

1

2

3 **Robust mouse tracking in complex environments**

4

using neural networks

5

6

7

8 Brian Q. Geuther, Sean P. Deats, Kai J. Fox, Steve A. Murray, Robert E. Braun,

9 Jacqueline K. White, Elissa J. Chesler, Cathleen M. Lutz, Vivek Kumar*

10

11

12

13

14

15 The Jackson Laboratory, Bar Harbor, ME, United States of America

16

17 * Corresponding author

18 Email: vivek.kumar@jax.org (VK)

19

20 **Abstract**

21 The ability to track animals accurately is critical for behavioral experiments. For video-based
22 assays, this is often accomplished by manipulating environmental conditions to increase
23 contrast between the animal and the background, in order to achieve proper
24 foreground/background detection (segmentation). However, as behavioral paradigms become
25 more sophisticated with ethologically relevant environments, the approach of modifying
26 environmental conditions offers diminishing returns, particularly for scalable experiments.
27 Currently, there is a need for methods to monitor behaviors over long periods of time, under
28 dynamic environmental conditions, and in animals that are genetically and behaviorally
29 heterogeneous. To address this need, we developed a state-of-the-art neural network-based
30 tracker for mice, using modern machine vision techniques. We test three different neural
31 network architectures to determine their performance on genetically diverse mice under varying
32 environmental conditions. We find that an encoder-decoder segmentation neural network
33 achieves high accuracy and speed with minimal training data. Furthermore, we provide a
34 labeling interface, labeled training data, tuned hyperparameters, and a pre-trained network for
35 the mouse behavior and neuroscience communities. This general-purpose neural network
36 tracker can be easily extended to other experimental paradigms and even to other animals,
37 through transfer learning, thus providing a robust, generalizable solution for biobehavioral
38 research.

39

40 **Author summary**

41 Accurate tracking of animals is critical for behavioral experiments, however tracking in complex
42 environments has been a long-standing issue in neurogenetics. If the environment changes
43 during the test or if occlusion occurs, then tracking using existing methods often fails. These
44 technological constraints limit the complexity of behavioral paradigms that can be carried out.
45 Here we use modern convolutional neural networks to overcome these limitations and design a
46 trainable mouse tracker for complex and dynamic environments. We test several neural network
47 architectures and show that a single trained network can track all strains of mice we have tested
48 consisting of various coat colors, body shapes, and behaviors. We provide a labeling interface,
49 labeled training data, tuned hyperparameters, and a pre-trained network for the mouse behavior
50 and neuroscience communities.

51 **Introduction**

52 Behavior is primarily an output of the nervous system in response to internal or external
53 stimuli. It is hierarchical, dynamic, and high dimensional, and is generally simplified for analysis
54 [1, 2]. For instance, the rich locomotor movement performed by a mouse that is captured in
55 video is routinely abstracted to either a simple point, a center of mass, or an ellipse for analysis.
56 In order to do this well with current methods, the experimental environment is simplified to obtain
57 optimal contrast between the mouse and background for proper segmentation. Segmentation, a
58 form of background subtraction, classifies pixels belonging to mice from background in video
59 and enables these high level abstractions to be mathematically calculated. During mouse
60 experimental assays, the arena background color is often changed depending on the animal's
61 coat color, potentially altering the behavior itself [3-5]. Making such changes comes at a cost, as
62 current video tracking technologies cannot be applied in complex and dynamic environments or
63 with genetically heterogeneous animals without a high level of user involvement, making both
64 long term experiments and large experiments unfeasible. As neuroscience and behavior moves
65 into an era of big behavioral data [2] and computational ethology [6], current tracking methods
66 are inadequate and improved methods are necessary that enable tracking animals in semi-
67 natural and dynamic environments over long periods of time. To address this shortfall, we
68 developed a robust scalable method of mouse tracking in an open field using modern
69 convolutional neural network architecture. Our trained neural network is capable of tracking all
70 commonly used strains of mice—including mice with different coat colors, body shapes, and
71 behaviors—under multiple experimental conditions without any user-involved adjustment of
72 tracking parameters. Thus we present a scalable and robust solution that allows tracking in
73 diverse experimental conditions.

74 **Results**

75 We first used existing tracking methods to track 59 different mouse strains in multiple
76 environments, and found them inadequate for our large-scale strain survey experiment (1,845
77 videos, 1,691 hours). Specifically, we tracked all the videos in this experiment using Ctrax [7], a
78 modern open-source tracking software package that uses background subtraction and blob
79 detection heuristics, and LimeLight (Actimetrics, Wilmette, IL), a commercially available tracking
80 software package that uses a proprietary tracking algorithm. Ctrax abstracts a mouse on a per
81 frame basis to five metrics: major and minor axis, x and y location of center of the mouse, and
82 the direction of the animal [7]. It utilizes the MOG2 background subtraction model, whereby the
83 software estimates both the mean and variation of the background of the video for use in
84 background subtraction. Ctrax uses the shape of the predicted foreground to fit ellipses.
85 LimeLight uses a single key-frame background model for segmentation and detection. Once a
86 mouse is detected using LimeLight, this software package abstracts the mouse to a center of
87 mass using a proprietary algorithm.

88 Our strain survey experiment includes videos of mice with different genetic backgrounds
89 causing expression of different coat colors including black, agouti, albino, grey, brown, nude,
90 and piebald (Fig. 1A, columns 1, 2, 3 and 4). We tracked all animals in the same open field
91 apparatus, which had a white background; this yielded good results for darker mice (black and
92 agouti mice), but poor results for lighter-colored (albino and grey mice) or piebald mice (Fig. 1A,
93 columns 1, 2, 3 and 4, S1 Video). Examples of ideal and actual tracking frames are shown for
94 the various coat colors (Fig. 1A, row 3 and 4 respectively).

95 **Fig. 1. Proposed solutions for our tracking problem.** (A) A representation of the environments
96 analyzed by existing approaches. A black mouse in a white open field achieves high foreground-
97 background contrast, and therefore actual tracking closely matches the ideal (column 1). Grey mice are
98 visually similar to the grey-colored arena walls and therefore often have their noses, which are grey,
99 removed while rearing on walls (column 2). Albino mice are visually very similar to the white arena floor
100 and are frequently not found during tracking (column 3). Piebald mice are broken in half by the tracking
101 software due to their patterned coat color (column 4). Placing a food cup, that is visually similar to the
102 mouse, into the arena causes tracking issues when the mouse climbs on top of the food cup (column 5).
103 Arenas with reflective surfaces also produce errors with tracking algorithms (column 6). (B) We identified

104 the cause of bad tracking as poor segmentation. However, testing a variety of difficult frames with multiple
105 background subtraction algorithms from the background subtraction library, we did not resolve this
106 segmentation issue. From top to bottom the background subtraction algorithms shown are: SuBSENSE,
107 Adaptive Median, Adaptive Background Learning, MultiCue BGS, and LOBSTER. (C) Our objective
108 tracking takes the form of an ellipse description of a mouse. For clarity, we show a cropped frame as input
109 into the networks, whereas the actual input is an unmarked full frame. (D) The structure of the
110 segmentation network architecture functions similarly to classical tracking approaches in which the
111 network predicts the segmentation mask for the mouse and then fits an ellipse to the predicted mask. (E)
112 The structure of the binned classification network architecture predicts a probability distribution of the
113 value for each ellipse-fit parameter, represented by the table where a max value is selected. Only three
114 parameters of the six ellipse-fit parameters are visually shown (X = center x-location, Major = major axis
115 length, Angle = direction of the mouse's nose). (F) The structure of the regression network architecture
116 directly predicts the 6 parameters used to describe an ellipse for tracking.

117 We also carried out video analysis of behavior in challenging environments including
118 both 24-hour experimental videos that added bedding and a food cup to our open field arena,
119 and videos from the open field experiment carried out as part of The Jackson Laboratory
120 KOMP2 (Knockout Mouse Phenotyping Project) [8] Phenotyping Center (Fig. 1A, column 5, 6,
121 respectively). In the 24-hour experiment, we collected data over multiple days in which mice
122 were housed in the open field with white paper bedding and food cup. The mice were kept in the
123 open field in this multiday data collection paradigm, and continuous recording was carried out in
124 light and dark conditions using an infrared light source. The bedding and food cups were moved
125 by the mouse and the imaging light source alternated between infrared and visible light over the
126 course of each day. The KOMP2 experiment uses a beam-break system in which mice are
127 placed in a clear acrylic arena with infrared beams on all sides. Since the floor of the arena is
128 clear acrylic, the surface of the table on which the arenas were placed shows through as dark
129 grey. In addition, one arena was placed on the junction between two tables, leaving the joint
130 visible. Further, the LED lights overhead caused a very high glare unique to each arena (S2
131 Video). This KOMP2 program has collected over five years of data using this system, and we
132 wanted to carry out video-based recording as an added analysis modality to detect gait affects
133 that cannot be identified by beam-break systems. Since environmental alterations could affect
134 the behavioral output and legacy data interpretation, we could not optimize or otherwise alter the
135 environment for video data collection. Instead, we simply added a camera on top of each arena.

136 Traditionally, contrast and reflection hurdles could be overcome by changing the environment
137 such that video data collection is optimized for analysis. For instance, to track albino mice, one
138 can increase contrast by changing the background color of the open field to black. However, the
139 color of the environment can affect the behavior of both mice and humans, and such
140 manipulations can potentially confound the experimental results [3, 4]. Regardless, such
141 solutions will not work for piebald mice in a standard open field, or any mice in either the 24-hour
142 data collection experiment or the KOMP2 arena.

143 We found that the combination of mouse coat colors and environments were difficult to
144 handle with Ctrax (S1 Video) and LimeLight. We optimized and fine-tuned Ctrax for each video
145 (Methods) in each of the three experiments and still found a significant number of frames with poor
146 tracking performance (Fig. 1A, row 4). Such optimization or tuning of background model was not
147 feasible with LimeLight. The frequency of poor tracking instances in an individual video
148 increased as the environment became less ideal for tracking. We discovered these errors in
149 ellipse fitting lead to larger errors in classifying behaviors using the Ctrax tracking output in
150 supervised classification using Janelia Automatic Animal Behavior Annotator (JAABA)[9]. Thus,
151 even the seemingly minor errors seen in grey and black mice (S1 Video) decreased
152 performance when the tracking data were used for behavior classification. Furthermore, the
153 distribution of the errors was not random; for example, tracking was highly inaccurate when mice
154 were in the corners, near walls, or on food cups (Fig. 1A, row 4), and less inaccurate when
155 animals were in the center (S1 Video). While it is feasible to discard poorly tracked frames, this
156 can lead to biased sampling and skewed biological interpretation.

157 We explored the cause of bad tracking across our three experiments and discovered
158 that, in most cases, improper tracking was due to poor segmentation of the mouse from the
159 background. This included both types of errors: Type I, instances when portions of the

160 background are included as the foreground (e.g. shadows), and Type II, instances when
161 portions of the mouse are removed from the foreground (e.g. albino mouse matching the
162 background color). Since Ctrax uses a single background model algorithm, we tested whether
163 other background model algorithms could improve tracking results. We tested 26 different
164 segmentation algorithms [10] and discovered that each of these traditional algorithms performs
165 well under certain circumstances and fail in others (Fig. 1B). Other available tracking software
166 packages including CADABRA [11], EthoVision [12], idTracker [13], MiceProfiler [14], MOTR
167 [15], Cleversys TopScan (<http://cleversysinc.com/CleverSysInc/>), Autotyping [16], and
168 Automated Rodent Tracker [17], all of which rely on background subtraction approaches for
169 tracking. Since all 26 background subtraction methods failed in some circumstances, we
170 postulate that our results for Ctrax and LimeLight will hold true for these other technologies. In
171 sum, although many video tracking solutions exist, none address the fundamental problem of
172 mouse segmentation appropriately and generally rely on environmental optimization to achieve
173 proper segmentation, therefore creating potential confounds with respect to robust data
174 sampling and analysis. Thus, we could not overcome the fundamental issue of proper mouse
175 segmentation in order to achieve high-fidelity mouse tracking with existing solutions.

176 A drawback in addition to the problem of inadequate mouse segmentation was the time
177 cost for fine-tuning Ctrax's settings or another background subtraction algorithm's parameters.
178 Fine-tuning the tracking settings for each video added significant time to our workflow when
179 analyzing thousands of videos. For example, in tracking data from the 24-hour experiment, when
180 mice were sleeping in one posture for an extended period of time, the mouse became part of the
181 background model and could not be tracked. Typical supervision, such as using the Ctrax
182 settings supervision protocol we outline in our methods, would take an experienced user 5
183 minutes of interaction for each hour of video to ensure high-quality tracking results. While this

184 level of user interaction is tractable for smaller and more restricted experiments, large-scale and
185 long-term experiments require a large time commitment to supervise the tracking performance.

186 We sought to overcome these difficulties by building a robust-next generation mouse
187 tracker that uses neural networks and achieves high performance under complex and dynamic
188 environmental conditions, is indifferent to coat color, and does not require persistent fine tuning
189 by the user. Convolutional neural networks are computational models that are composed of
190 multiple spatial processing layers that learn representations of data with multiple levels of
191 abstraction. These methods have dramatically improved the state-of-the-art in speech
192 recognition, visual object recognition, object detection, and many other domains such as drug
193 discovery and genomics [18]. One of the key advantages of neural networks is that once an
194 efficient network with suitable hyperparameters has been developed, it can easily be extended
195 to other tasks by simply adding appropriate training data [19]. Thus, we sought to build a highly
196 generalizable solution for mouse tracking.

197 We tested three primary neural network architectures for solving this visual tracking
198 problem (Fig. 1D-E). Each approach attempted to describe the location of the animal through six
199 variables: x and y location of the mouse in the matrix, major and minor axes of the mouse, and
200 the angle the head is facing (Fig. 1C). To avoid the discontinuity of equivalent repeating angles,
201 the networks predict the sine and cosine of the angle.

202 The first architecture is an encoder-decoder segmentation network that predicts a
203 foreground-background segmented image from a given input frame (Fig. 1D). This network
204 predicts on a pixel-wise basis whether there is a mouse or no mouse, with the output being a
205 segmentation mask. The segmentation mask identifies all the pixels in the image that belong to
206 the mouse. The primary structure of this architecture starts with a feature encoder, which
207 abstracts the input image down into a small-spatial-resolution set of features. The encoded

208 features are then passed to a feature decoder that converts this set of features back into the
209 same shape as the original input image. Additionally, the encoded features are also passed to
210 three fully connected layers to predict which cardinal direction the ellipse is facing. We trained
211 this feature decoder to produce a foreground-background segmented image. After the network
212 produces this segmented image, we applied an ellipse-fitting algorithm for tracking (Note A in S1
213 Information).

214 The second network architecture is a binned classification network, whereby a probability
215 distribution across a pre-defined range of possible values is predicted for each of the 6 ellipse-fit
216 parameters (Fig. 1E). This network architecture begins with a feature encoder that abstracts the
217 input image down into a small-spatial-resolution set of features. The encoded features are
218 flattened and connected to additional fully connected layers whose output shape is determined
219 by the desired resolution of the output. For instance, at a desired resolution of 1 pixel for the x-
220 coordinate location of the mouse, there are 480 possible x-values to select from for a 480 x 480
221 px image. As such, the network contains 480 values (bins) to select from, one bin for each x-
222 column in the 480 x 480 px image. When the network is run, the largest value in each heatmap
223 is selected as the most probable value of the corresponding parameter. Each desired output
224 parameter is realized as an independent set of trainable fully connected layers connected to the
225 encoded features.

226 The third architecture is a regression network that predicts the numerical ellipse values
227 directly from the input image (Fig. 1F). The network architecture begins with a feature encoder
228 that abstracts the input down into a small spatial resolution. These encoded features are then
229 flattened and connected to fully connected layers to produce an output shape of 6, the number
230 of values that we ask the network to predict to fit an ellipse. We tested a variety of currently
231 available general purpose feature encoders, and present data from the feature encoder Resnet

232 V2 [20] with 200 convolutional layers, which achieved the best performing results for this
233 architecture.

234 To test the neural network architectures, we built a training dataset of 16,234 training
235 images and 568 separate validation images across multiple mouse strains and experimental
236 setups (Note B in S1 Information). Annotated training images were augmented eightfold during
237 training by applying reflections. Additionally, training images were further augmented by adding
238 small random changes in contrast, brightness, and rotations to make the network robust to minor
239 fluctuations in input data. We created an OpenCV-based labeling interface for creating our
240 training data (Methods) that allows us to quickly label foreground and background, and fit an
241 ellipse (S1 Fig.). This labeling interface can be used to quickly generate annotated training data
242 in order to adapt any network to new experimental conditions through transfer learning.

243 Our network architectures were built, trained, and tested in Tensorflow v1.0, an open-
244 source software library for designing applications that use neural networks [21]. Training
245 benchmarks presented were conducted on the Nvidia P100 GPU architecture. We tuned the
246 hyperparameters through several training iterations. After the first training of networks, it was
247 observed that the networks performed poorly under particular circumstances that had not been
248 included in the annotated data, including mid-jump, odd postures, and urination in the arena. We
249 identified and incorporated these difficult frames into our training dataset to further improve
250 performance. A full description of the network architecture definitions and training
251 hyperparameters are available (Methods, Table A in S1 Information). Overall, training and
252 validation loss curves indicated that each of the three network architectures trains to a
253 performance with an average error between 1 and 2 pixels (Fig. 2A). The encoder-decoder
254 segmentation architecture converged to a validation error of 0.9px (Fig. 2 A, B, C). Surprisingly,
255 upon inspection of the validation curve for the binned classification network we found that it

256 displayed unstable loss curves, indicating overfitting and poor generalization (Fig. 2B, E). The
257 regression architecture converged to a validation error of 1.2 px, showing a better training than
258 validation performance (Fig. 2A, B, D).

259 **Fig. 2. Neural network performance metrics.** (A-E) Performance of our tested network architectures
260 during trainings. (A) Training curves show comparable performances of the three architectures during
261 training, independent of the network architecture. (B) Validation curves show different performances
262 across the three network architectures. The encoder-decoder segmentation network performs the best.
263 (C, D, E) Comparison of training and validation performance curves, by network architecture type. (C)
264 Performance increases for validation in our encoder-decoder segmentation network architecture. (D)
265 Performance decreases for validation in our regression network architecture, but a good generalization
266 performance is maintained by asymptotically converging to a value. (E) The binned classification network
267 architecture becomes unstable at 55 epochs of training, even though the training curve shows continued
268 improved performance at this timepoint. (F) Comparing our encoder-decoder segmentation network
269 architecture with a beam break system, we observe a high correlation. Each point represents an individual
270 video tracked using both our neural network and a beam break system. Our network performs
271 consistently, even though the arenas are visually different from one another. We identify two videos of
272 individual mice that deviate from this trend (red arrows). (G) Predictions from two approaches yield
273 high agreement on environments with high contrast between the mouse and background (black, grey, and
274 piebald mice in the white background open-field assay). As the segmentation problem becomes more
275 computationally difficult, the relative error increases (albino mice in the white background open-field assay
276 black mice in the 24-hr assay, KOMP2 experiment). Due to low activity in the 24-hr setup, minor errors in
277 tracking have a large influence on measurements of the total distance traveled. Points indicate individuals
278 in a group, bars indicate mean \pm standard deviation. (H) Relative standard deviation of the minor axis
279 maintains a high correlation when the mouse and environment have a high contrast (black mice in the
280 white background open-field assay). When segmentation includes shadows, includes reflections, or
281 removes portions of the mouse, the minor axis length is not properly predicted and increases the relative
282 standard deviation (grey, piebald, and albino mice in the white background open-field assay, black mice in
283 the 24-hour assay, KOMP2 experiment). Points indicate individuals in a group, bars indicate mean \pm -
284 standard deviation.

285 Not only does the encoder-decoder segmentation architecture perform well, but it also is
286 computationally efficient for GPU compute, requiring an average processing time of 5-6ms per
287 frame. With the encoder-decoder segmentation architecture, our video data could be processed
288 at a rate of up to 200 frames per second (fps) (6.7X realtime) on a Nvidia P100, which is a
289 server-grade GPU.; and a rate of up to 125 fps (4.2X realtime) on a Nvidia TitanXP, a consumer-
290 grade GPU. This high processing speed is likely due to the structure of the encoder-decoder
291 segmentation architecture, as it is only 18 layers deep and contains only 10.6 million trainable
292 parameters. In comparison, Resnet V2 200, the feature extractor that gave the best results for
293 the regression architecture, is a large and deep network with over 200 layers and 62.7 million

294 trainable parameters and leads to a substantially longer processing time per frame (33.6ms on a
295 Nvidia P100). Other pre-built general-purpose networks [22] achieve similar or worse
296 performances at a tradeoff of faster compute time. Thus, regression networks are an accurate
297 but computationally expensive solution.

298 We also tested the minimum training dataset size required to train the encoder-decoder
299 segmentation network, by randomly subsetting our training dataset to smaller numbers of
300 annotated images (10,000 to 500) and training the network from the beginning. Surprisingly, we
301 obtained good results from a network trained with only 2,500 annotated images, a task that
302 takes approximately three hours to generate with our labeling interface (S2 Fig.). Given the
303 computational efficiency, accuracy, and training stability of the encoder-decoder segmentation
304 architecture, and the small training dataset size that it requires, we concluded that this
305 architecture is optimal for our needs. We used this trained neural network to predict the location
306 of mice for entire videos and compare tracking performance with other non-neural network
307 approaches including a beam-break system (KOMP2) and a video tracking system (Ctrax).

308 We evaluated the quality of the encoder-decoder segmentation neural network tracking
309 architecture by inferring entire videos from mice with disparate coat colors and data collection
310 environments (Fig. 1A) and visually evaluating the quality of the tracking. We also compared this
311 neural network-based tracking architecture with an independent modality of tracking, the
312 KOMP2 beam-break system (Fig. 1A, column 6). We tracked 2,002 videos of individual mice
313 comprising 700 hours of video from the KOMP2 experiment using the encoder-decoder
314 segmentation neural network architecture and compared the results with the tracking data
315 obtained using the KOMP2 beam-break system (Fig. 2F). These data comprised mice of 232
316 knockout lines on the C57BL/6NJ background that were tested in 20-minute open field assay in
317 2016 and 2017. Since each KOMP2 arena has slightly different background due to the

318 transparent and reflective walls, we compared tracking performances of the two approaches for
319 each of the eight testing arenas used in the 2016 and 2017 KOMP2 open-field assays (Fig. 2F,
320 colors shows arena), and compared tracking performances for all the arenas combined (Fig. 2F,
321 black line). We observed a very high correlation between the total distance traveled in the open
322 field as measured by the two approaches across all eight KOMP2 testing arenas ($R = 96.9\%$,
323 Fig. 2F). We observed two animals with high discordance from this trend (Fig. 2F, red arrows).
324 Observation of the video showed odd behaviors for both animals, with a waddle gait in one and
325 a hunched posture in the other (S2 Video). We postulate that these behaviors led to abnormal
326 beams breaks causing erroneously high total distances traveled measured via the beam break
327 system. This example highlights an important advantage of the neural network, as it is
328 unaffected by the behavior of the animal.

329 We then compared the performance of our trained segmentation neural network with the
330 performance of Ctrax across a broad selection of videos from the various testing environments
331 and coat colors previously tracked using Ctrax and LimeLight (Fig. 1A). We wish to emphasize
332 that we compared the performance of our network with that of Ctrax because Ctrax is one of the
333 best conventional tracking software packages that allows fine tuning of the many tracking
334 settings, is open source, and provides user support. Given the results with the 26 background
335 subtraction approaches (Fig. 1B), we expected similar or worse performances from other
336 tracking systems. We tracked 72 videos, broken into 6 groups (Fig. 1A) with 12 animals per
337 group, with both our trained encoder-decoder segmentation neural network and Ctrax. The
338 settings for Ctrax were fine-tuned for each of the 72 videos, as described in 'Ctrax Settings
339 Supervision Protocol' in Methods. Videos from the 24-hr experiment showing that animals that
340 were sleeping continually for the full video duration (one hour) were manually omitted from
341 comparison, as Ctrax will incorporate the mouse as part of the background model. We

342 calculated a cumulative relative error of total distance traveled between Ctrax and our neural
343 network (Fig. 2G). Specifically, for every minute in the video, we compared the distance-traveled
344 prediction of the neural network with that of Ctrax. This metric measures the accuracy of center
345 of mass tracking of each mouse. Tracking for black, gray, and piebald mice in the white-
346 background open-field apparatus showed errors less than 4%; however, significantly higher
347 levels of error were seen in albino mice in the open-field arena with a white floor (14%), black
348 mice in the 24-hour arena (27%), and black mice in the KOMP2 testing arena (10%) (Fig. 2G
349 and S1 Video). Thus, we could not adequately track albino mice in the open-field arena with a
350 white floor, black mice in the 24-hour arena, or black mice in the KOMP2 testing arena without
351 the neural network tracker.

352 We also observed, using Ctrax, that when foreground segmentation prediction is
353 incorrect, such as when shadows are included in the prediction, the ellipse fit does not correctly
354 represent the posture of the mouse (S1 Video). In these cases, even though the center of mass
355 tracking was acceptable, the ellipse fit itself was highly variable. Modern machine learning
356 software for behavior recognition, such as the Janelia Automatic Animal Behavior Annotator
357 (JAABA)[9], utilize the time series of ellipse fit tracking for classification of behaviors. We
358 quantitated the stability of ellipse tracking through measuring the relative standard deviation of
359 the minor axis and comparing approaches. This metric shows the least variance across all sizes
360 of laboratory mice, as the width of an individual mouse remains similar through a wide range of
361 postures expressed in behavioral assays when tracking is accurate. We observed a high level of
362 tracking variation with grey and piebald mice in the white open field arena (Fig. 2H) even though
363 there is low cumulative relative error of total distance traveled (Fig. 2G). As expected, we
364 observed a high relative standard deviation of the minor axis for albino mice (white open field
365 arena) and KOMP2 tracking. Thus, for both center of mass tracking and variance of ellipse fit we

366 find that the neural network tracker outperforms traditional background subtraction-based
367 trackers.

368 Having established the encoder-decoder segmentation neural network as a highly
369 accurate tracker, we tested its performance using two large behavioral experiments. For the first
370 experiment, we generated white-surfaced open-field video data with 1,845 mice, including 58
371 strains of mice including mice with diverse coat colors, piebald mice, nude mice, and obese
372 mice; and covering a total of 1,691 hours (Fig. 3A). This dataset consists of 47 inbred strains
373 and 11 isogenic F1 strains and is the largest open-field dataset generated, based on the data in
374 the Mouse Phenome Database[23]. Using a single trained network without any user tuning, we
375 were able to track all mice with high accuracy. We visually checked mice from a majority of the
376 strains for fidelity of tracking and observed excellent performance. The activity phenotypes that
377 we observed agree with previously published datasets of mouse open-field behavior[23]. For the
378 second dataset, we tracked 24-hour video data collected for four C57BL/6J and two BTBR T⁺
379 Itpr3^{fl}/J mice (Fig. 1A, column 5). These mice were housed with bedding and a food cup over
380 multiple days during which the food changed location and under 12:12 light-dark conditions.
381 Video data were recoded using visible and infrared light sources. We tracked activity across all
382 animals under these conditions using the same encoder-decoder segmentation neural network
383 architecture used for the first experiment, and observed very good performance under light and
384 dark conditions (Fig. 3B, light and dark blue points, respectively). As expected, we observed
385 daily activity rhythm with high levels of locomotor activity during the dark phase (Fig. 3B, red
386 curve).

387 **Fig. 3. Highly scalable tracking with a single neural network.** (A) A large strain survey showing
388 genetically diverse animals traced with our encoder-decoder segmentation network. 1,845 animals
389 including 58 inbred and F1 isogenic strains, totaling 1,691 hours of video, were processed by a single
390 trained neural network without any user-involved fine-tuning. Total distance traveled in a 55-minute open
391 field assay is shown. Points indicate individuals in a strain, bars indicate mean +/- standard deviation. Two
392 reference mouse strains are shown in bold, C57BL/6J and C57BL/6NJ (B) Daily activity rhythms were
393 observed in six animals continuously tracked over 4 days in a dynamic environment with our encoder-

394 decoder segmentation neural network. Points indicate distance traveled in an epoch. Red line indicates
395 polynomial fit showing daily activity rhythms.

396 **Discussion**

397 Video-based tracking of animals in complex environments has been a long-standing
398 challenge in the field of animal behavior [24]. Current state-of-the-art animal-tracking systems do
399 not address the fundamental issue of animal segmentation and rely heavily on visual contrast
400 between the foreground and background for accurate tracking. As a result, the user must restrict
401 the environment to achieve optimal results. Here we describe a modern neural network-based
402 tracker that is able to function in complex and dynamic environments. Our network addresses a
403 fundamental issue in tracking—foreground and background segmentation—by using a trainable
404 neural network. We test three different architectures and find that an encoder-decoder
405 segmentation network architecture achieves the highest level of accuracy and functions at a
406 high speed (over 6X real time). Furthermore, we provide a labeling interface that allows the user
407 to train a new network for their specific environment by labeling as few as 2,500 images, which
408 takes approximately 3 hours. We compare our network to two existing solutions and find that it
409 vastly outperforms them in complex environments. We expect similar results with any off-the-
410 shelf system that utilizes traditional background subtraction approaches. In fact, when we tested
411 26 different background subtraction methods we discovered that each failed under certain
412 circumstances. However, a single neural network architecture functions for all coat colors of
413 mice under multiple environments without the need for fine tuning or user input. Our machine
414 learning approach enables long-term tracking under dynamic environmental conditions with
415 minimal user input, thus establishing the basis of the next generation of tracking architecture for
416 behavioral research.

417 **Materials and methods**

418 **Experimental arenas**

419 **Open Field Arena**

420 Our open field arena measures 52cm by 52cm by 23cm. The floor is white PVC plastic
421 and the walls are grey PVC plastic. To aid in cleaning maintenance, a white 2.54cm chamfer
422 was added to all the inner edges. Illumination is provided by an LED ring light (Model: F&V
423 R300). The ring light was calibrated to produce 600 lux of light in each of our 24 arenas.

424 **24-Hour monitoring open field arena**

425 We augmented 6 of our open field arenas for multiple day testing. We set our overhead
426 LED lighting to a standard 12:12 light-dark cycle. ALPHA-dri was placed into the arena for
427 bedding. To provide food and water, a single Diet Gel 76A food cup was placed in the arena.
428 This nutritional source was monitored and replaced when depleted. Each arena was illuminated
429 at 250 lux during the day and <5 lux during the night. For recording videos during the night,
430 additional IR LED (940nm) lighting was added.

431 **KOMP2 open field arena**

432 In addition to our custom arenas, we also benchmarked our approach on a commercially
433 available system. The Accuscan Versamax Activity Monitoring Cages is constructed using clear
434 plastic walls. As such, visual tracking becomes very difficult due to the consequent reflections.
435 The cage measures 42cm by 42cm by 31cm. Lighting for this arena was via LED illumination at
436 100-200 lux.

437 **Video acquisition**

438 **Imaging hardware**

439 All data was acquired using the same imaging equipment. Data was acquired at

440 640x480px resolution, 8-bit monochrome depth, and 30fps using Sentech cameras (Model:
441 STC-MB33USB) and Computar lenses (Model: T3Z2910CS-IR). Exposure time and gain were
442 controlled digitally using a target brightness of 190/255. Aperture was adjusted to its widest so
443 that lower analog gains were used to achieve the target brightness. This in turn reduced
444 amplification of baseline noise. Files were saved temporarily on a local hard drive using the “raw
445 video” codec and “pal8” pixel format. Our typical assays run for two hours, yielding a raw video
446 file of approximately 50GB. Overnight, we use FFmpeg software (<https://www.ffmpeg.org/>) to
447 apply a 480x480px crop, de-noise filter, and compress using the mpeg4 codec (quality set to
448 max), which yields a compressed video size of approximately 600MB.

449 One camera and lens was mounted approximately 100cm above each arena to alleviate
450 perspective distortion. Zoom and focus were set manually to achieve a zoom of 8px/cm. This
451 resolution both minimizes the unused pixels on our arena border and yields approximately 800
452 pixels area per mouse. Although the KOMP2 arena is slightly smaller, the same zoom of 8px/cm
453 target was utilized.

454 **Ctrax settings supervision protocol**

455 Ctrax contains a variety of settings to enable optimization of tracking [7]. The authors of
456 this software strongly recommend, first and foremost, ensuring that the arena is set up under
457 specific criteria to ensure good tracking. In most of our tests, we intentionally use an
458 environment in which Ctrax is not designed to perform well (e.g., albino mice on a white
459 background). That being said, with well-tuned parameters, a good performance is still
460 achievable. However, with a large number of settings to manipulate, Ctrax can easily require
461 substantial time to achieve a good tracking performance. Here, we describe our protocol for
462 setting up Ctrax for tracking mice in our environments.

463 First, we create a background model. The core of Ctrax is based on background
464 subtraction, and thus a robust background model is essential for functionality. Models function
465 optimally when the mouse is moving. To create the background model, we seek to a segment of
466 the video in which the mouse is clearly moving, and we sample frames from that section. This
467 ensures that the mouse is not included in the background model. This approach significantly
468 improves Ctrax's tracking performance on our 24-hour data, as the mouse moves infrequently
469 due to sleeping and would typically be incorporated into the background model.

470 The second step is to set the settings for background subtraction. Here, we use the
471 Background Brightness normalization method with a Std Range of 254.9 to 255.0. The
472 thresholds applied to segment out the mouse are tuned on a per-video basis, as slight changes
473 in exposure and coat color will influence the performance. To fine-tune these thresholds, we
474 apply starting values based on previous videos analyzed and adjust values by checking multiple
475 portions of the video. Every video is inspected for proper segmentation on difficult frames, such
476 as the mouse rearing on the wall. Additionally, we apply morphological filtering to both remove
477 minor noise in the environment as well as remove the tails of mice for fitting an ellipse. We use
478 an opening radius of 4 and a closing radius of 5.

479 Lastly, we manually set a variety of tracking parameters that Ctrax enables to ensure that
480 the observations are in fact mice. For optimal time efficiency, these parameters were tuned well
481 once and then used for all other mice tracked. If a video was performing noticeably poorly, the
482 general settings were tweaked to improve performance. For the shape parameters, we
483 computed bounds based on two standard deviations from an individual black mouse video. We
484 lowered the minimum values further because we expected that certain mice would perform
485 poorly on the segmentation step. This allows Ctrax to still find a good location of the mouse
486 despite not being able to segment the entire mouse. This approach functions well, as all of our

487 setups have the same zoom of 8, and the mice tested are generally the same shape. Motion
488 settings are very lenient, because our experimental setup tracks only one mouse in the arena at
489 a time. Under the observation parameters, we primarily utilize the “Min Area Ignore” setting to
490 filter out detections larger than 2,500 pixels. Under the hindsight tab, we use the “Fix Spurious
491 Detections” setting to remove detections with a length shorter than 500 frames.

492 **Training sets**

493 **Labeling software**

494 We annotated our own training data using custom software that was written to
495 accommodate obtaining the necessary labels. We used the OpenCV library (<https://opencv.org/>)
496 to create an interactive watershed-based segmentation and contour-based ellipse-fit. Using the
497 software GUI we developed, the user left-clicks to mark points as the foreground (a mouse) and
498 right-clicks to label other points as the background (S1 Fig.). Upon a keystroke, the watershed
499 algorithm is executed to predict a segmentation and ellipse. If users need to make edits to the
500 predicted segmentation and ellipse, they can simply mark additional areas and run the
501 watershed again. When the predictions are of sufficiently high quality, users then select the
502 direction of the ellipse. They do this by selecting one of four cardinal directions: up, down, left,
503 right. Since the exact angle is selected by the ellipse-fitting algorithm, users need only to identify
504 the direction ± 90 degrees. Once a direction is selected, all the relevant data is saved to disk and
505 users are presented with a new frame to label. Full details on the software controls can be found
506 in the software documentation.

507 The objective of our annotated dataset is to identify good ellipse-fit tracking data for
508 mice. While labeling data, we optimized the ellipse-fit such that the ellipse was centered on the
509 mouse’s torso with the major axis edge approximately touching the nose of the mouse.

510 Frequently, the tail was removed from the segmentation mask to provide a better ellipse-fit. For
511 training networks for inference, we created three annotated training sets. Each training dataset
512 includes a reference frame (input), segmentation mask, and ellipse-fit. Each training set was
513 generated to track mice in a different environmental setup.

514 **Neural network models**

515 The neural networks we trained fall into three categories: segmentation, regression, and
516 binning. Our tested models can be viewed visually in Fig. 1D-F.

517 The first network architecture is modeled after a typical encoder-decoder structure for
518 segmentation (Fig. 1D). The first half of the network (encoder) utilizes 2D convolutional layers
519 followed by batch normalization, a ReLu activation, and 2D max pooling layers. We use a
520 starting filter size of 8 that doubles after every pooling layer. The kernels used are of shape 5x5
521 for 2D convolution layers and 2x2 for max pooling layers. Our input is of shape 480x480x1 and
522 after six of these repeated layers, the resulting shape is 15x15x128. We apply another 2D
523 convolutional layer (kernel 5x5, 2x filters) followed by a 2D max pool with a different kernel of
524 3x3 and stride of 3. One final 2D convolutional layer is applied to yield our feature bottleneck
525 with a shape of 5x5x512. This feature bottleneck is then passed to both the segmentation
526 decoder and angle predictor. The segmentation decoder reverses the encoder using strided
527 transpose 2D convolutional layers and carries over pre-downsampled activations through
528 summation junctions. It should be noted that this decoder does not utilize ReLu activations. After
529 the layers return to the 480x480x8 shape, we apply one additional convolution, with a kernel
530 size of 1x1, to merge the depth into two images: background prediction and foreground
531 prediction. We apply a softmax function across this depth. From the feature bottleneck, we also
532 create a prediction for angle prediction. We achieve this by applying two 2D convolution layers
533 with batch normalization and ReLu activations (kernel size 5x5, feature depths 128 and 64).

534 From here, we flatten and use one fully connected layer to yield a shape of four neurons, which
535 function to predict the quadrant in which the mouse's head is facing. Since the angle is predicted
536 by the mask, we need only to select the correct direction (± 180 deg). The four possible
537 directions that the network can select are 45-135, 135-225, 225-315 and 315-45 degrees on a
538 polar coordinate grid. These boundaries were selected to avoid discontinuities in angle
539 prediction.

540 The second network architecture is a binned regression approach (Fig. 1E). Instead of
541 predicting the parameters directly, the network instead selects the most probable value from a
542 selection of binned possible values. The major difference between this structure and a
543 regression structure is that the binned regression network training relies on a cross entropy loss
544 function whereas a regression network relies on a mean squared error loss function. Due to
545 memory limitations, we tested only custom VGG-like networks with reduced feature dimensions.
546 Our best-performing network is structured with two 2D convolutional layers followed by a 2D
547 max pooling layer. The kernels used are of shape 3x3 for 2D convolutional layers and 2x2 for 2D
548 max pooling layers. We start with a filter depth of 16 and double after every 2D max pool layer.
549 This two convolutional plus max pool sequence is repeated five times to yield a shape of
550 15x15x256. This layer is flattened and connected to a fully connected layer for each output
551 ellipse-fit parameter. The shape of each output is dictated by the desired resolution and range of
552 the prediction. For testing purposes, we observed only the center location and trained with a
553 range of the entire image (0-480). Additional outputs, such as angle prediction, could simply be
554 added as additional output vectors.

555 The third network architecture is modeled after a typical regression predictor structure
556 (Fig. 1F). While the majority of regression predictors realize the solution through a bounding box,
557 an ellipse simply adds one additional parameter: the angle of the mouse's head direction. Since

558 the angle is a repeating series with equivalence at 360deg and 0deg, we transform the angle
559 parameter into its sine and cosine components. This yields a total of six parameters regressed
560 from the network. The first half of this network encodes a set of features relevant to correctly
561 predicting the six parameters. From the encoded feature set, we flatten the network and applied
562 a fully convolutional layer to regress the parameters for the ellipse-fit. We tested a wide variety
563 of pre-built feature detectors including Resnet V2 50, Resnet V2 101, Resnet V2 200, Inception
564 V3, Inception V4, VGG, and Alexnet. In addition to these pre-built feature detectors, we also
565 surveyed a wide array of custom networks. Of these general purpose feature encoders and
566 custom networks, Resnet V2 200 performed the best.

567 **Neural network training**

568 This section describes all of the procedures pertaining to training our neural network
569 models. The three procedures described here are training set augmentation, training
570 hyperparameters, and a benchmark for training set size.

571 Training set augmentation has been an important aspect of training neural networks
572 since Alexnet [25]. We utilize a handful of training set augmentation approaches to achieve good
573 regularization performance. Since our data is from a birds-eye view, it is straightforward to apply
574 horizontal, vertical, and diagonal reflections for an immediate 8x increase in our equivalent
575 training set size. Additionally, at runtime, we apply small rotations and translations for the entire
576 frame. Rotation augmentation values are sampled from a uniform distribution. Finally, we apply
577 noise, brightness, and contrast augmentations to the frame. The random values used for these
578 augmentations are selected from a normal distribution.

579 Hyperparameters, such as training learn rate and batch size, were selected
580 independently for each network architecture trained. While larger networks, such as Resnet V2

581 200, can run into memory limitations for batch sizes at an input size of 480x480, good learn rate
582 and batch size were experimentally identified using a grid search approach [26]. Table A in S1
583 Information summarizes all the hyperparameters selected for training these network
584 architectures.

585 We also benchmarked the influence of training set size on network generalization in
586 order to determine the approximate amount of annotated training data required for good network
587 performance of the encoder-decoder segmentation network architecture (S2 Fig.). We tested
588 this benchmark by shuffling and randomly sampling a subset of the training set. Each
589 subsampled training set was trained and compared to an identical validation set. While the
590 training curves appear indistinguishable, the validation curves trained with fewer than 2,500
591 training annotations diverge from the group. This suggests that the training set is no longer large
592 enough to allow the network to generalize well. While the exact number of training samples will
593 ultimately rely on the difficulty of the visual problem, a recommended starting point would be
594 around 2,500 training annotations.

595 **Animals used**

596 All animals were obtained from The Jackson Laboratory production colonies. Adult mice
597 aged 8 to 14 weeks were behaviorally tested in accordance with approved protocols from The
598 Jackson Laboratory Institutional Animal Care and Use Committee guidelines. Open field
599 behavioral assays were carried out as previously described [27]. Briefly, group-housed mice
600 were weighed and allowed to acclimate in the testing room for 30-45 minutes before the start of
601 video recording. Data from the first 55 minutes of activity are presented here. Where available, 8
602 males and 8 females were tested from each inbred strain and F1 isogenic strain.

603 **Code and training set availability**

604 Neural network training and inference code as well as annotated datasets will become
605 available upon publication.

606 **Acknowledgments**

607 We thank the members of the Kumar laboratory for suggestions and editing of the
608 manuscript. We thank JAX Information Technology team members Edwardo Zaborowski, Shane
609 Sanders, Rich Brey, David McKenzie, and Jason Macklin for infrastructure support; and we
610 thank KOMP2 behavioral testers James Clark, Pamela Fraungruber, Rose Presby, Zachery
611 Seavey, and Catherine Witmeyer. This work used the National Science Foundation (NSF)
612 Extreme Science and Engineering Discovery Environment (XSEDE) XStream service at
613 Stanford University through allocation TG-DBS170004. Funding also included the NIH grant
614 DA041668 from NIDA, a Brain and Behavioral Foundation Young Investigator Award, and a JAX
615 Director's Innovation Fund to V. K. We thank Nvidia GPU grant program for providing a TitanXP
616 for our work.

617 **References**

- 618 1. Egnor SE, Branson K. Computational Analysis of Behavior. *Annu Rev Neurosci.* 2016;39:217-36.
619 doi: 10.1146/annurev-neuro-070815-013845. PubMed PMID: 27090952.
- 620 2. Gomez-Marin A, Paton JJ, Kampff AR, Costa RM, Mainen ZF. Big behavioral data: psychology,
621 ethology and the foundations of neuroscience. *Nat Neurosci.* 2014;17(11):1455-62. doi: 10.1038/nn.3812.
622 PubMed PMID: 25349912.
- 623 3. Valdez P, Mehrabian A. Effects of color on emotions. *Journal of experimental psychology:*
624 *General.* 1994;123(4):394.
- 625 4. Kuleskaya N, Voikar V. Assessment of mouse anxiety-like behavior in the light-dark box and
626 open-field arena: role of equipment and procedure. *Physiol Behav.* 2014;133:30-8. doi:
627 10.1016/j.physbeh.2014.05.006. PubMed PMID: 24832050.
- 628 5. Dell AI, Bender JA, Branson K, Couzin ID, de Polavieja GG, Noldus LP, et al. Automated image-
629 based tracking and its application in ecology. *Trends in ecology & evolution.* 2014;29(7):417-28.
- 630 6. Anderson DJ, Perona P. Toward a science of computational ethology. *Neuron.* 2014;84(1):18-31.

- 631 7. Branson K, Robie AA, Bender J, Perona P, Dickinson MH. High-throughput ethomics in large
632 groups of *Drosophila*. *Nature methods*. 2009;6(6):451.
- 633 8. Austin CP, Battey JF, Bradley A, Bucan M, Capecchi M, Collins FS, et al. The knockout mouse
634 project. *Nat Genet*. 2004;36(9):921-4. doi: 10.1038/ng0904-921. PubMed PMID: 15340423; PubMed
635 Central PMCID: PMCPMC2716027.
- 636 9. Kabra M, Robie AA, Rivera-Alba M, Branson S, Branson K. JAABA: interactive machine learning
637 for automatic annotation of animal behavior. *Nat Methods*. 2013;10(1):64-7. doi: 10.1038/nmeth.2281.
638 PubMed PMID: 23202433.
- 639 10. Sobral A, editor BGSLibrary: An opencv c++ background subtraction library. IX Workshop de
640 Visao Computacional; 2013.
- 641 11. Dankert H, Wang L, Hoopfer ED, Anderson DJ, Perona P. Automated monitoring and analysis of
642 social behavior in *Drosophila*. *Nature methods*. 2009;6(4):297.
- 643 12. Noldus LP, Spink AJ, Tegelenbosch RA. EthoVision: a versatile video tracking system for
644 automation of behavioral experiments. *Behavior Research Methods, Instruments, & Computers*.
645 2001;33(3):398-414.
- 646 13. Pérez-Escudero A, Vicente-Page J, Hinz RC, Arganda S, De Polavieja GG. idTracker: tracking
647 individuals in a group by automatic identification of unmarked animals. *Nature methods*. 2014;11(7):743.
- 648 14. De Chaumont F, Coura RD-S, Serreau P, Cressant A, Chabout J, Granon S, et al. Computerized
649 video analysis of social interactions in mice. *Nature methods*. 2012;9(4):410.
- 650 15. Ohayon S, Avni O, Taylor AL, Perona P, Egnor SR. Automated multi-day tracking of marked mice
651 for the analysis of social behaviour. *Journal of neuroscience methods*. 2013;219(1):10-9.
- 652 16. Patel TP, Gullotti DM, Hernandez P, O'Brien WT, Capehart BP, Morrison III B, et al. An open-
653 source toolbox for automated phenotyping of mice in behavioral tasks. *Frontiers in behavioral*
654 *neuroscience*. 2014;8:349.
- 655 17. Hewitt BM, Yap MH, Hodson-Tole EF, Kennerley AJ, Sharp PS, Grant RA. A novel automated
656 rodent tracker (ART), demonstrated in a mouse model of amyotrophic lateral sclerosis. *Journal of*
657 *neuroscience methods*. 2017.
- 658 18. LeCun Y, Bengio Y, Hinton G. Deep learning. *Nature*. 2015;521(7553):436-44. doi:
659 10.1038/nature14539. PubMed PMID: 26017442.
- 660 19. Yosinski J, Clune J, Bengio Y, Lipson H, editors. How transferable are features in deep neural
661 networks? *Advances in neural information processing systems*; 2014.
- 662 20. He K, Zhang X, Ren S, Sun J, editors. Identity mappings in deep residual networks. *European*
663 *Conference on Computer Vision*; 2016: Springer.
- 664 21. Abadi M, Barham P, Chen J, Chen Z, Davis A, Dean J, et al., editors. TensorFlow: A System for
665 Large-Scale Machine Learning. OSDI; 2016.
- 666 22. Zoph B, Vasudevan V, Shlens J, Le QV. Learning transferable architectures for scalable image
667 recognition. *arXiv preprint arXiv:170707012*. 2017.

- 668 23. Bogue MA, Grubb SC, Walton DO, Philip VM, Kolishovski G, Stearns T, et al. Mouse Phenome
669 Database: an integrative database and analysis suite for curated empirical phenotype data from laboratory
670 mice. *Nucleic Acids Res.* 2018;46(D1):D843-D50. doi: 10.1093/nar/gkx1082. PubMed PMID: 29136208;
671 PubMed Central PMCID: PMC5753241.
- 672 24. Egnor SR, Branson K. Computational analysis of behavior. *Annual review of neuroscience.*
673 2016;39:217-36.
- 674 25. Krizhevsky A, Sutskever I, Hinton GE, editors. Imagenet classification with deep convolutional
675 neural networks. *Advances in neural information processing systems*; 2012.
- 676 26. Larochelle H, Erhan D, Courville A, Bergstra J, Bengio Y, editors. An empirical evaluation of deep
677 architectures on problems with many factors of variation. *Proceedings of the 24th international conference*
678 *on Machine learning*; 2007: ACM.
- 679 27. Kumar V, Kim K, Joseph C, Thomas LC, Hong H, Takahashi JS. Second-generation high-
680 throughput forward genetic screen in mice to isolate subtle behavioral mutants. *Proceedings of the*
681 *National Academy of Sciences.* 2011;108(Supplement 3):15557-64.
682

683 **Supporting information**

684 **S1 Information. (Note A) Description of the ellipse-fit function. (Note B) Annotated**
685 **dataset descriptions. (Table A) Training hyperparameters.**

686 **S1 Fig. Example of our labeling GUI software.** (A) Our software allows the user to
687 zoom into the the region of interest for annotation (mouse) and placed two marks: one
688 for foreground (green) and one for background (red). (B) Upon a keystroke, the software
689 provides the resulting segmentation (magenta), ellipse-fit (cyan), and the old
690 background annotations (yellow).

691 **S2 Fig. Training set size scaling benchmark.** We benchmarked how the training-set
692 size influences the performance of a trained encoder-decoder segmentation network.
693 Full training set includes 16,234 annotated frames. (A) Training-set size does not impact
694 training set error rate. (B) Validation performance converges to the same value above
695 2,500 training samples, but the error rate increases when 1,000 or fewer training

696 samples are used. (C-F) Validation accuracy outperforms training accuracy when 2,500
697 or more training samples are used. (G) Validation accuracy begins to show signs of
698 weak generalization by only matching, and not exceeding, training accuracy at 1,000
699 training samples. (H) A network trained using only 500 training samples is clearly
700 overtraining, shown by the diverging and increasing validation error rate.

701 **S1 Video. Comparison of mouse tracking.** A comparison of mouse tracking across a
702 variety of coat colors and environments using both our proposed encoder-decoder
703 segmentation neural network (red) and Ctrax (blue). (0-22s) Black mice and (22-45s)
704 grey mice in a white environment have strong agreement across approaches. When
705 rearing on the wall, Ctrax starts to not properly fit the ellipse. (45-66s) Piebald mice in a
706 white environment have strong tracking concordance, but depending upon the unique
707 coat pattern may have incorrect shape predicted by Ctrax. (66-90s) Albino mice on a
708 white background are a difficult problem for background subtraction approaches (Ctrax),
709 while a neural network approach tracks appropriately. (90-112s) Black mice in the 24-hr
710 setup, which contains bedding and a food cup, are difficult for background subtraction
711 approaches (Ctrax) to create adequate background models for tracking. A neural
712 network approach learns to handle this difficulty. (112-134s) Black mice in the KOMP2
713 arena, which has reflective floors and walls, poses a difficult situation for background
714 subtraction approaches (Ctrax). A neural network approach learns to not include
715 reflections without any tuning of parameters. Playback for all clips in this video are at
716 half-speed to better observe and compare tracking performance.

717 **S2 Video. KOMP2 observed odd behavior.** A 1-minute sample from the two off-

718 diagonal KOMP2 videos. In the first clip (0-62s), we observe a high degree of waddle in
719 the animal's gait as well as odd stride frequency. In the second clip (62-125s), we
720 observe a hunched posture during locomotion as well as a frequent sideways motion.
721 Red ellipse denotes our neural network tracker prediction.





