} to re-identify animals
64 when temporal-continuity-based tracking does not work, for example when animals enter or
65 exit a scene. This enables SIPEC to identify primates across weeks and to outperform the
66 identification module of idtracker.ai⁴ both within-session and across sessions (see also
67 Discussion) as well as primnet³¹. SIPEC:PoseNet performs top-down multi-animal pose
68 estimation which we compared to DeepLabCut (DLC)². SIPEC:BehaveNet uses an Xception³²
69 network in combination with a temporal convolution network (TCN)^{33,34} to classify behavioral
70 events directly from raw pixels. To rapidly train our modules, we use image augmentation³⁵ as
71 well as transfer-learning³⁶, optimized specifically for each task. SIPEC enables researchers to
72 identify behaviors of multiple animals in complex and changing environments over multiple
73 days or weeks in 3D space, even from a single camera with relatively little labeling, in contrast
74 to other approaches that use heavily equipped environments and large amounts of labelled data⁸.

75
76 To accelerate the reusability of SIPEC, we share the network weights among all four modules
77 for mice and primates, which can be directly used for analyzing new animals in similar
78 environments without further training or serve as pre-trained networks to accelerate training of
79 networks in different environments.

80

81 Results

82 Our algorithm performs segmentation (SIPEC:SegNet) followed by identification
83 (SIPEC:IdNet), behavioral classification (SIPEC:BehaveNet) and finally pose estimation
84 (SIPEC:PoseNet) from video frames (Fig. 1). These four artificial neural networks, trained for
85 different purposes, can also be used individually or combined in different ways (Fig. 1a). To
86 illustrate the utility of this feature, Fig. 1b shows the output of pipelining SIPEC:SegNet and
87 SIPEC:IdNet to track the identity and location of 4 primates housed together (Fig. 1b, Supp.
88 Video 1). Fig. 1c shows the output of pipelining SIPEC:SegNet and SIPEC:PoseNet to do
89 multi-animal pose estimation in a group of 4 mice.

90 **Segmentation module SIPEC:SegNet.** SIPEC:SegNet (see Methods, Supp. Fig. 12) is based
91 on the Mask-RCNN architecture²⁷, which we optimized for analyzing multiple animals and
92 integrated into SIPEC. We further applied transfer learning³⁶ onto the weights of the Mask-
93 RCNN ResNet-backbone³⁷ pre-trained on the Microsoft Common Objects in Context (COCO
94 dataset)³⁸ (see Methods for SIPEC:SegNet architecture and training). Moreover, we applied
95 image augmentation³⁵ to increase network robustness against invariances, e.g. rotational
96 invariance and therefore increase generalizability.

97 *Segmentation performance on individual mice and groups of 4.* We first examined the
98 performance of SIPEC:SegNet on top-view video recordings of individual mice, behaving in
100 an open-field test (OFT). While segmenting black mice on a blank background could be
101 achieved by thresholding alone, we still included this task for completeness. 8 mice were freely
102 behaving for 10 minutes in the TSE Multi Conditioning System's OFT arena, previously
103 described in Sturman et al.²⁰. We labeled the outlines of mice in a total of 23 frames using the
104 VGG image annotator³⁹ from videos of randomly selected mice. To evaluate the performance,
105 we used 5-fold cross-validation (CV). We assessed the segmentation performance on images
106 of individual mice, where SIPEC:SegNet achieved a mean-Average Precision (mAP) of 1.0 ± 0 (mean \pm s.e.m., see Methods for metric details). We performed a videoframe ablation study
107 to determine how many labeled frames (outline of the animal, see Supp. Fig. 1) are needed for
108 SIPEC:SegNet to reach peak performance (Supp. Fig. 2). While randomly selecting an
109 increasing amount of training frames, we measured performance using CV. For single-mouse
110 videos, we find that our model achieves 95% of mean peak performance (mAP of 0.95 ± 0.05)
111 using as few as a total of 3 labeled frames for training. To the existing 23 labeled single-mouse
112 frames, we added 57 labeled 4-plex frames, adding to a total of 80 labeled frames. Evaluated
113 on a 5-fold CV, SIPEC:SegNet achieves an mAP of 0.97 ± 0.03 (Fig. 2b). For segmentation in
114 groups of 4 mice, we performed an ablation study as well and found that SIPEC:SegNet
115 achieves better than 95% of the mean peak performance (mAP of 0.94 ± 0.05) using as few as
116 only 16 labeled frames. To assess the overlap between prediction and ground truth, we report
117 IoU and dice coefficient metrics as well (Fig. 2b).

118 *Segmentation performance of groups of primates.* To test SIPEC:SegNet for detecting instances
119 of primates within a group, we annotated 191 frames from videos on different days (Day 1, Day
120 9, Day 16, Day 18). As exemplified in Fig. 2a, the network handles even difficult scenarios
121 very well: representative illustrations include ground-truth as well as predictions of moments

124 in which multiple primates are moving rapidly while strongly occluded at varying distances
125 from the camera. SIPEC:SegNet achieved a mAP of 0.91 ± 0.03 (mean \pm s.e.m.) using 5-fold
126 CV. When we performed the previously described ablation study, SIPEC:SegNet achieved 95%
127 of mean peak performance (mAP of 0.87 ± 0.03) with only 30 labeled frames (Fig. 2b). To
128 assess the overlap between prediction and ground truth, we report IoU and dice coefficient
129 metrics as well (Fig. 2c).

130 **Pose estimation module SIPEC:PoseNet.** We also added a pose estimation network, built on
131 an encoder-decoder architecture⁴⁰ with an EfficientNet⁴¹ backbone, to SIPEC for performing
132 pose estimation (SIPEC:PoseNet) (see Methods, Supp. Fig. 11). SIPEC:PoseNet can be used to
133 perform pose estimation on N animals (with N the total number of animals or less), yielding K
134 different coordinates for previously defined landmarks on each animal's body. The main
135 advantage of SIPEC:PoseNet in comparison to previous approaches is that it receives its inputs
136 from SIPEC:SegNet (top-down pose estimation): While bottom-up approaches such as DLC²
137 require grouping of pose estimates to individuals, our top-down approach makes the assignment
138 of pose estimates to individual animals trivial, as inference is performed on the masked image
139 of an individual animal and pose estimates within that mask are assigned to that particular
140 individual (Fig. 1c). We labeled frames with 13 standardized body parts of individual mice in
141 an OFT similarly to Sturman et. al.²⁰ to train and test the performance of SIPEC:PoseNet against
142 that of DLC². SIPEC:PoseNet achieves a Root-Mean-Squared-Error (RMSE) (see Methods) of
143 2.9 pixels in mice (Fig. 2d) for a total of 96 labeled training frames, while DLC² achieves a 3.9
144 pixel RMSE². Previously published pose estimation methods for single animals can easily be
145 substituted into our pipeline to perform multi-animal pose estimation in conjunction with
146 SIPEC:SegNet.

147

148 **Identification module SIPEC:IdNet.** The identification network (SIPEC:IdNet) (see
149 Methods, Supp. Fig. 10) allows the determination of the identity of individual animals. Given
150 SIPEC:IdNet receives input as a series (T time points) of cropped images of N (with N the total
151 number of animals or less) individuals from SIPEC:SegNet, the output of SIPEC:IdNet are N
152 identities. The input images from SIPEC:SegNet are scaled to the same average size (see
153 Methods) before being fed into SIPEC:IdNet. We designed a feedforward classification neural
154 network, which utilizes a DenseNet²⁸-backbone pre-trained on ImageNet⁴². This network serves
155 as a feature-recognition network on single frames. We then utilize past and future frames by
156 dilating the mask around the animal with each timestep. The outputs of the feature-recognition
157 network on these frames are then integrated over T timesteps using a GRU (see Methods for
158 architecture and training details). SIPEC:IdNet can integrate information from none to many
159 temporally-neighboring frames based on a particular application's accuracy and speed
160 requirements. We used spatial area dropout augmentations to increase robustness against
161 occlusions⁴³. We developed an annotation tool for a human to assign identities of individual
162 animals, in a multi-animal context, to segmentation masks in videoframes, which capture
163 primates from different perspectives (Supp. Fig. 3). This tool was used for annotating
164 identification data in the following sections. Below we compared SIPEC:IdNet's performance
165 to that of the current state-of-the-art i.e. the identification module of idTracker.ai⁴ and the
166 primnet³¹ network for primate re-identification. primnet³¹ relies on faces of individuals being
167 clearly visible for re-identification, which in our case is not possible for most of the video

168 frames. idTracker.ai⁴ is a self-supervised algorithm for tracking the identity of individual
169 animals within a single session. Particularly in complex or enriched home-cage environments,
170 where animals are frequently obstructed as they move underneath/behind objects or enter/exit
171 the scene and background or lighting conditions change constantly, temporally based tracking
172 and identification as idtracker.ai performs it becomes impossible. We evaluated the
173 identification performance of SIPEC:IdNet across sessions with the identification module of
174 idTracker.ai, providing each network with identical training and testing data. While idtracker.ai
175 behaves self-supervised, the identification module it uses to distinguish animals is trained with
176 the labels generated by idTracker.ai's cascade algorithm in a supervised fashion. Apart from re-
177 identifying animals across sessions using SIPEC:IdNet, SIPEC:SegNet segmentation masks
178 can be used via greedy-mask matching (see Methods) to track the identities of animals
179 temporally as well (Supp. Videos 2-4) or to smooth the outputs of SIPEC:IdNet as a secondary
180 step, that can boost performance for continuous video sequences, but this advantage was not
181 used in the following evaluations for mice and primates.
182

183 *Identification of mice in an open-field test.* We first evaluated the performance of SIPEC:IdNet
184 in identifying 8 individual mice. We acquired 10-minute-long videos of these mice behaving
185 in the previously mentioned OFT (see Methods for details). While for the human observer,
186 these mice are difficult to distinguish (Supp. Fig. 4), our network copes rather well. We used
187 5-fold CV to evaluate the performance, i.e. splitting the 10-minute videos into 2-minute long
188 ones, while using one fold for testing and the rest to train the network. Since this data is
189 balanced, we use the accuracy metric for evaluation. We find that SIPEC:IdNet achieves $99 \pm$
190 0.5% (mean and s.e.m.) accuracy, while the current state of the art idTracker.ai⁴ only achieves
191 $87 \pm 0.2\%$ accuracy (Fig. 2e). The ablation study shows that only 650 labeled frames (frame
192 and identity of the animal) are sufficient for the SIPEC:IdNet to achieve 95% of its mean peak
193 performance (Fig. 2f). We tested how this performance translates to identifying the same
194 animals during the subsequent days (Supp. Fig. 5). We find that identification performance is
195 similarly high on the second day $86 \pm 2\%$, using the network trained on day 1. Subsequently,
196 we tested identification robustness with respect to the interventions on day 3. Following a
197 forced swim test, the identification performance of SIPEC:IdNet, trained on data of day 1,
198 dropped dramatically to $4 \pm 2\%$. This indicates that features utilized by the network to identify
199 the mice are not robust to this type of intervention, i.e. their behavior and outlook is altered by
200 the stress and residual water on the fur significantly.
201

202 *Identification of individual primates in a group.* To evaluate SIPEC: IdNet's performance on
203 identifying individual primates within a group, we used the SIPEC:SegNet-processed videos
204 of the 4 macaques (see Section "Segmentation performance of groups of primates"). We
205 annotated frames from 7 videos taken on different days, with each frame containing multiple
206 individuals, yielding approximately 2200 labels for cutouts of individual primates. We used
207 leave-one-out CV with respect to the videos in order to test SIPEC:IdNet generalization across
208 days. Across sessions SIPEC:IdNet reaches an accuracy of $78 \pm 3\%$ (mean \pm s.e.m.) while
209 idTracker.ai⁴ achieves only $33 \pm 3\%$ and primnet³¹ $34 \pm 3\%$ (Fig. 2e), where the human expert
210 (i.e. ground truth) had the advantage of seeing all the video frames and the entire cage (i.e. the
211 rest of the primates). We did a separate evaluation of the identification performance on "typical
212 frames" i.e., the human expert can correctly identify the primates using single frames. In this
213 case, SIPEC:IdNet achieved a performance of 86 ± 3 (Supp. Fig. 6). The identification labels
214 can then be further enhanced by greedy mask-match-based tracking (see Methods for details).
215 Supp. Video 1 illustrates the resulting performance on a representative video snippet. We
216 perform here an ablation study as well, which yields 95% of mean peak performance at 1504
217 annotated training samples (Fig. 2g).

218

219 **Behavioral classification module SIPEC:BehaveNet.** SIPEC:BehaveNet (see Methods,
220 Supp. Fig. 13) offers researchers a powerful means to recognize specific animal behaviors
221 directly from raw pixels using a single neuronal net framework. SIPEC:BehaveNet uses video
222 frames of N individuals over T time steps to classify the animals' actions. The video frames of
223 the N individuals are generated by SIPEC:SegNet. If only a single animal is present in the
224 video, SIPEC:BehaveNet can be used directly without SIPEC:SegNet. We use a recognition
225 network to extract features from single frames analysis, based on the Xception³² network
226 architecture. We initialize parts of the network with ImageNet⁴ weights. These features are
227 then integrated over time by a TCN^{33,34} to classify the animal's behavior in each frame (see
228 Methods for architecture and training details).

229

230 *SIPEC behavior recognition outperforms DLC-based approach.* We compare our raw-pixel-
231 based approach to Sturman et al.²⁰, who recently demonstrated that they can classify behavior
232 based on DLC² generated features. On top of a higher classification performance with fewer
233 labels, SIPEC:BehaveNet does not require annotation and training for pose estimation if the
234 researcher is interested in behavioral classification alone. The increased performance with
235 fewer labels comes at the cost of a higher computational demand since we increased the
236 dimensionality of the input data by several orders of magnitude (12 pose estimates vs. 16384
237 pixels). We used the data and labels from Sturman et al.²⁰ on 20 freely behaving mice in an
238 OFT to test our performance. The behavior of these mice was independently annotated by 3
239 different researchers on a frame-by-frame basis using the VGG video annotation tool³⁹.
240 Annotations included the following behaviors: supported rears, unsupported rears, grooming
241 and none (unlabeled/default class). While Sturman et al.²⁰ evaluated the performance of their
242 behavioral event detection by averaging across chunks of time, evaluating the frame-by-frame
243 performance is more suitable for testing the actual network performance since it was trained
244 the same way. Doing such frame-by-frame analysis shows that SIPEC:BehaveNet has fewer
245 false positives as well as false negatives with respect to the DLC-based approach of Sturman

et al.²⁰. We illustrate a representative example of the performance of both approaches for each of the behaviors with their respective ground truths (Fig. 3a). We further resolved spatially the events that were misclassified by Sturman et al., that were correctly classified by SIPEC:BehaveNet and vice versa (Fig. 3b). We calculated the percentage of mismatches, that occurred in the center or the surrounding area. For grooming events mismatches of Sturman et al.²⁰ and SIPEC:BehaveNet occurs similarly often in the center $41 \pm 12\%$ (mean and s.e.m.) and $42 \pm 12\%$ respectively. For supported and unsupported rearing events Sturman et al.²⁰ has more mismatches occurring in the center compared to SIPEC:BehaveNet (supported rears: $40 \pm 4\%$ and $37 \pm 6\%$, unsupported rears: $12 \pm 2\%$ and $7 \pm 2\%$). This indicates that the misclassifications of the pose estimation-based approach are more biased towards the center than the ones of SIPEC:BehavNet. To quantify the behavioral classification over the whole time course of all videos of 20 mice, we used leave-one-out CV (Fig. 3c). We used macro-averaged F1-score as a common metric to evaluate a multi-class classification task and Pearson correlation (see Methods for metrics) to indicate the linear relationship between the ground truth and the estimate over time. For the unsupported rears/grooming/supported rears behaviors SIPEC:BehaveNet achieves F1-Scores of 0.6 ± 0.16 / 0.49 ± 0.21 / 0.84 ± 0.04 (values reported as mean \pm s.e.m.) respectively, while the performance of the manually intensive Sturman et al.²⁰'s approach reaches only 0.49 ± 0.11 / 0.37 ± 0.2 / 0.84 ± 0.03 , leading to a significantly higher performance of SIPEC:BehaveNet for the unsupported rearing ($F1: p=1.689 \times 10^{-7}$, Wilcoxon paired-test was used as recommended⁴⁴) as well as the grooming ($F1: p=6.226 \times 10^{-4}$) behaviors. While we see a higher precision only in the classification of supported rears in the DLC-based approach, SIPEC:BehaveNet has an improved recall for the supported rears as well as improved precision and recall for the other behaviors (Supp. Fig. 7a). As expected, more stereotyped behaviors with many labels like supported rears yield higher F1. In comparison, less stereotypical behaviors like grooming with fewer labels have lower F1 for SIPEC:BehaveNet and the DLC-based approach. Additionally, we computed the mentioned metrics on a dataset with shuffled labels to indicate chance performance for each metric as well as computed each metric when tested across human annotators to indicate an upper limit for frame-by-frame behavioral classification performance (Supp. Fig. 7b). While the overall human-to-human F1 is 0.79 ± 0.07 (mean \pm s.e.m.), SIPEC:BehaveNet classifies with an F1 of 0.71 ± 0.07 . We then grouped behaviors by integrating the classification over multiple frames as described in Sturman et al.²⁰. This analysis results in a behavior count per video. For these per video behavior counts, we found no significant difference between human annotators, SIPEC:BehaviorNet and Sturman et al.²⁰ (Tukey's multiple comparison test, Supp. Fig. 15). Such classification and counting of specific behaviors per video are commonly used to compare the number of occurrences of behaviors across experimental groups. Using such analysis, Sturman et al.²⁰ demonstrate how video-based analysis outperforms commonly used commercial systems. Moreover, we also tested combining the outputs of pose estimation-based classification together with the raw-pixel model (Combined Model in Methods, Supp. Fig. 7). Lastly, we performed a frame ablation study and showed that SIPEC:BehaveNet needs only 114 minutes, less than 2 hours, of labeled data to reach peak performance in behavior classification (Fig. 3d).

288

289 **Socially interacting primate behavior classification.** We used the combined outputs of
290 SIPEC:SegNet and SIPEC:IdNet, smoothed by greedy match-based tracking, to generate
291 videos of individual primates over time (see Methods for details). To detect social events, we
292 used SIPEC:SegNet to generate additional video events covering "**pairs**" of primates. An
293 interaction event was detected whenever the masks of individual primates came sufficiently
294 close (see Methods). We were able to rapidly annotate these videos again using the VGG video
295 annotation tool³⁹ (overall 80 minutes of video are annotated from 3 videos, including the
296 individual behaviors of object interaction, searching, **social grooming** and none (background
297 class)). We then trained SIPEC:BehaveNet to classify individuals' frames and merged frames
298 of pairs of primates socially interacting over time. We used grouped 5-fold stratified CV over
299 all annotated video frames, with labeled videos being the groups. Overall SIPEC:BehaveNet
300 achieved a macro-F1 of 0.72 ± 0.07 (mean \pm s.e.m.) across all behaviors (Fig. 4a). This
301 performance is similar to the earlier mentioned mouse behavioral classification performance.
302 The increased variance compared to the classification of mouse behavior is expected as imaging
303 conditions, as previously mentioned, are much more challenging and primate behaviors are
304 much less stereotyped compared to mouse behaviors. This can be likely compensated with more
305 training data.

306

307 **Tracking position of primates in 3D without stereo-vision.** By performing SIPEC:SegNet
308 and SIPEC:IdNet inference on a full one-hour video, we built a density map of positions of
309 individuals within the husbandry (Fig. 1a). Without stereo-vision, one cannot optically acquire
310 depth information. Instead, we used the output masks of SIPEC:SegNet and annotated the
311 positions of the primates in 300 frames using a 3D model (Supp. Fig. 8). Subsequently, we
312 generated 6 features using Isomap⁴⁵ and trained a multivariate linear regression model to predict
313 the 3D positions of the primates (Fig. 4b). Using 10-fold CV, our predicted positions using only
314 single camera have an overall RMSE of only 0.43 ± 0.01 m (mean \pm s.e.m.), that is of $0.27 \pm$
315 0.01 m in x-direction or 6% error w.r.t the room dimension in x-direction; 0.26 ± 0.01 m / 7%
316 and 0.21 ± 0.01 m / 7% for the y and z coordinates respectively. If an annotation is impossible,
317 quasi depth estimates can be calculated through the mask size alone and correlate highly with
318 the actual depth (Supp. Fig. 14).

319 **Discussion**

320 We have presented SIPEC, a novel pipeline, using specialized deep neural networks to perform
321 segmentation, identification, behavioral classification, and pose estimation on individual and
322 interacting animals. With SIPEC we address multiple key challenges in the domain of
323 behavioral analysis. Our **SIPEC:SegNet** enables the segmentation of animals with only 3-30
324 labels (Fig. 2a,b,c). In combination with greedy-mask matching, SIPEC:SegNet can be used to
325 track animals' identities within one session similar to idtracker.ai, but even in complex
326 environments with changing lighting conditions, where idtracker.ai fails (Supp. Video 1).

327 Subsequently, **SIPEC:BehaveNet** enables animal behavior recognition directly from raw video
328 data. Raw-video classification has the advantage of not requiring pre-processing adjustments
329 or feature engineering to specific video conditions. Moreover, we show that learning task-
330 relevant features directly from the raw video can lead to better results than pose-estimation-

331 based approaches which train a classifier on top of the detected landmarks. In particular, we
332 demonstrate that our network outperforms a state-of-the-art pose estimation approach¹³ on a
333 well-annotated mouse behavioral dataset (Fig. 3) and reaches human-level performance for
334 counting behavioral events (Supp. Fig. 15). Thus, pose-estimation can be skipped if researchers
335 are solely interested in classifying behavior. We note that our raw-pixel approach increases the
336 input-dimensionality of the behavior classification network and therefore uses more
337 computational resources and is slower than pose-estimation-based approaches.

338 **SIPEC:IdNet** identifies primates in complex environments across days with high accuracy.
339 SIPEC:SegNet enhances SIPEC:IdNet's high identification performance through mask-
340 matching-based tracking and integration of identities through time. We demonstrate that
341 identification accuracy is significantly higher than that of the identification module of state-of-
342 art idtracker.ai and primnet³¹ (Fig. 2e). We note, however, that identification using deep nets is
343 not robust to interventions that affect mice's appearance strongly immediately after the
344 intervention (such as forced swim test, Supp. Fig. 5). However, even without any interventions,
345 expert human observers have difficulty identifying mice of such similar size and color. The
346 effects of different interventions on the recognition performances of deep net architectures
347 should be studied in the future. Finally, **SIPEC:PosNet** enables top-down pose estimation of
348 multiple animals in complex environments, making it easy to assign pose estimates to
349 individual animals with higher performance than DLC (Fig. 2d).

350 All approaches are optimized through augmentation and transfer learning, significantly
351 speeding up learning and reducing labeling compared to the other approaches we tested on the
352 mouse and non-human primate datasets. We also performed ablation studies for each of the
353 networks to estimate the number of labels necessary for successful training. The number of
354 labels necessary can change depending on the dataset, for example, if the background, etc. are
355 more complex each network could require more annotated frames to be trained successfully.
356 To perform well under the complex video conditions for non-human primates, SIPEC:SegNet
357 needs about 30 labels, SIPEC:IdNet about 1500 labels and SIPEC:BehaveNet less than 2 hours
358 of annotated video (Fig. 2c,g; Fig. 4a).

359 SIPEC can be used to study the behavior of primates and their social interactions over longer
360 periods in a naturalistic environment, as we demonstrated for social grooming (Fig. 4a). In
361 addition, after initial training of SIPEC modules, they can automatically output a behavioral
362 profile for each individual in a group, over days or weeks and therefore also be used to quantify
363 the changes in behaviors of individuals in social contexts over time. Since SIPEC is fully
364 supervised, it may be difficult to scale it to large colonies with hundreds of animals, such as
365 bees and ants. However, SIPEC is well suited for most other animal species beyond insects.

366 Finally, we show how SIPEC enables 3D localization and tracking from a single-camera view,
367 yielding an off-the-shelf solution for home-cage monitoring of primates, without the need for
368 setting stereo-vision setups (Fig. 4b). Estimating the 3D position requires the experimenter to
369 create a 3D model and annotate 3D data. However, we show a quasi-3D estimate can be
370 generated directly from the mask size, without manual annotation, that correlates highly with
371 the actual position of the animal (Supp. Fig. 14).

372 Behaviors which were not recognized and annotated by the researcher and therefore not learned
373 by the neural network could be picked up using complementary unsupervised approaches^{12,13}.
374 The features-vectors, embedding individual behaviors, created by SIPEC:BehaveNet can be
375 used as input to unsupervised approaches, which can help align the outputs of unsupervised
376 approaches with human annotation. Moreover, the output of other modules (SIPEC:SegNet,
377 SIPEC:IdNet and SIPEC:PoseNet) can also be used after such unsupervised approaches to
378 analyse individual animals.

379 **Data Availability**

380 Mouse data from Sturman et al.²⁰ is available under <https://zenodo.org/record/3608658>. Primate
381 data is available upon reasonable request from authors. Exemplary data for training is available
382 through our github repository.
383

384 **Code Availability**

385 We provide the code for SIPEC at <https://github.com/SIPEC-Animal-Data-Analysis/SIPEC>
386 (<https://doi.org/10.5281/zenodo.5927367>) and the GUI for the identification of animals
https://github.com/SIPEC-Animal-Data-Analysis/idtracking_gui.
388

389 **Acknowledgments**

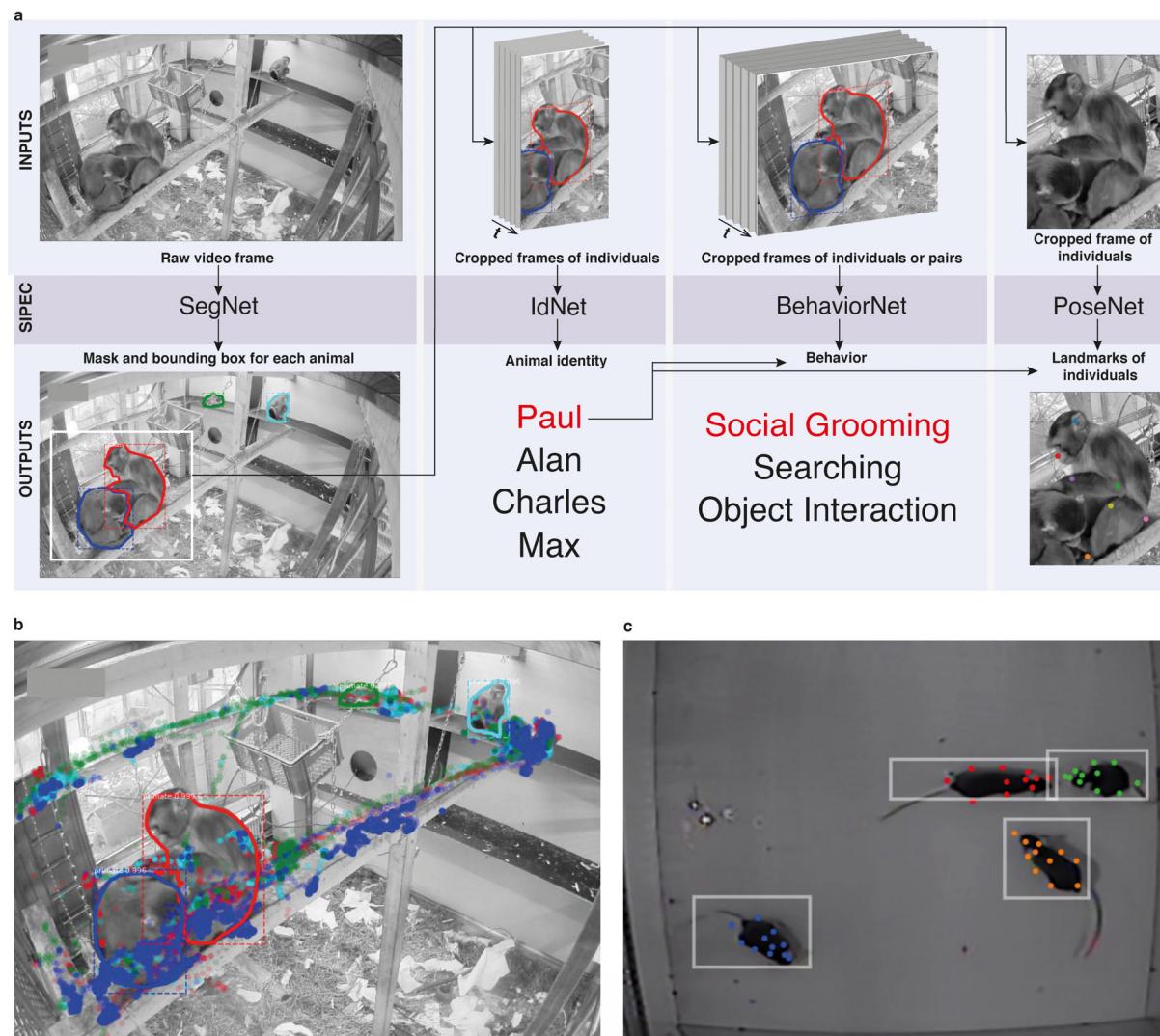
390 This project was funded by the Swiss Federal Institute of Technology (ETH) Zurich and the
391 European Research Council (ERC) under the European Union's Horizon 2020 research and
392 innovation program (grant agreement No 818179), SNSF (CRSII5_198739/1 to MFY;
393 310030_172889/1 to JB, PP00P3_157539 to VM) ETH Research Grant (ETH-20 19-1 to JB),
394 3RCC (OC-2019-009 to JB and MFY) , the Simons Foundation (awards 328189 and 543013 to
395 VM) and the Botnar Foundation (to JB). We would like to thank Petra Tornmalm and Victoria
396 de La Rochefoucauld for annotating primate data and feedback on primate behavior. We would
397 like to thank Paul Johnson, Baran Yasar, Bifeng Wu, and Aagam Shah for helpful discussions
398 and feedback.
399

400 **Author contributions**

401 M.M. developed, implemented, and evaluated the SIPEC modules and framework. J.Q.
402 developed segmentation filtering, tracking, and 3D-estimation. M.M., W.B., and M.F.Y. wrote
403 the manuscript. M.M., O.S., LvZ., S.K., W.B., V.M., J.B., and M.F.Y. conceptualized the study.
404 All authors gave feedback on the manuscript.

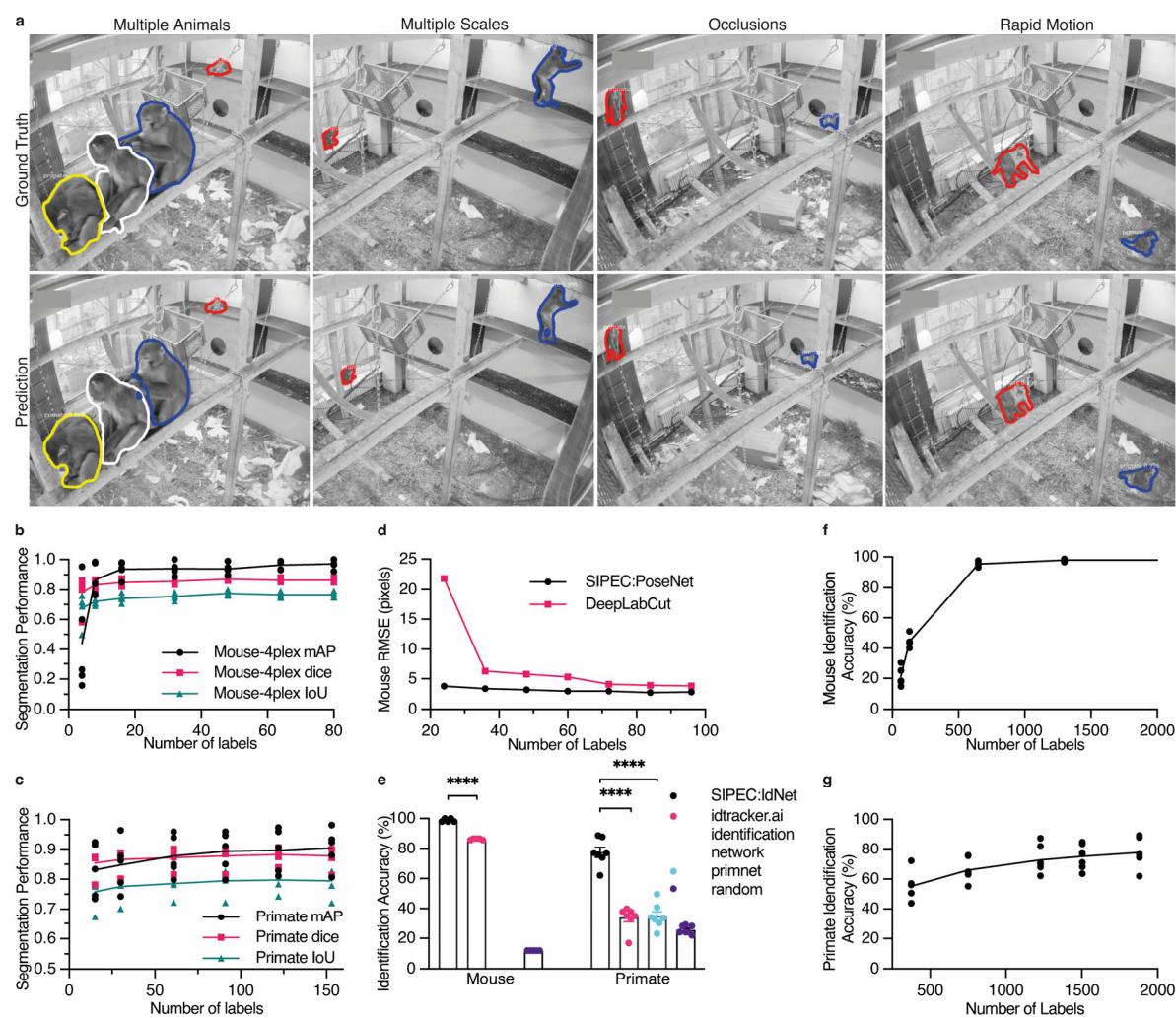
405 **Competing interests**

406 The authors declare no competing interests.
407



408

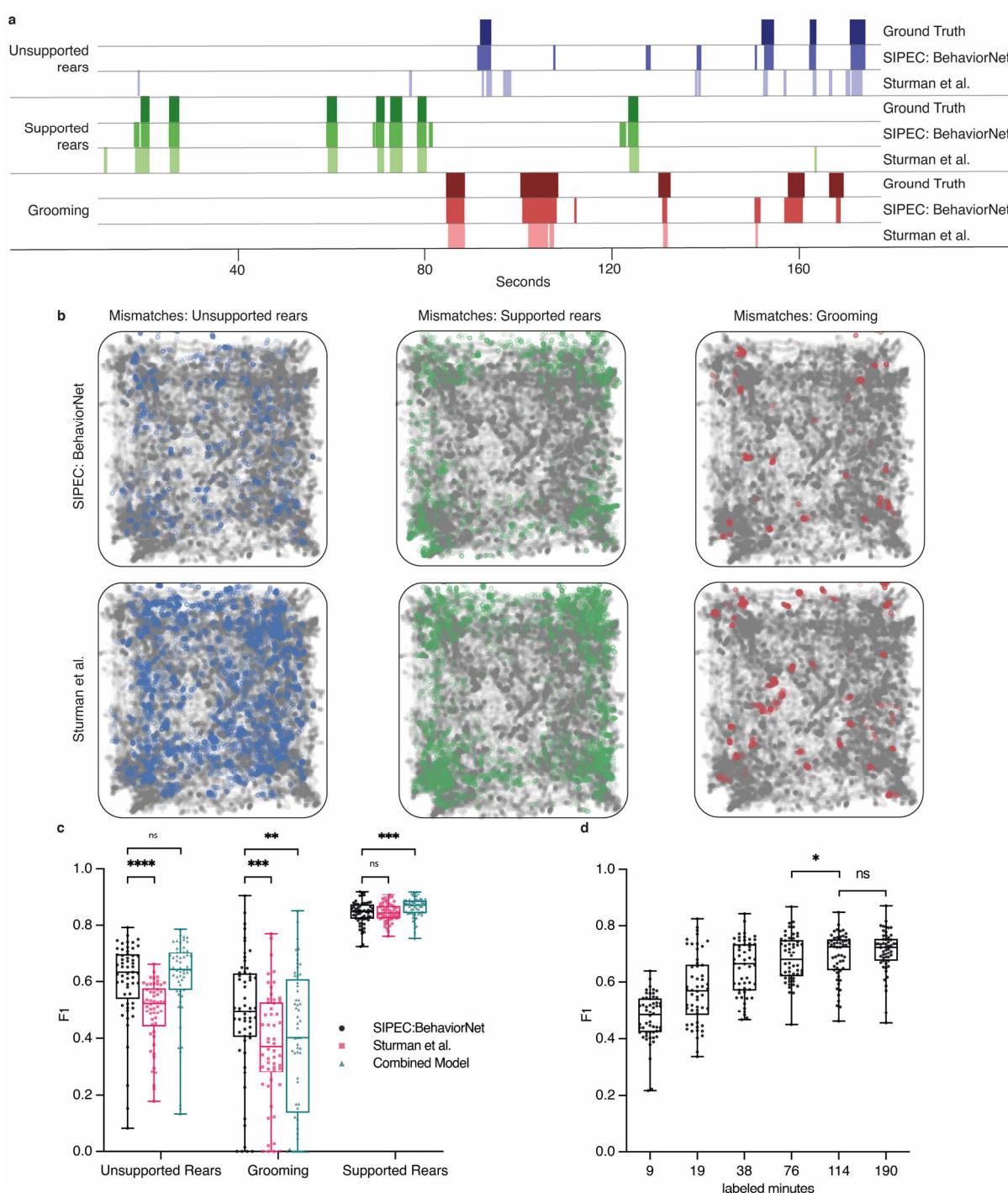
409 **Fig. 1 | Overview of the SIPEC workflow and modules.** a) From a given video, instances of
410 animals are segmented with the segmentation network (SIPEC:SegNet), indicated by masked
411 outline as well as bounding boxes. Subsequently, individuals are identified using the
412 identification network (SIPEC:IdNet). For each individual, the pose and behavior can be
413 estimated/classified using the pose estimation network (SIPEC:PoseNet) and the behavioral
414 identification network (SIPEC:BehaveNet), respectively. b) Outcome of SIPEC:SegNet, and
415 SIPEC:IdNet modules are overlaid on a representative videoframe. Time-lapsed positions of
416 individual primates (center of mass) are plotted as circles with respective colors. c) Outputs of
417 SIPEC:SegNet (boxes) and SIPEC:PoseNet (colored dots) on a representative videoframe of
418 mouse open-field data.



419

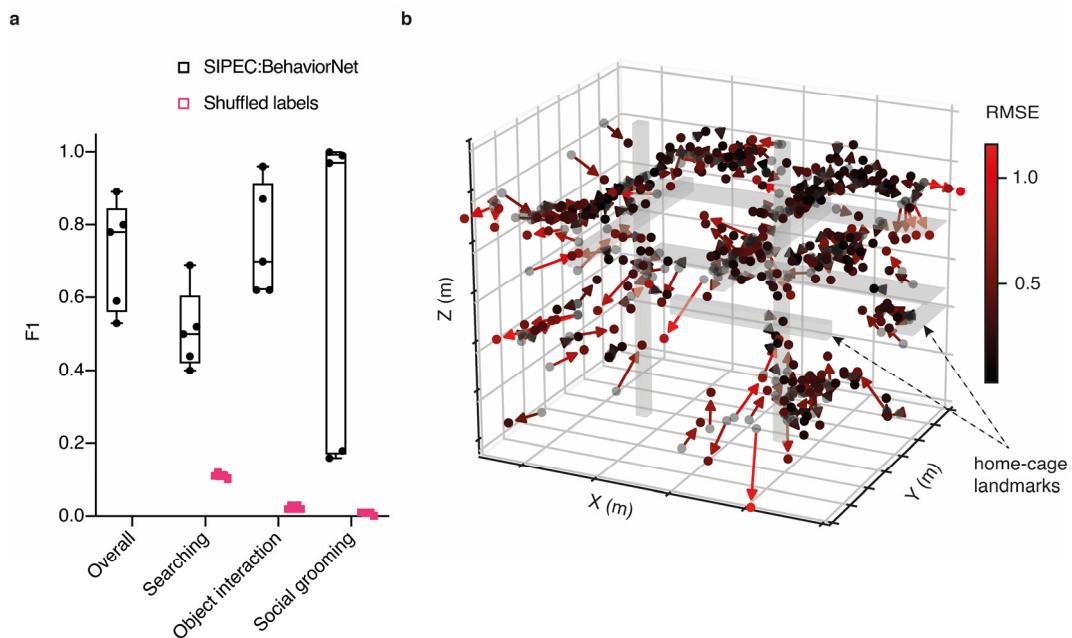
420 **Fig. 2 | Performance of the segmentation (SIPEC:SegNet), pose estimation**
 421 **(SIPEC:PoseNet), and identification (SIPEC:IdNet) modules under demanding video**
 422 **conditions and using few labels.** a) Qualitative comparison of ground truth (top row) versus
 423 predicted segmentation masks (bottom row) under challenging conditions; multiple animals, at
 424 varying distances from the camera, under strong visual occlusions, and in rapid motions. b) For
 425 mice, SIPEC:SegNet performance in mAP (mean average precision), dice (dice coefficient),
 426 and IoU (intersection over union) as a function of the number of labels. The lines indicate the
 427 means for 5-fold CV while circles, squares, triangles indicate the mAP, dice, and IoU,
 428 respectively, for individual folds. c) For primates, SIPEC:SegNet performance in mAP, dice,
 429 and IoU as a function of the number of labels. The lines indicate the means for 5-fold CV while
 430 circles, squares, triangles indicate the mAP, dice, and IoU, respectively, for individual folds. d)
 431 The performance of SIPEC:PoseNet in comparison to DeepLabCut measured as RMSE in
 432 pixels on single mouse pose estimation data. e). Comparison of identification accuracy for
 433 SIPEC:IdNet module, idtracker.ai⁴, primnet³¹ and randomly shuffled labels (chance
 434 performance). 8 videos from 8 individual mice and 7 videos across 4 different days from 4
 435 group-housed primates are used. f) For mice, the accuracy of SIPEC:IdNet as a function of the
 436 number of training labels used. The black lines indicate the mean for 5-fold CV with individual
 437 folds displayed. g) For primates, the accuracy of SIPEC:IdNet as a function of the number of

438 training labels used. The black lines indicate the mean for 5-fold CV with individual folds
 439 displayed. All data is represented by mean, showing all points.



440
 441 **Fig. 3 | SIPEC:BehaveNet outperforms pose-estimation (DeepLabCut) based approach**
 442 **(Sturman et al.²⁰)**. a) Comparison of behavioral classification by human annotator (ground
 443 truth), SIPEC:BehaveNet, and Sturman et al.²⁰ b) Errors in the classification of mouse behavior
 444 in the open arena for SIPEC:BehaveNet versus Sturman et al. Each colored dot represents a
 445 behavioral event that is incorrectly classified by that method (while correctly classified by the
 446 other) with respect to the ground truth. none-classified (background class) positions of mice are
 447 indicated as grey dots. c) Frame-by-frame classification performance per video (n=20 mice)

448 compared to ground truth. d) SIPEC:BehaveNet classification performance as a function of
449 labeled minutes. All data is represented by a Tukey box-and-whisker plot, showing all points.
450 Wilcoxon paired test: * p <= 0.05; *** p <= 0.001; **** p <= 0.0001.
451
452



453 **Fig. 4 | SIPEC can recognize social interactions of multiple primates and infer their 3D**
454 **positions using a single camera.** a) Performance of SIPEC:BehaveNet for individual and
455 social behaviors with respect to ground truth evaluated using grouped 5-fold CV. Behaviors
456 include searching, object interaction and social grooming; while the performance is measured
457 using F1. F1 on shuffled labels is included for comparison. All data is represented by a
458 minimum-to-maximum box-and-whisker plot, showing all points. b) Evaluation of 3D position
459 estimates of primates in home-cage. Black spots mark annotated positions (n=300) while
460 predicted positions are marked as red-hued spots at the end of the solid arrows (color-coded
461 using a red gradient with brighter red indicating higher RMSE of predicted to true position).

462

463 **Methods**

464 **Animals.** C57BL/6J (C57BL/6JRj) mice (male, 2.5 months of age) were obtained from Janvier
465 (France). Mice were maintained in a temperature- and humidity-controlled facility on a 12-h
466 reversed light-dark cycle (lights on at 08:15 am) with food and water ad libitum. Mice were
467 housed in groups of 5 per cage and used for experiments when 2.5–4 months old. For each
468 experiment, mice of the same age were used in all experimental groups to rule out confounding
469 effects of age. All tests were conducted during the animals' active (dark) phase from 12–5 pm.
470 Mice were single housed 24 h before behavioral testing in order to standardize their
471 environment and avoid disturbing cage mates during testing. The animal procedures of these
472 studies were approved by the local veterinary authorities of the Canton Zurich, Switzerland,
473 and carried out in accordance with the guidelines published in the European Communities
474 Council Directive of November 24, 1986 (86/609/EEC).

475 **Acquisition of mouse data.** For mouse behavioral data and annotation, we refer to Sturman et
476 al.²⁰. For each day, we randomized the recording chamber of mice used. On days 1-2, we
477 recorded animals 1-8 individually. On day 3, for measuring the effect of interventions on
478 performance, mice were forced-swim-tested in water for 5 minutes immediately before the
479 recording sessions.

480 **Acquisition of primate data.** 4 male rhesus macaques were recorded with a 1080p camera
481 within their home-cage. The large indoor room was about 15m². Videos were acquired using a
482 Bosch Autodome IP starlight 7000 HD camera with 1080p resolution at 50 Hz.

483 **Annotation of segmentation data.** To generate training data for segmentation training, we
484 randomly extracted frames of mouse and primate videos using a standard video player. Next,
485 we used the VIA video annotator³⁹ to draw outlines around the animals.

486 **Generation and annotation of primate behavioral videos.** For creating the dataset, 3 primate
487 videos of 20-30 minutes were annotated using the VIA video annotator³⁹. These videos were
488 generated by previous outputs of SIPEC:SegNet and SIPEC:IdNet. Frames of primates,
489 identified as the same over consecutive frames, were stitched together to create individualized
490 videos. To generate videos of social interactions, we dilated the frames of each primate in each
491 frame and checked if their overlap crossed a threshold, in which case we recalculated the COM
492 of those two masks and center-cropped the frames around them. Labeled behaviors included
493 'searching', 'object interacting', 'social grooming' and 'none' (background class).

494 **Tracking by segmentation and greedy mask-matching.** Based on the outputs of the
495 segmentation masks, we implemented greedy-match-based tracking. For a given frame the
496 bounding box of a given animal is assigned to the bounding box previous frames with the
497 largest spatial overlap, with a decaying factor for temporally distant frames. The resulting
498 overlap can be used as a confidence of SIPEC:SegNet based tracking of the individual. This
499 confidence can be used as a weight when using the resulting track identities to optionally
500 smooth the labels that SIPEC:IdNet.

501 **Identification labeling with the SIPEC toolbox.** As part of SIPEC we release a GUI that
502 allows to label for identification when multiple animals are present (Supp. Fig. 3). To use the
503 GUI, SIPEC:SegNet has to be trained and inference has to be performed on videos to be identity
504 labeled. SIPEC:SegNet results can then be loaded from the GUI and overlaid with the original
505 videos. Each box then marks an instance of the species that is to be labeled in green. For each
506 animal, a number on the keyboard can be defined, which corresponds to the permanent ID of
507 the animal. This keyboard number is then pressed, and the mask-focus jumps to the next mask
508 until all masks in that frame are annotated. Subsequently, the GUI jumps to the next frame in
509 either regular intervals or randomly throughout the video, as predefined by the user. Once a
510 predefined number of masks is reached, results are saved, and the GUI is closed.

511 **SIPEC top-down workflow.** For a given image, if we assume that N individuals (with N the
512 total number of animals or less) are in the field of view (FOV), the output of SIPEC:SegNet is
513 N segmentations or masks of the image. This step is mandatory if the analysis is for multiple
514 animals in a group since subsequent pipeline parts are applied to the individual animals. Based
515 on the masks, the individual animals' center of masses (COMs) are calculated as a proxy for
516 the animals' 2D spatial positions. Next, we crop the original image around the COMs of each
517 animal, thus reducing the original frame to N COMs and N square-masked cutouts of the
518 individuals. This output can then be passed onto other modules.

519 **SIPEC:SegNet network architecture and training.** SIPEC:SegNet was designed by
520 optimizing the Mask R-CNN architecture. We utilized a ResNet101 and feature pyramid
521 network (FPN)⁴⁶ as the basis of a convolutional backbone architecture. These features were fed
522 to the region proposal network (RPN), which applies convolutions onto these feature maps and
523 proposes regions of interest (ROIs). Subsequently, these are passed to a ROIAlign layer, which
524 performs feature pooling, while preserving the pixel-correspondence in the original image. Per
525 level of this pyramidal ROIAlign layer, we assign an ROI feature map from the different layers
526 of the FPN feature maps. Multiple outputs are generated from the FPN, one of which is
527 classifying if an animal is identified. The regressor head of the FPN returns bounding-box
528 regression offsets per ROI. Another fully convolutional layer, followed by a per-pixel sigmoid
529 activation, performs the mask prediction, returning a binary mask for each animal ROI. The
530 network is trained using stochastic gradient descent, minimizing a multi-task loss for each ROI:

531
$$L = L_{mask} + L_{regression} + L_{class}$$

532 where L_{mask} is the average binary cross-entropy between predicted and ground truth
533 segmentation mask, applied to each ROI. $L_{regression}$ is a regression loss function applied to the
534 coordinates of the bounding boxes, modified to be outlier robust as in the original Fast R-CNN
535 paper⁴⁷. L_{class} is calculated for each of the proposed ROIs (or anchors) as a logarithmic loss of
536 non-animal vs animal. The learning rate was adapted by an animal specific schedule and
537 training was done iteratively, by first training the output layers for some epochs and then
538 incrementally including previous blocks in the training process. SIPEC:SegNet outputs
539 segmentation masks and bounding boxes to create cutouts or masked cutouts of individual
540 animals to be used by one of the downstream modules.

541 **SIPEC:IdNet network architecture and training.** SIPEC:IdNet was based on the DenseNet
542 architecture²⁸ for frame-by-frame identification. It consists of 4 dense blocks, which consist of
543 multiple sequences of a batch normalization layer, a ReLU activation, and a convolution. The
544 resulting feature maps are concatenated to the outputs of the following sequences of layers
545 (skip-connections). The resulting blocks are connected through transitions, that are
546 convolutional followed by pooling layers. After the last dense block, we connect an average
547 pooling layer to a Dropout⁴⁸ layer with a dropout rate of 0.5 followed by the softmax
548 classification layer. For the recurrent SIPEC:IdNet, we remove the softmax layer and feed the
549 output of the average pooling layers for each time point into a batch normalization layer⁴⁹
550 followed by 3 layers of bidirectional gated recurrent units^{29,30} with leaky ReLU activation^{50,51}
551 (alpha=0.3) followed by a Dropout⁴⁸ layer with rate 0.2 followed by the softmax layer. The
552 input for SIPEC:IdNet is the output cutouts of individuals, generated by SIPEC:SegNet (for the
553 single-animal case background-subtracted thresholding and centered-cropping would also
554 work). For the recurrent case, the masks of past or future frames are dilated with a frames per
555 second (FPS) dependent factor that increases with distance in time in order to increase the field
556 of view. We first pre-trained the not-recurrent version of SIPEC:IdNet using Adam⁵² with an
557 lr=0.00025, a batch size of 16 and using a weighted cross-entropy loss. We used a learning rate
558 scheduler in the following form:

559
$$L_{E+1} = \frac{L_E}{k^E} (2)$$

560 where E stands for epoch, using a k=1.5. Subsequently, we removed the softmax layer and
561 fixed the network's weights. We then trained the recurrent SIPEC:IdNet again using Adam⁵²
562 and an lr=0.00005, k=1.25 and a batch size of 6.

563 **SIPEC:BehaveNet network architecture and training.** SIPEC:BehaveNet was constructed
564 as a raw-pixel action recognition network. It consists of a feature recognition network that
565 operates on a single frame basis and a network, which integrates these features over time. The
566 feature recognition network (FRN) is based on the Xception³² architecture, consisting of an
567 entry, middle, and exit flow. The entry flow initially processes the input with convolution and
568 ReLU blocks. Subsequently, we pass the feature maps through 3 blocks of separable
569 convolution layers, followed by ReLU, separable convolution, and a max-pooling layer. The
570 outputs of these 3 blocks are convolved and concatenated and passed to the middle flow. The
571 middle flow consists of 8 blocks of ReLU layers followed by a separable convolution layer.
572 The Exit receives the feature maps from the middle flow and passes it one more entry-flow-
573 like block, followed by separable convolution and ReLU units. Finally, these features are
574 integrated by a global average pooling layer, followed by a dense layer and passed through the
575 softmax activation. This FRN was first pre-trained on a frame-by-frame basis using an
576 lr=0.00035, gradient clipping norm of 0.5, and batch size=36 using the Adam⁵² optimizer. We
577 reduced the original Xception architecture by the first 17 layers for mouse data to speed up the
578 computation and reduce overfitting. After training the FRN, the outputting dense and softmax
579 layers were removed, and all weights were fixed for further training. The FRN-features were
580 integrated over time by a non-causal Temporal Convolution Network³³. It is non-causal because,
581 for classification of behavior at time point t , it combines features from $[t-n, t+n]$ with n being

582 the number of timesteps, therefore looking backward in time and forward. In this study, we
583 used an n of 10. The FRN features are transformed by multiple TCN blocks of the following
584 form: 1D-Convolution followed by batch normalization, a ReLU activation and spatial dropout.
585 The optimization was performed using Adam⁵² as well with a learning rate of 0.0001 and a
586 gradient clipping norm of 0.5, trained with a batch size of 16.

587 **Loss adaptation.** To overcome the problem of strong data imbalance (most frames are
588 annotated as 'none', i.e. no labeled behavior), we used a multi-class adaptation technique Focal
589 loss⁵³, commonly used for object detection, and adapt it for action recognition, to discount the
590 contribution of the background class to the overall loss:

591

$$L_{focal} = -\alpha(1 - p_t)^\gamma \log p_t$$

592 We used a gamma = 3.0 and an alpha = 0.5. For evaluation, we used the commonly used *F1*
593 metric to assess multi-class classification performance while using *Pearson Correlation* to
594 assess temporal correlation.

595 **SIPEC:PoseNet network architecture and training.** Combined with SIPEC:SegNet we can
596 perform top-down pose estimation with SIPEC:PoseNet. That means, instead of the pose
597 estimation network outputting multiple possible outputs for one landmark, corresponding to
598 different animals, we can first segment different animals and then run SIPEC:PoseNet per
599 animal on its cropped frame. In principle, every architecture can now be run on the cropped
600 animal frame, including DLC². The SIPEC:PoseNet architecture is based on an encoder-
601 decoder design⁴⁰. In particular, we used EfficientNet⁴¹ as a feature detection network for a
602 single frame. Subsequently, these feature maps are deconvolved into heatmaps that regress
603 towards the target location of that landmark. Each deconvolutional layer is followed by a batch
604 normalization layer and a ReLU activation function layer. For processing target images for
605 pose-regression, we convolved pose landmark locations in the image with a 2D Gaussian
606 kernel. Since there were many frames with an incomplete number of labels, we defined a
607 custom cross-entropy-based loss function, which was 0 for non-existing labels.

608

$$L_{incomplete} = \begin{cases} \text{CrossEntropy} \\ 0, \text{if labels does not exist} \end{cases}$$

609 **Combined Model.** To test performance effects of doing a pose-estimation-based classification
610 in conjunction with SIPEC:BehaveNet, we pre-trained SIPEC:PoseNet (with classification
611 layer on top) as well as SIPEC:BehavNet individually. Subsequently removed the output layers
612 and fixed the weights of the individual networks and trained a joint output model, which
613 combined inputs of each stream followed by a batch normalization layer, a dense layer (64
614 units), and a ReLU activation layer. The resulting units were concatenated into a joint tensor
615 followed by a batch normalization layer, a dense layer (32 units), and a ReLU activation layer.
616 This layer was followed by a dense layer with 4 units for the 4 behavioral classes and softmax
617 activation function. This combined model was trained using Adam⁵² with a lr=0.00075. We

618 further offer to use optical flow as an additional input, which has been shown to enhance action
619 recognition performance⁵⁴.

620 **Implementation and Hardware.** For all neural network implementations, we used
621 Tensorflow⁵⁵ and Keras⁵⁶. Computations were done on either NVIDIA RTX 2080 Ti or V100
622 GPUs.

623 **3D location labeling.** To annotate the 3D location of a primate, we firstly create a precise
624 model of the physical room (Supp. Fig. 8) using Blender. For a given mask-cutout of a primate,
625 we place an artificial primate at an approximate location in the 3D model. We can then directly
626 read out the 3D position of the primate. 300 samples are annotated, covering the most frequent
627 parts of the primate positions.

628 **3D location estimation.** To regress the animal positions in 3D, we trained a manifold
629 embedding using Isomap⁴⁵ using the mask size (normalized sum of positively classified pixels),
630 the x and y pixel positions and their pairwise multiplications as features. We used the resulting
631 6 Isomap features, together with the inverse square root of the mask size, mask size and x-y-
632 position in pixel space to train an ordinary least squares regression model to predict the 3D
633 position of the animal.

634 **Metrics used.** Abbreviations used: Pearson – Pearson Correlation, RMSE – Root mean squared
635 error, IoU – intersection over union, mAP – mean average precision, dice – dice coefficient.

$$636 \quad Pearson_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

$$637 \quad RMSE = \sqrt{\frac{\sum_{n=1}^N (\hat{y}_n - y_n)^2}{N}}$$

$$638 \quad precision = \frac{TP}{TP + FP}$$

$$639 \quad recall = \frac{TP}{TP + FN}$$

640 Where TP denotes True Positives, FP False Positives, TN True Negatives, and FN False
641 Negatives.

642

$$643 \quad F1 = 2 \cdot \frac{precision \cdot recall}{precision + recall}$$

644

$$IoU(M_{GT}, M_P) = \frac{M_{GT} \cap M_P}{M_{GT} \cup M_P}$$

645 Where M_{GT} denotes the ground truth mask and M_P the predicted one. We now calculate the
646 mAP for detections with an IoU > 0.5 as follows:

647

$$mAP = \sum_{n=0} (\mathbf{r}_{n+1} - \mathbf{r}_n) \rho_{interp}(\mathbf{r}_{n+1})$$

648 With

649

$$\rho_{interp}(r_{n+1}) = \max_{\tilde{r}: \tilde{r} \geq r_{n+1}} \rho(\tilde{r})$$

650 Where $\rho(r)$ denotes precision measure at a given recall value.

651

$$dice = \frac{2 * M_{GT} \cap M_P}{|M_{GT}| + |M_P|}$$

652

653 References

- 654
655 1. Datta, S. R., Anderson, D. J., Branson, K., Perona, P. & Leifer, A. Computational Neuroethology: A Call to
656 Action. *Neuron* **104**, 11–24 (2019).
- 657 2. Mathis, A. et al. DeepLabCut: markerless pose estimation of user-defined body parts with deep learning.
658 (2018).
- 659 3. Geuther, B. Q. et al. Robust mouse tracking in complex environments using neural networks. *Communications
660 biology* **2**, 124 (2019).
- 661 4. Romero-Ferrero, F., Bergomi, M. G., Hinz, R. C., Heras, F. J. & de Polavieja, G. G. idtracker. ai: tracking all
662 individuals in small or large collectives of unmarked animals. *Nature methods* **16**, 179 (2019).
- 663 5. Forys, B., Xiao, D., Gupta, P., Boyd, J. D. & Murphy, T. H. Real-time markerless video tracking of body parts in
664 mice using deep neural networks. *bioRxiv* 482349 (2018).
- 665 6. Pereira, T. D. et al. Fast animal pose estimation using deep neural networks. *Nature methods* **16**, 117 (2019).
- 666 7. Graving, J. M. et al. DeepPoseKit, a software toolkit for fast and robust animal pose estimation using deep
667 learning. *eLife* **8**, e47994 (2019).
- 668 8. Bala, P. C. et al. Automated markerless pose estimation in freely moving macaques with OpenMonkeyStudio.
669 *Nature Communications* **11**, 4560 (2020).
- 670 9. Günel, S. et al. DeepFly3D, a deep learning-based approach for 3D limb and appendage tracking in tethered,
671 adult *Drosophila*. *eLife* **8**, e48571 (2019).
- 672 10. Chen, Z. et al. AlphaTracker: A Multi-Animal Tracking and Behavioral Analysis Tool. 2020.12.04.405159 (2020)
673 doi:10.1101/2020.12.04.405159.
- 674 11. Lauer, J. et al. Multi-animal pose estimation and tracking with DeepLabCut.
675 <http://biorxiv.org/lookup/doi/10.1101/2021.04.30.442096> (2021) doi:10.1101/2021.04.30.442096.
- 676 12. Wiltschko, A. B. et al. Mapping sub-second structure in mouse behavior. *Neuron* **88**, 1121–1135 (2015).
- 677 13. Hsu, A. I. & Yttri, E. A. B-SOI_D: An Open Source Unsupervised Algorithm for Discovery of Spontaneous
678 Behaviors. <http://biorxiv.org/lookup/doi/10.1101/770271> (2019) doi:10.1101/770271.
- 679 14. Berman, G. J., Choi, D. M., Bialek, W. & Shaevitz, J. W. Mapping the stereotyped behaviour of freely moving
680 fruit flies. *Journal of The Royal Society Interface* **11**, 20140672 (2014).

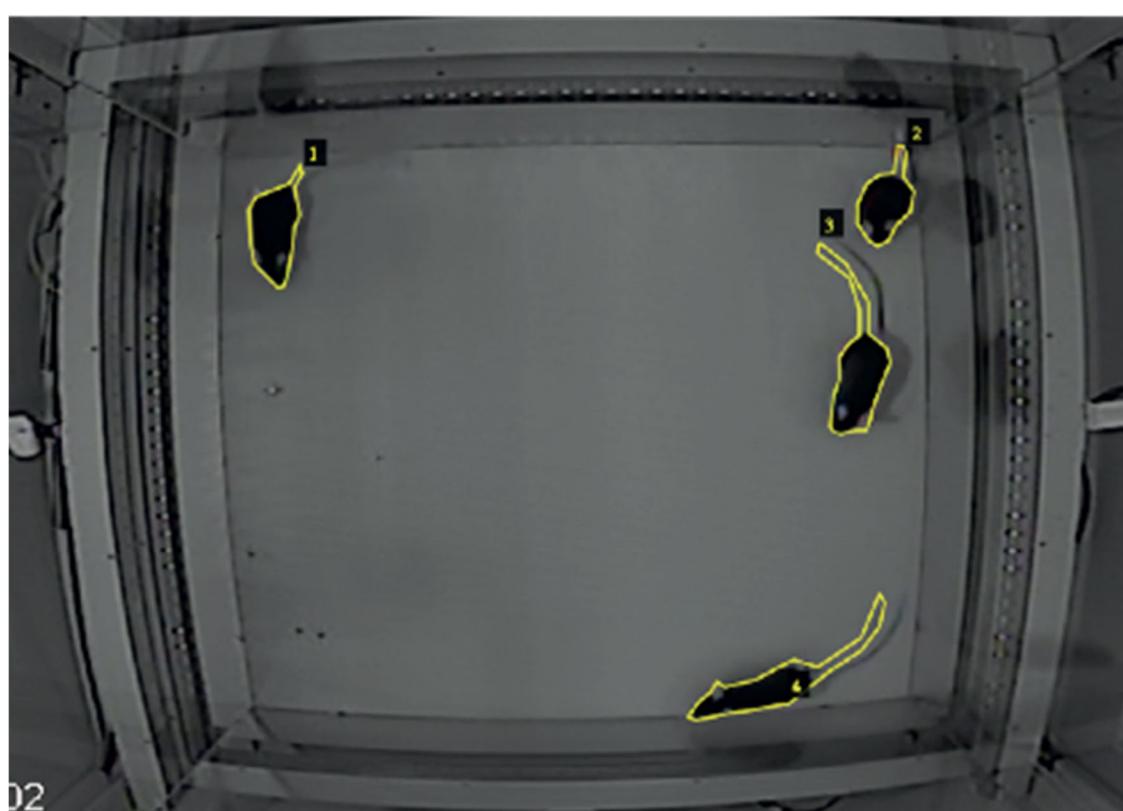
- 681 15.Whiteway, M. R. *et al.* Partitioning variability in animal behavioral videos using semi-supervised variational
682 autoencoders. *PLOS Computational Biology* **17**, e1009439 (2021).
- 683 16.Calhoun, A. J., Pillow, J. W. & Murthy, M. Unsupervised identification of the internal states that shape natural
684 behavior. *Nat Neurosci* **22**, 2040–2049 (2019).
- 685 17.Batty, E. *et al.* BehaveNet: nonlinear embedding and Bayesian neural decoding of behavioral videos. *12.*
- 686 18.Nilsson, S. R. *et al.* *Simple Behavioral Analysis (SimBA) – an open source toolkit for computer classification of*
687 *complex social behaviors in experimental animals.*
688 <http://biorxiv.org/lookup/doi/10.1101/2020.04.19.049452> (2020) doi:10.1101/2020.04.19.049452.
- 689 19.Segalin, C. *et al.* *The Mouse Action Recognition System (MARS): a software pipeline for automated analysis of*
690 *social behaviors in mice.* <http://biorxiv.org/lookup/doi/10.1101/2020.07.26.222299> (2020)
691 doi:10.1101/2020.07.26.222299.
- 692 20.Sturman, O. *et al.* Deep learning-based behavioral analysis reaches human accuracy and is capable of
693 outperforming commercial solutions. *Neuropsychopharmacology* **1**–13 (2020) doi:10.1038/s41386-020-
694 0776-y.
- 695 21.Nourizonoz, A. *et al.* EthoLoop: automated closed-loop neuroethology in naturalistic environments. *Nature*
696 *Methods* **17**, 1052–1059 (2020).
- 697 22.Branson, K., Robie, A. A., Bender, J., Perona, P. & Dickinson, M. H. High-throughput ethomics in large groups
698 of Drosophila. *Nat Methods* **6**, 451–457 (2009).
- 699 23.Dankert, H., Wang, L., Hoopfer, E. D., Anderson, D. J. & Perona, P. Automated monitoring and analysis of
700 social behavior in Drosophila. *Nat Methods* **6**, 297–303 (2009).
- 701 24.Kabra, M., Robie, A. A., Rivera-Alba, M., Branson, S. & Branson, K. JAABA: interactive machine learning for
702 automatic annotation of animal behavior. *Nature methods* **10**, 64 (2013).
- 703 25.Jhuang, H. *et al.* Automated home-cage behavioural phenotyping of mice. *Nat Commun* **1**, 68 (2010).
- 704 26.Hayden, B. Y., Park, H. S. & Zimmermann, J. Automated pose estimation in primates. *American Journal of*
705 *Primateology* **n/a**, e23348.
- 706 27.He, K., Gkioxari, G., Dollár, P. & Girshick, R. Mask r-cnn. in *Proceedings of the IEEE international conference*
707 *on computer vision* 2961–2969 (2017).
- 708 28.Huang, G., Liu, Z., Van Der Maaten, L. & Weinberger, K. Q. Densely connected convolutional networks. in
709 *Proceedings of the IEEE conference on computer vision and pattern recognition* 4700–4708 (2017).

- 710 29.Cho, K. *et al.* Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine
711 Translation. *arXiv:1406.1078 [cs, stat]* (2014).
- 712 30.Chung, J., Gulcehre, C., Cho, K. & Bengio, Y. Empirical Evaluation of Gated Recurrent Neural Networks on
713 Sequence Modeling. *arXiv:1412.3555 [cs]* (2014).
- 714 31.Deb, D. *et al.* Face Recognition: Primates in the Wild. *arXiv:1804.08790 [cs]* (2018).
- 715 32.Chollet, F. Xception: Deep learning with depthwise separable convolutions. in *Proceedings of the IEEE*
716 *conference on computer vision and pattern recognition* 1251–1258 (2017).
- 717 33.Oord, A. van den *et al.* Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499* (2016).
- 718 34.Bai, S., Kolter, J. Z. & Koltun, V. An empirical evaluation of generic convolutional and recurrent networks for
719 sequence modeling. *arXiv preprint arXiv:1803.01271* (2018).
- 720 35.Jung, A. B. *et al.* *imgaug*. (2020), GitHub repository: <https://github.com/aleju/imgaug>.
- 721 36.Yosinski, J., Clune, J., Bengio, Y. & Lipson, H. How transferable are features in deep neural networks? in
722 *Advances in neural information processing systems* 3320–3328 (2014).
- 723 37.He, K., Zhang, X., Ren, S. & Sun, J. Deep residual learning for image recognition. in *Proceedings of the IEEE*
724 *conference on computer vision and pattern recognition* 770–778 (2016).
- 725 38.Lin, T.-Y. *et al.* Microsoft coco: Common objects in context. in *European conference on computer vision* 740–
726 755 (Springer, 2014).
- 727 39.Dutta, A. & Zisserman, A. The VIA Annotation Software for Images, Audio and Video. in *Proceedings of the*
728 *27th ACM International Conference on Multimedia* (ACM, 2019). doi:10.1145/3343031.3350535.
- 729 40.Xiao, B., Wu, H. & Wei, Y. Simple Baselines for Human Pose Estimation and Tracking. *arXiv:1804.06208 [cs]*
730 (2018).
- 731 41.Tan, M. & Le, Q. V. EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks.
732 *arXiv:1905.11946 [cs, stat]* (2020).
- 733 42.Krizhevsky, A., Sutskever, I. & Hinton, G. E. Imagenet classification with deep convolutional neural networks.
734 in *Advances in neural information processing systems* 1097–1105 (2012).
- 735 43.Vidal, M., Wolf, N., Rosenberg, B., Harris, B. P. & Mathis, A. Perspectives on Individual Animal Identification
736 from Biology and Computer Vision. *Integrative and Comparative Biology* **61**, 900–916 (2021).
- 737 44.Demšar, J. Statistical comparisons of classifiers over multiple data sets. *Journal of Machine learning research*
738 **7**, 1–30 (2006).

- 739 45.Tenenbaum, J. B. A Global Geometric Framework for Nonlinear Dimensionality Reduction. *Science* **290**, 2319–
740 2323 (2000).
- 741 46.Lin, T.-Y. *et al.* Feature Pyramid Networks for Object Detection. in *2017 IEEE Conference on Computer Vision
742 and Pattern Recognition (CVPR)* 936–944 (IEEE, 2017). doi:10.1109/CVPR.2017.106.
- 743 47.Girshick, R. Fast R-CNN. in *2015 IEEE International Conference on Computer Vision (ICCV)* 1440–1448 (2015).
744 doi:10.1109/ICCV.2015.169.
- 745 48.Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I. & Salakhutdinov, R. Dropout: A Simple Way to Prevent
746 Neural Networks from Overfitting. *30*.
- 747 49.Ioffe, S. & Szegedy, C. Batch Normalization: Accelerating Deep Network Training by Reducing Internal
748 Covariate Shift. *arXiv:1502.03167 [cs]* (2015).
- 749 50.Maas, A. L., Hannun, A. Y. & Ng, A. Y. Rectifier Nonlinearities Improve Neural Network Acoustic Models. *6*.
- 750 51.Xu, B., Wang, N., Chen, T. & Li, M. Empirical Evaluation of Rectified Activations in Convolutional Network.
751 *arXiv:1505.00853 [cs, stat]* (2015).
- 752 52.Kingma, D. P. & Ba, J. Adam: A Method for Stochastic Optimization. *arXiv:1412.6980 [cs]* (2017).
- 753 53.Lin, T.-Y., Goyal, P., Girshick, R., He, K. & Dollár, P. Focal Loss for Dense Object Detection. *arXiv:1708.02002
754 [cs]* (2018).
- 755 54.Bohnslav, J. P. *et al.* *DeepEthogram: a machine learning pipeline for supervised behavior classification from
756 raw pixels.* <http://biorxiv.org/lookup/doi/10.1101/2020.09.24.312504> (2020)
757 doi:10.1101/2020.09.24.312504.
- 758 55.Abadi, M. *et al.* TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems. *19*.
- 759 56.Chollet, F. *Keras.* (2015), GitHub repository: <https://github.com/fchollet/keras> (2015)..
- 760
- 761

762 **Supplementary**

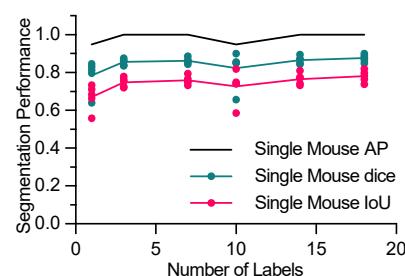
763



764

765 **Supplementary Fig. 1 | Segmentation annotation illustration.** An exemplary frame of mice
766 in OFT with manually annotated outlines.

767



768

769 **Supplementary Fig. 2 | Mouse single segmentation.** For mice, SIPEC:SegNet performance
770 in mAP, dice and IoU for single mouse as a function of the number of labels. The lines indicate
771 the means for 5-fold CV while circles, squares, triangles indicate the mAP, dice, and IoU,
772 respectively, for individual folds. All data is represented by mean, showing all points.

773

774

775

776

777

778

779



780
781
782
783
784

Supplementary Fig. 3 | Identification Graphical User Interface. Mask-box results from SIPEC:SegNet is overlaid over frames in blue and can be labeled one by one. The current box to be labeled is in green. A simple keyboard input scheme is provided within the GUI. Names of individuals and the number of masks to be labeled can be set by the user.

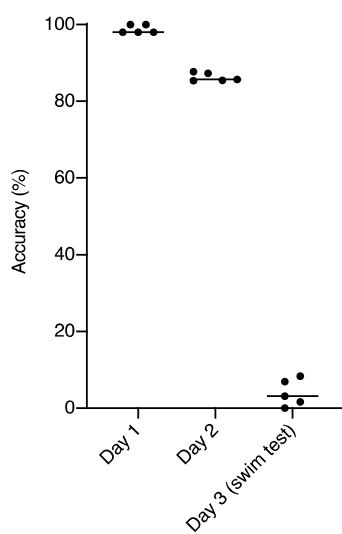
785
786



787
788 **Supplementary Fig. 4 | Example frames of the 8 distinct mice.**
789
790
791

792

793



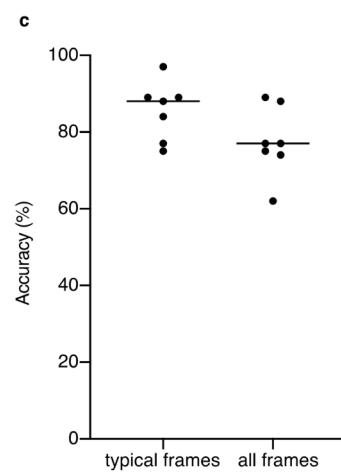
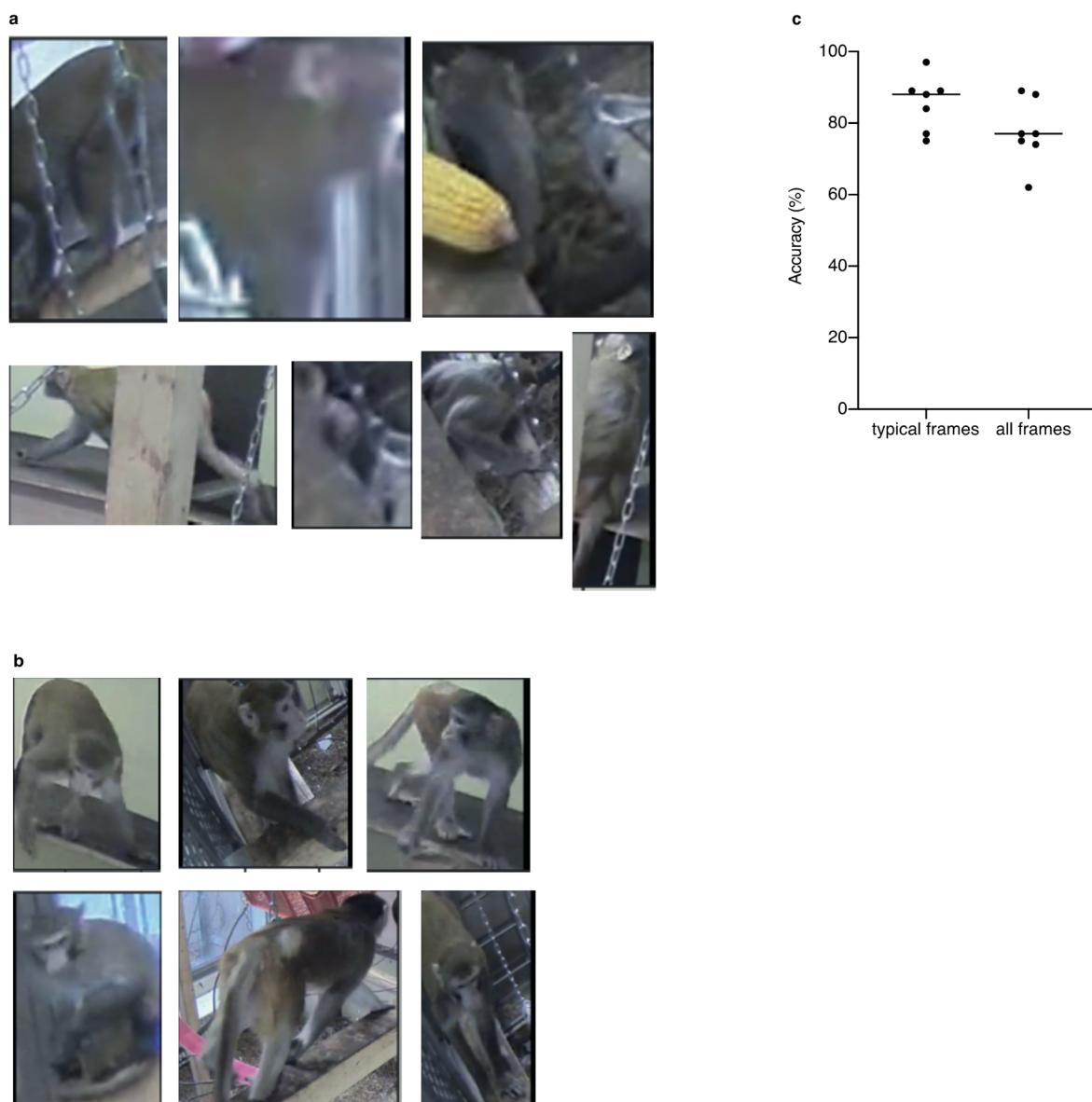
794

795 Supplementary Fig. 5 | Identification performance of mice across days and interventions.

796 Identification accuracy across days for models trained on day 1. While the performance for the
797 day the model is trained on is very high it drops when tested on day 2, but is still significantly
798 above chance level. When tested on day 3, after a forced swim test intervention, the
799 performance drops significantly. All data is represented by mean, showing all points.

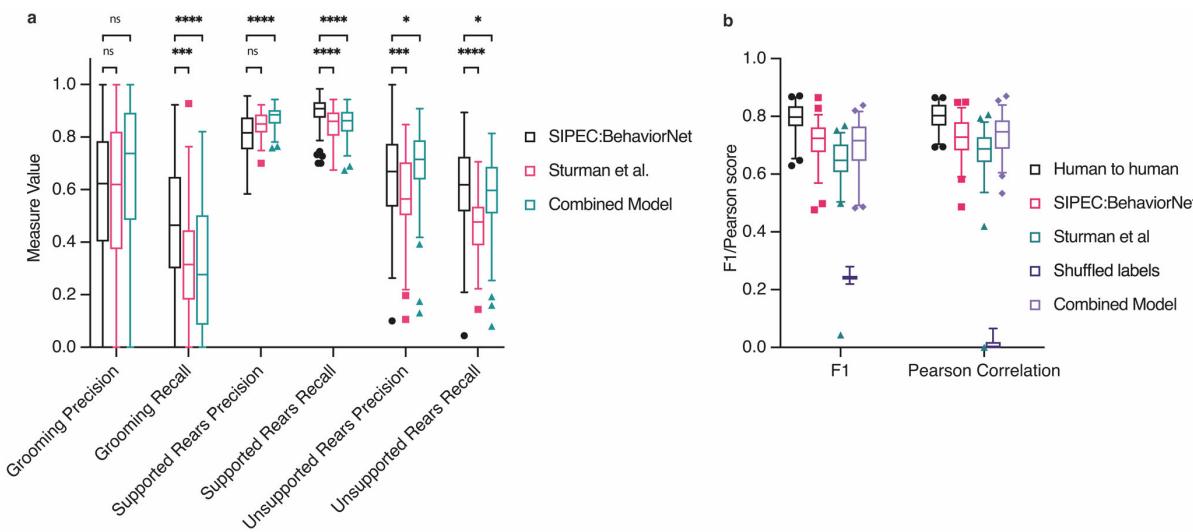
800

801



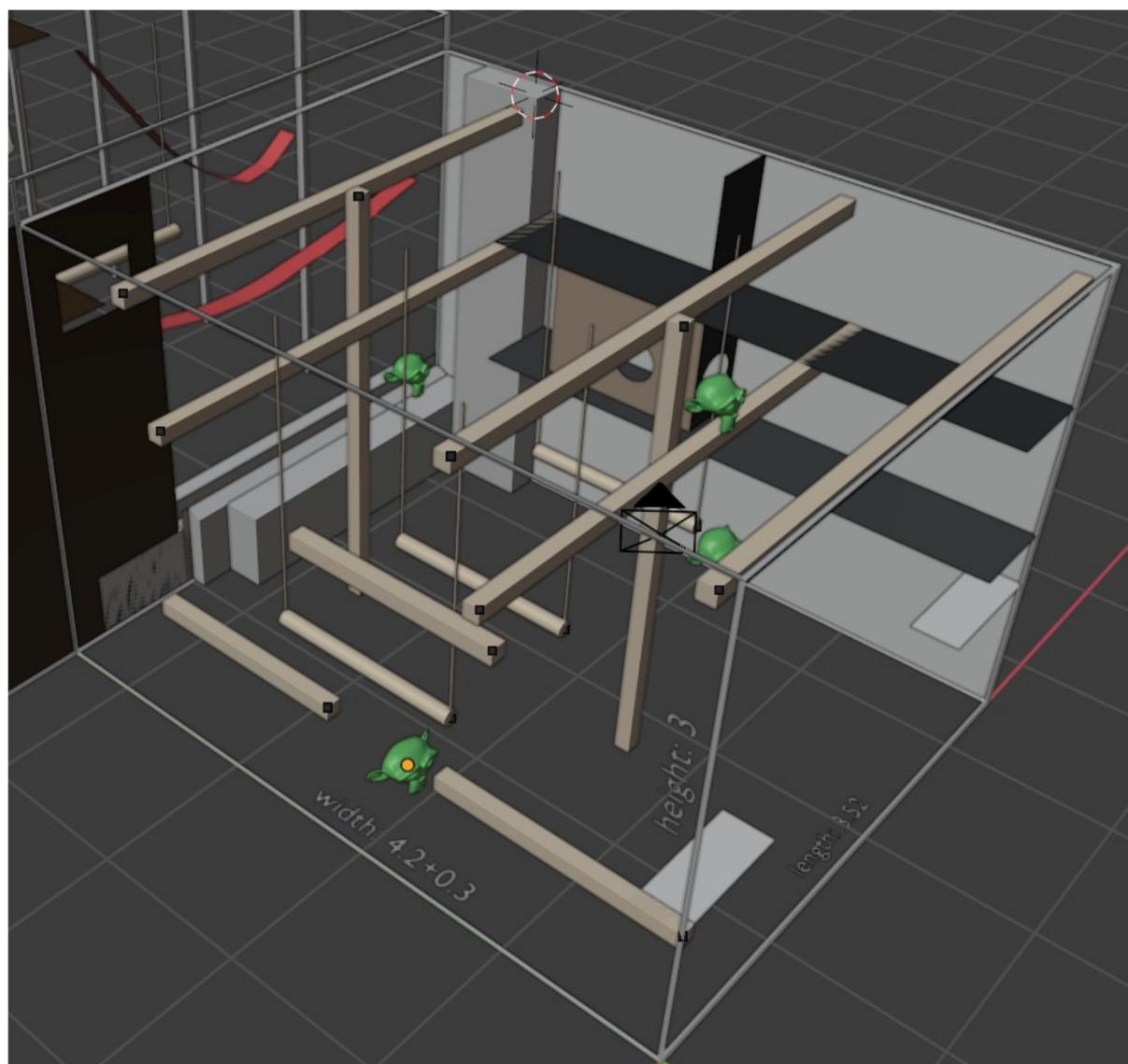
802
803
804
805
806
807
808
809
810
811
812
813
814

Supplementary Fig. 6 | Identification of typical vs difficult frames. a) Very difficult exemplary frames, which are also beyond human single-frame recognition, are excluded for the ‘typical’ frame evaluation. b) Exemplary frames used for the ‘typical’ frame analysis. c) Identification performance is significantly higher on ‘typical’ frames than on all frames. All data is represented by mean, showing all points.



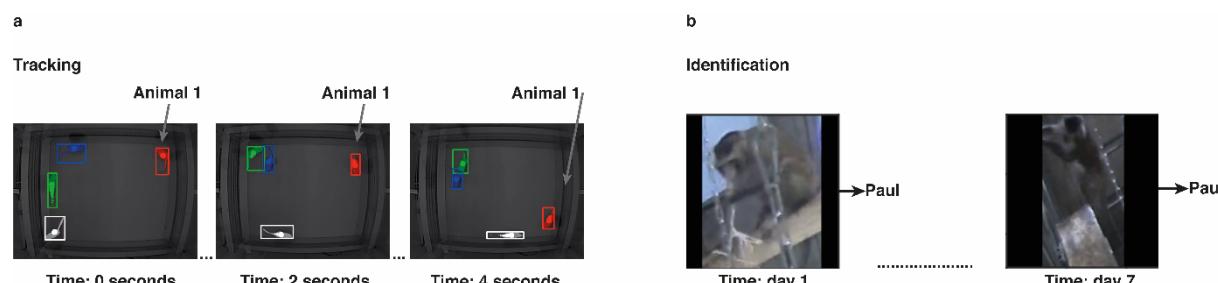
815
816
817
818
819
820
821
822
823
824
825
826

Supplementary Fig. 7 | Additional behavioral evaluation. a) Overall increased F1 score is caused by an increased recall in case of grooming events and precision for unsupported rearing events. b) Comparison of F1 values as well as Pearson Correlation of SIPEC:BehaveNet to human-to-human performance as well as combined model. Using pose estimates in conjunction with raw-pixel classification increases precision in comparison with solely raw-pixel classification while suffering from a decrease in recall. All data is represented by a Tukey box-and-whisker plot, showing all points.



827
828
829
830

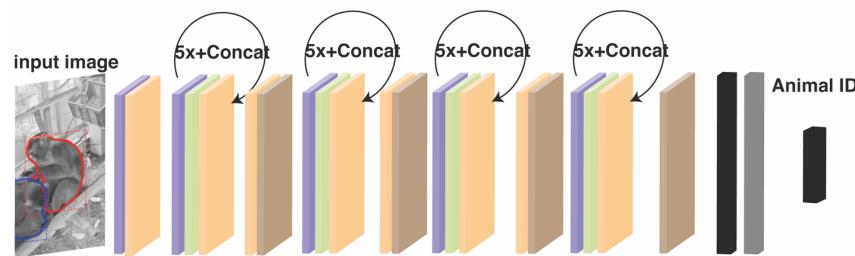
Supplementary Fig. 8 | 3D model used for annotation of primate 3D-location data.



831
832
833
834
835
836
837
838

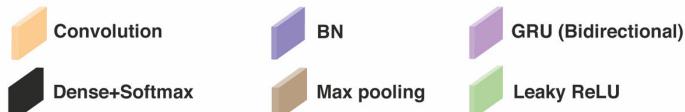
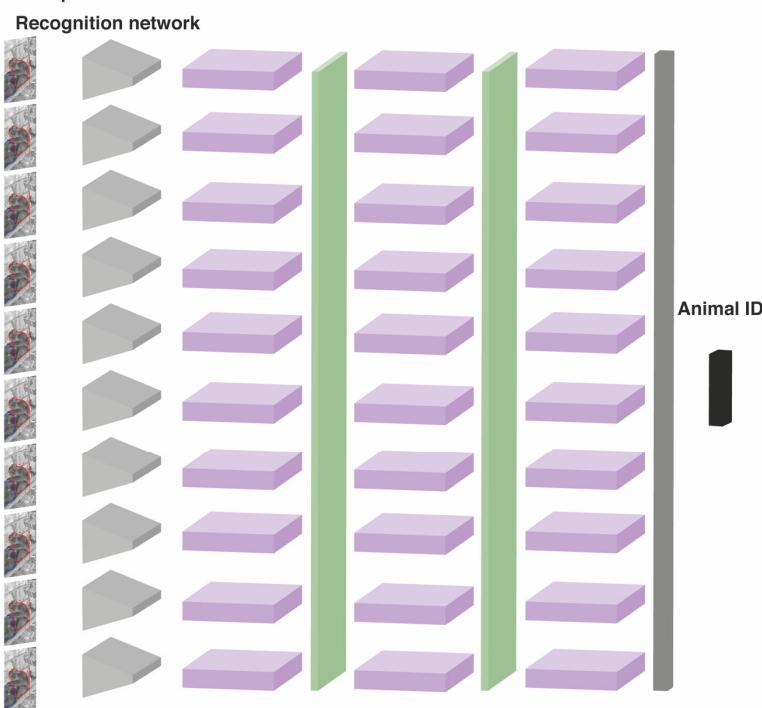
Supplementary Fig. 9 | Comparison tracking and identification. a) **Tracking** describes the process of following each individual animal in a group of animals within one session in a given field of view. b) **Identification** describes the ability to identify an individual from a single frame or a few consecutive frames across multiple sessions that could be apart hours, days, or months. This entails difficulties such as varying lighting conditions, occlusions, changes in the appearance of animals over time.

a) SIPEC:IdNet, single frame



b) SIPEC:IdNet, recurrent

temporal sequence



839

840

Supplementary Fig. 10 | SIPEC IdNet Architecture.

841

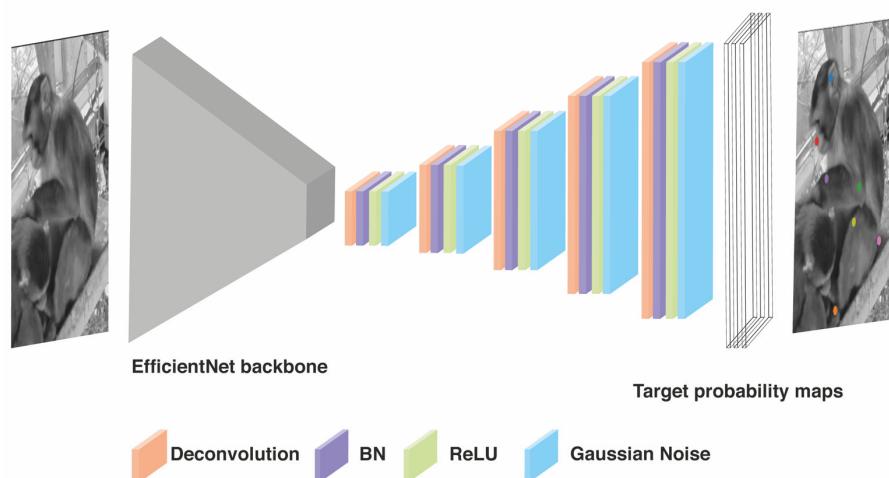
842

843

844

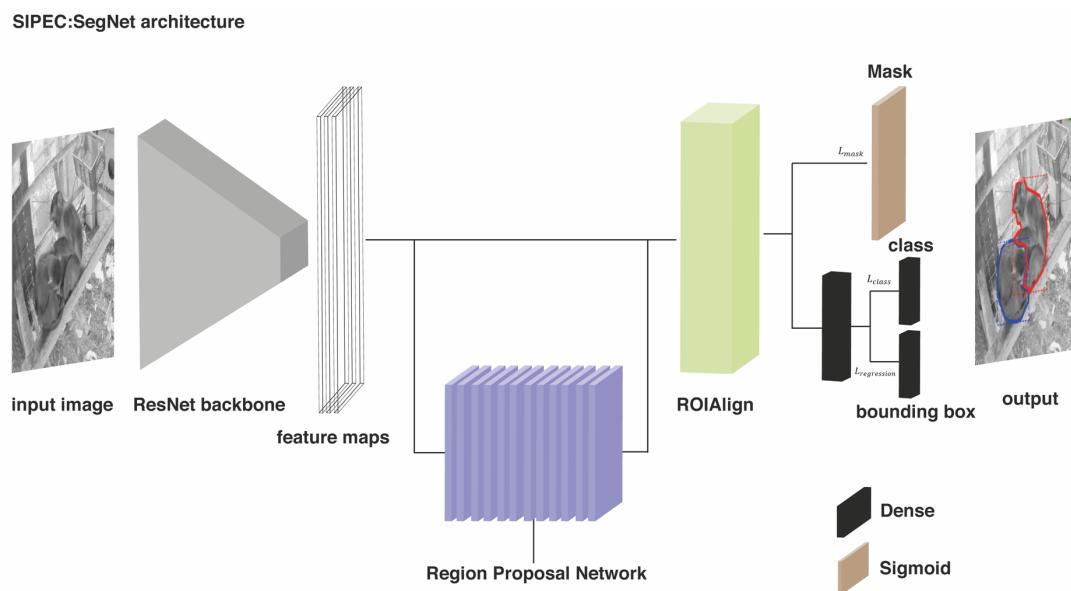
845

SIPEC:PoseNet architecture



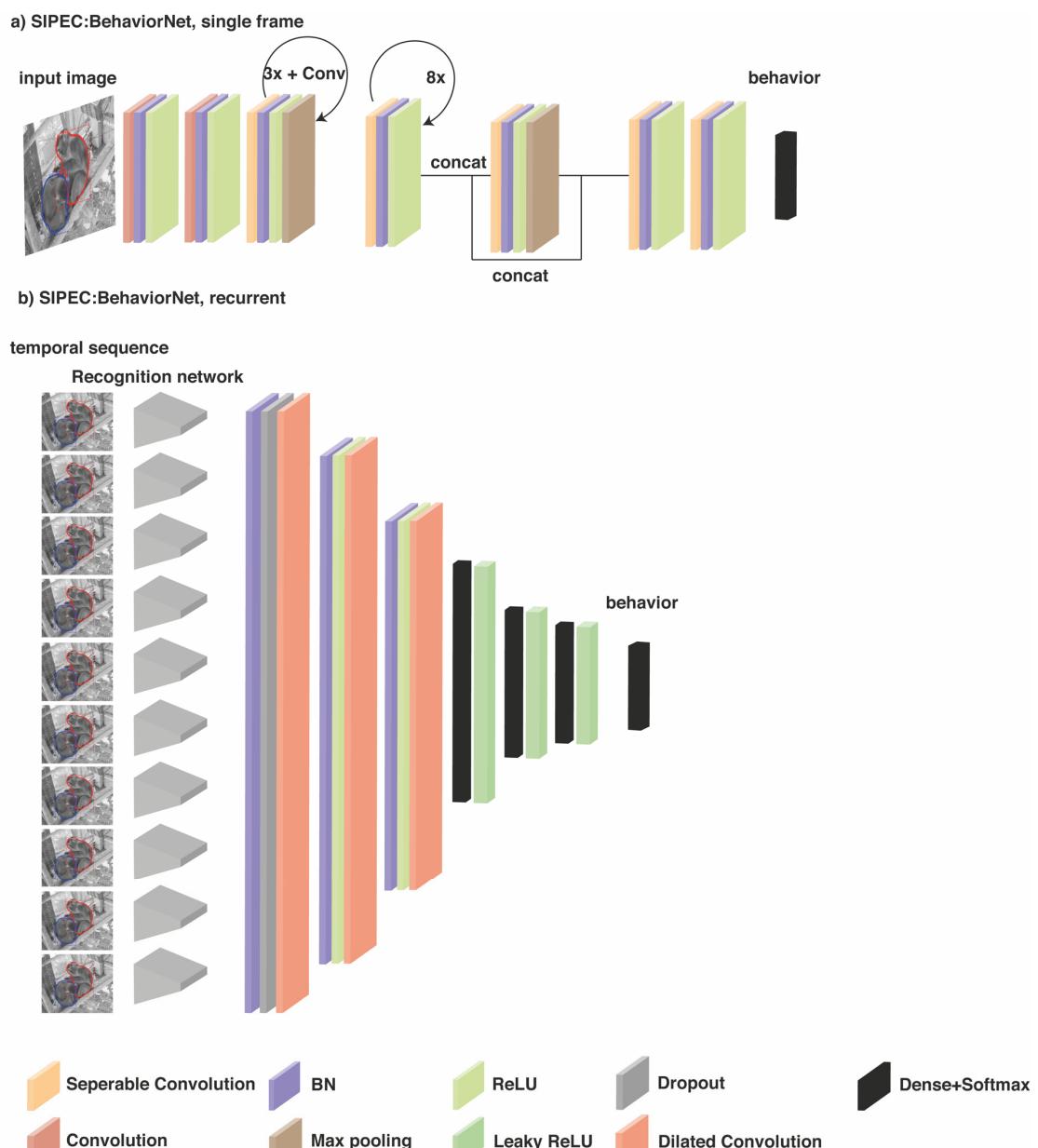
846
847
848

Supplementary Fig. 11 | SIPEC PoseNet Architecture.



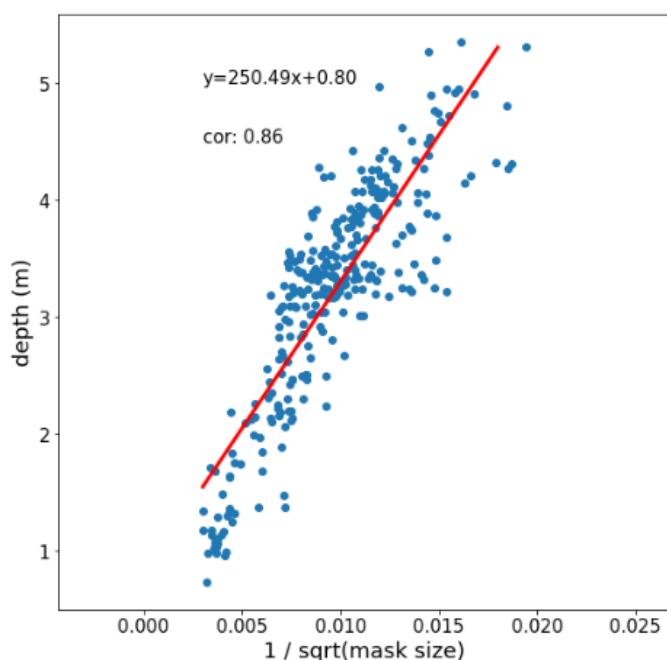
849
850
851
852
853

Supplementary Fig. 12 | SIPEC SegNet Architecture.



854
855

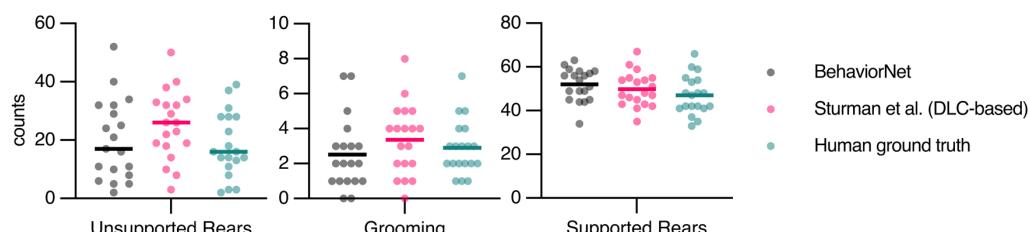
Supplementary Fig. 13 | SIPEC BehaviorNet Architecture.



856
857
858
859
860
861

Supplementary Fig. 14 | 3D depth estimates based on mask size. The inverse of the square root of the mask size (based on SIPEC:SegNet output) highly correlates with the depth of the individual in 3D space.

862
863
864
865
866
867
868
869
870
871
872



Supplementary Fig. 15 | Comparison of counts of behaviors between SIPEC:BehaviorNet, pose estimation based approach and human raters. Unsupported and supported rears and grooming events were counted per video for n=20 different mice videos. Behaviors were integrated over multiple frames, as described in Sturman et al.¹⁴. Behavioral counts of 3 different human expert annotators were averaged (in legend as ‘human ground truth’). No significant differences were found for comparing the number of behaviors between SIPEC:BehaviorNet and human annotators or Sturman et al.¹⁴ and human annotators (Tukey’s multiple comparison test). All data is represented by mean, showing all points.

Species	Network	Training (seconds/epoch)	Epoch s	Total training (min)	Inference (seconds/frame)
Primate	SegNet	133	100	222	0.4
Primate	IdNet(single frame)	134	10	22	0.35
Primate	IdNet(recurrent)	60	20	20	0.6
Primate	PoseNet	34	900	510	0.07

Primate	BehaveNet(single frame)	15	50	5	0.08
Primate	BehaveNet(recurrent)	190	10	31	0.6
Mouse	SegNet	40	100	67	0.14
Mouse(@1000 frames)	IdNet(single frame)	432	10	72	0.1
Mouse(@1000 frames)	IdNet(recurrent)	502	10	83	0.19
Mouse	PoseNet	20	2450	817	0.03
Mouse	BehaveNet(single frame)	547	10	91	0.05
Mouse	BehaveNet(recurrent)	1163	10	194	0.2

873 **Supplementary Tab. 1 | Training and inference times**

874 All measures are done with an NVIDIA RTX 2080 Ti and represent average values.

875

876

877 **Supplementary Video 1 | Illustration of SIPEC:SegNet and SIPEC:IdNet in primate**

878 homecage environment.

879 Short exemplary video of behaving primates in their homecage environment. SIPEC:SegNet is
880 used to mask different primates and SIPEC:IdNet is used to identify them. During obstructions,
881 the identity of a primate can alter but SIPEC:IdNet quickly recovers the correct identity over
882 the next frames, as it becomes more visible and therefore better identifiable.

883

884 **Supplementary Video 2 | Comparison of SIPEC and idtracker.ai for mice.**

885 Comparison for tracking 4 mice by idtracker.ai (Left) and by SIPEC(Right). We used publicly
886 available data from idtracker.ai (<https://drive.google.com/drive/folders/1Vua7zd6VuH6jc-NAd1U5iey4wU5bNrm4>) as well as idtracker.ai's publicly available inference results
887 (<https://www.youtube.com/watch?v=ANsThSPgBFM>) for a tracking comparison. **Left video:**
888 The tracking of idtracker.ai exhibits prolonged label switching errors where the label of two or
889 more animals gets swapped for some time. **Right Video:** Tracking is performed by
890 SIPEC:SegNet in conjunction with greedy-mask matching to track the identities of animals. In
891 this example video, SIPEC is more robust to these kinds of errors than idtracker.ai. (see also
892 Supp. Video 4).

893

894 **Supplementary Video 3 | Tracking of 4 mice by SIPEC in an open-field test.**

895 The masks generated by SIPEC:SegNet in conjunction with greedy-mask matching are used to
896 robustly track identities of four mice in an open-field test (see Methods).

897

898 **Supplementary Video 4 | SIPEC tracking over 52-minute video.**

899 We used publicly available data from idtracker.ai
900 (<https://drive.google.com/drive/folders/1Vua7zd6VuH6jc-NAd1U5iey4wU5bNrm4>) and
901 tracked 4 mice. The masks generated by SIPEC:SegNet in conjunction with greedy-mask
902 matching are used to robustly track identities of four mice in an open-field test (see Methods).

903