1    _**Supplementary material to:**_

2    **Automatic mapping of multiplexed social receptive fields**

3    **by deep learning and GPU-accelerated 3D videography**
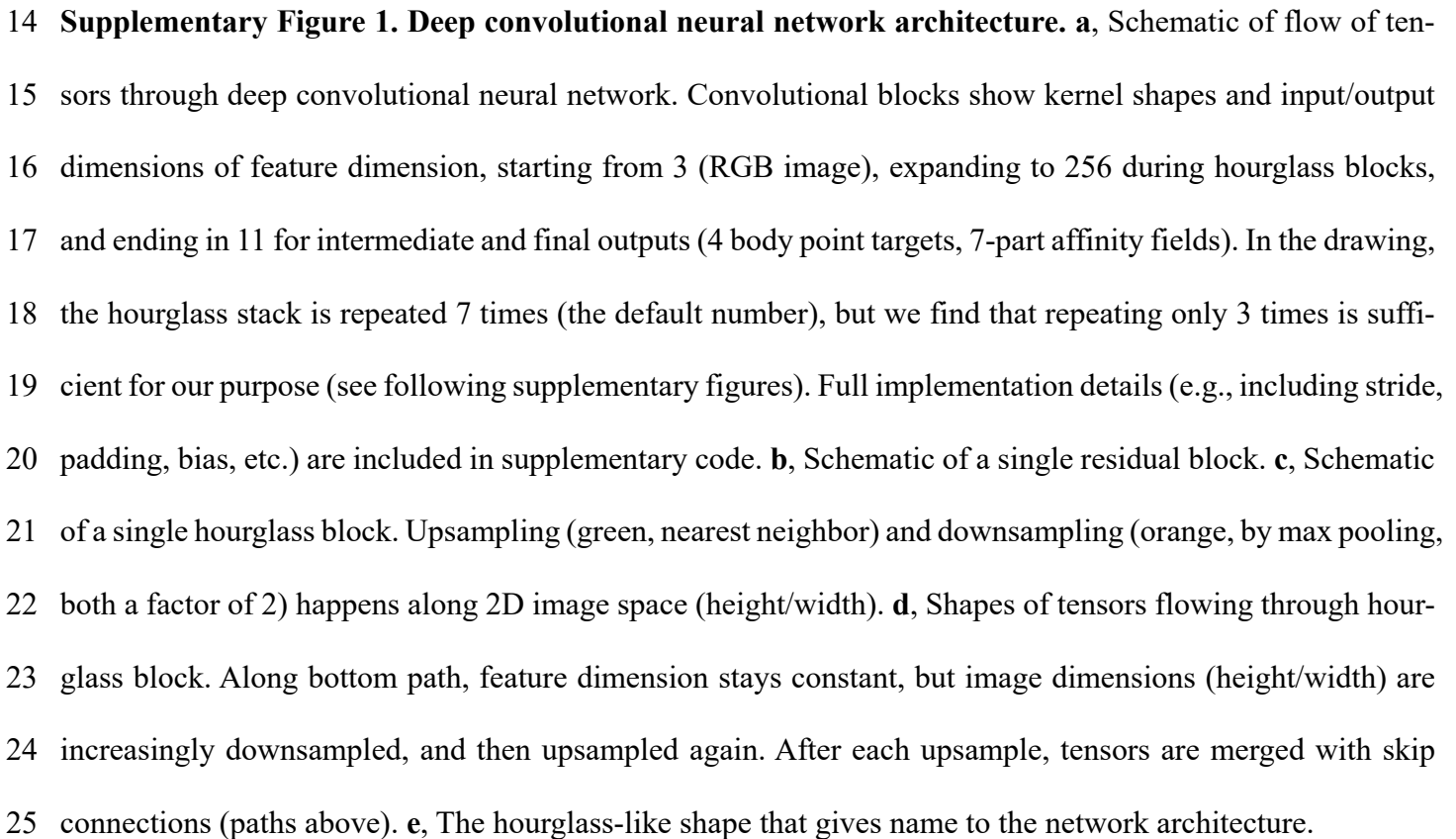
4

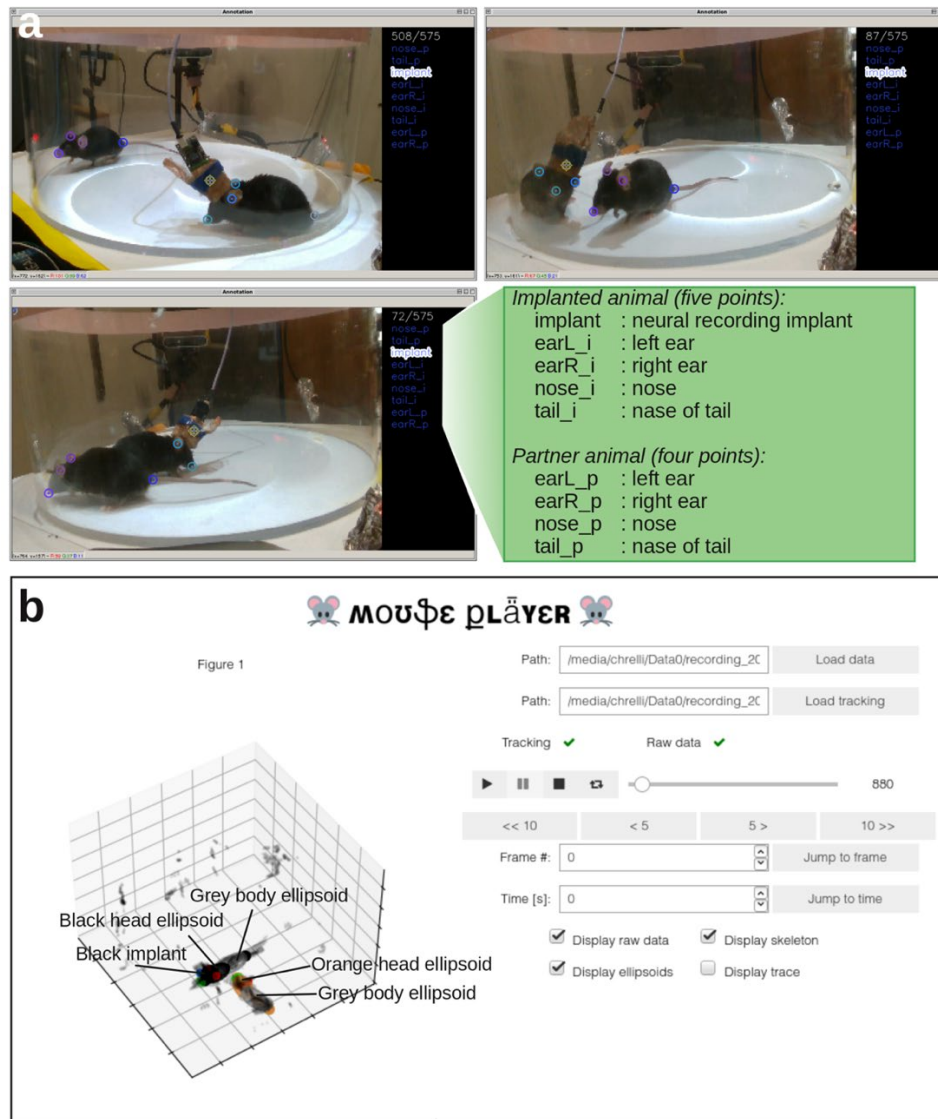5    Christian L. Ebbesen[1,2],* & Robert C. Froemke[1,2],*

6

7    [1] Skirball Institute of Biomolecular Medicine, Neuroscience Institute, Departments of Otolaryngology, Neu-

8    roscience and Physiology, New York University School of Medicine, New York, NY, 10016, USA.

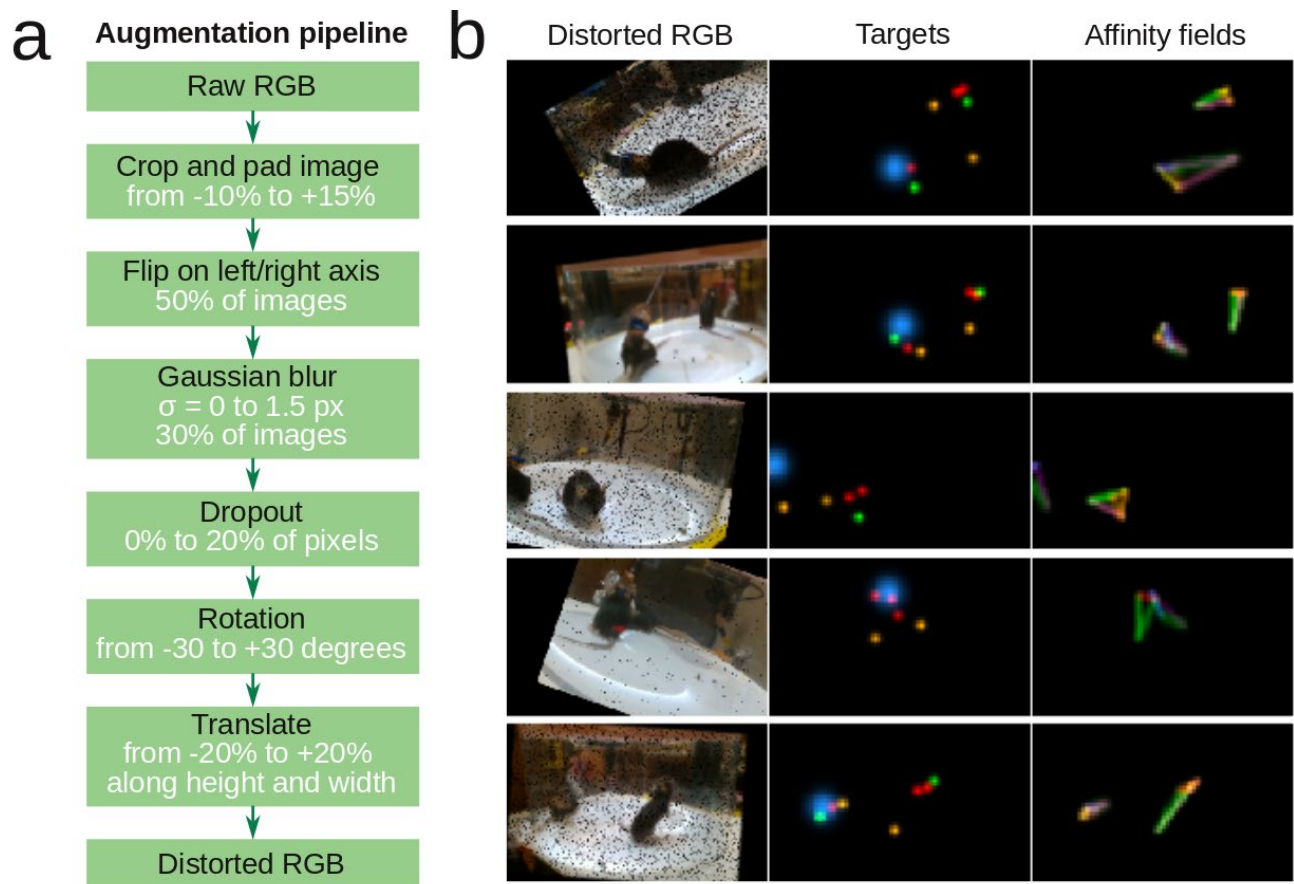9    [2] Center for Neural Science, New York University, New York, NY, 10003, USA.

10   * Correspondence to: C.L.E. (christian.ebbesen@nyumc.org) or R.C.F. (robert.froemke@med.nyu.edu)

11  **Supplementary Figures and Legends**



13

14  **Supplementary Figure 1. Deep convolutional neural network architecture. a**, Schematic of flow of ten-

15  sors through deep convolutional neural network. Convolutional blocks show kernel shapes and input/output

16  dimensions of feature dimension, starting from 3 (RGB image), expanding to 256 during hourglass blocks,

17  and ending in 11 for intermediate and final outputs (4 body point targets, 7-part affinity fields). In the drawing,

18  the hourglass stack is repeated 7 times (the default number), but we find that repeating only 3 times is suffi-

19  cient for our purpose (see following supplementary figures). Full implementation details (e.g., including stride,

20  padding, bias, etc.) are included in supplementary code. **b**, Schematic of a single residual block. **c**, Schematic

21  of a single hourglass block. Upsampling (green, nearest neighbor) and downsampling (orange, by max pooling,

22  both a factor of 2) happens along 2D image space (height/width). **d**, Shapes of tensors flowing through hour-

23  glass block. Along bottom path, feature dimension stays constant, but image dimensions (height/width) are

24  increasingly downsampled, and then upsampled again. After each upsample, tensors are merged with skip

25  connections (paths above). **e**, The hourglass-like shape that gives name to the network architecture.

**Implanted animal (five points):**
- implant : neural recording implant
- earL_i : left ear
- earR_i : right ear
- nose_i : nose
- tail_i : nase of tail

**Partner animal (four points):**
- earL_p : left ear
- earR_p : right ear
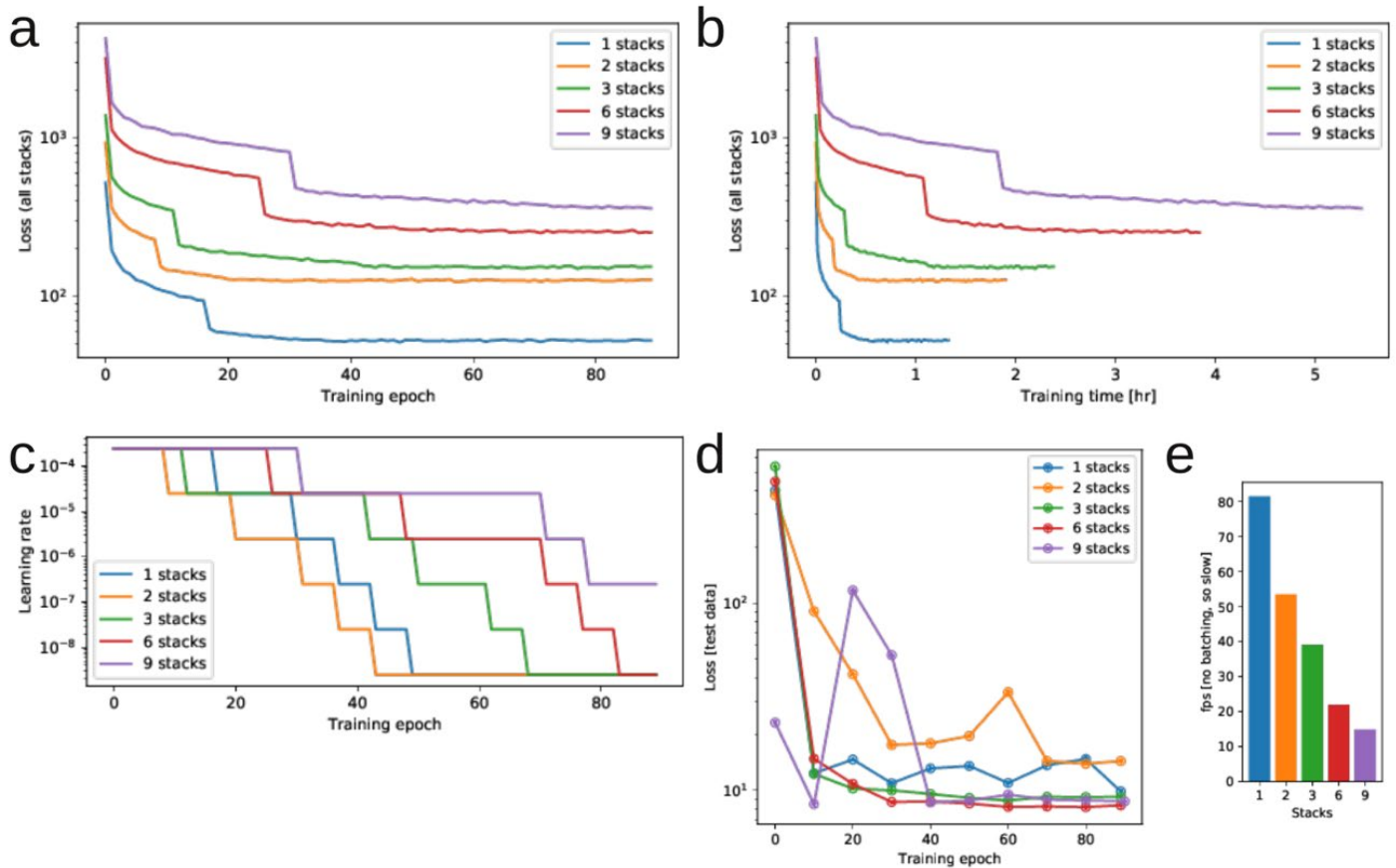- nose_p : nose
- tail_p : nase of tail

🐭 MOU♇E ♇LÄYER 🐭

27   **Supplementary Figure 2. GUIs for labeling of training data for the neural network and for viewing**

28   **tracked data. a**, For training the network to recognize body parts, we must generate labeled frames by manual

29   annotation. For each frame, 1-5 body parts are labeled on the implanted animal and 1-4 body parts on the

30   partner animal. This can be done with any annotation software; we used a modified version of the free 'Deep-

31   PoseKit-Annotator' (Graving et al., 2019) (https://github.com/jgraving/DeepPoseKit-Annotator/) included in

32   the supplementary code. This software allows easy labeling of the necessary points, and pre-packages training

33   data for use in our training pipeline. Body parts are indexed by i/p for implanted/partner animal ('nose_p' is

34   the nose of the partner animal, for example). **b,** GUI for viewing and quality control of tracked behavior (raw

35   data, body skeleton, ellipsoid surfaces and time trajectory) running in an interactive Jupyter notebook.
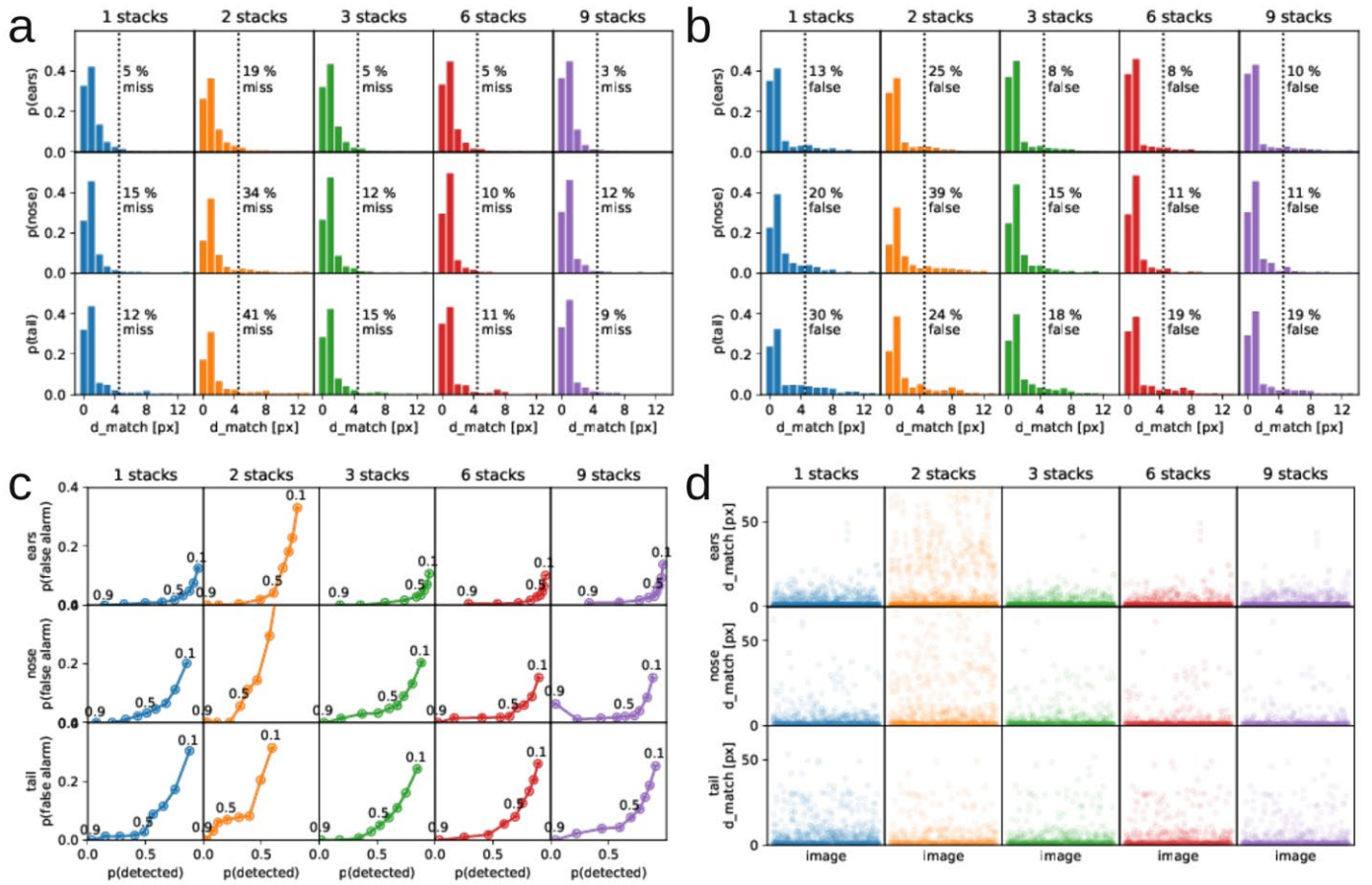
**a Augmentation pipeline**

Raw RGB
↓
Crop and pad image
from -10% to +15%
↓
Flip on left/right axis
50% of images
↓
Gaussian blur
σ = 0 to 1.5 px
30% of images
↓
Dropout
0% to 20% of pixels
↓
Rotation
from -30 to +30 degrees
↓
Translate
from -20% to +20%
along height and width
↓
Distorted RGB

**b** Distorted RGB    Targets    Affinity fields

37  **Supplementary Figure 3. Augmentation pipeline for network training. a**, Flowchart of augmentation

38  pipeline used to generate distorted frames during network training. **b**, Examples of distorted labeling frames

39  generated by augmentation pipeline, as well as corresponding body part targets and affinity fields used dur-

40  ing training.

41

42

**Supplementary Figure 4. Network training history and performance as a function of hourglass stacks.**
**a**, Loss history as a function of training epoch, for networks with 1, 2, 3, 6, and 9 hourglass stacks. **b,** same as panel a, but as a function of training time. **c,** Learning rate schedule (automatically adjusted) as a function of training epoch. **d,** Test loss, as a function of training epoch, for networks with 1, 2, 3, 6, and 9 hourglass stacks. Beyond 3 stacks, there was little improvement in the training loss. **e,** Relative inference time with no batching (since the batch size will have to be smaller for a network with more stacks – for real use, we used batched inference).

**Supplementary Figure 5. Network sensitivity and precision as a function of hourglass stacks. a,** Number of 'missed' keypoints, as a function of number of network stacks, after full training, for ear, nose and tail keypoints (the three rows). We count a hand-labeled body keypoint as found by the network as 'detected' if it was within 5 pixels of a keypoint suggested by the network. **b,** Same as **a**, but counting the number of 'false detections'. We defined a false detection, as a keypoint suggested by the network, which was more than 5 pixels from a corresponding hand-labeled keypoint. **c,** Plots of false alarm rate (p(false alarm)) as a function of detected rate (p(detected)) for all keypoints in the test data, plotted for different probability cutoffs (from 0.1 to 0.9 in 0.1 steps, indicated by dots on the curves). **d,** Shortest distance from a proposed keypoint location to a hand-labeled keypoint location (d_match), across images in the test data, for ear, nose and tail keypoints (rows) as a function of image stacks (columns).
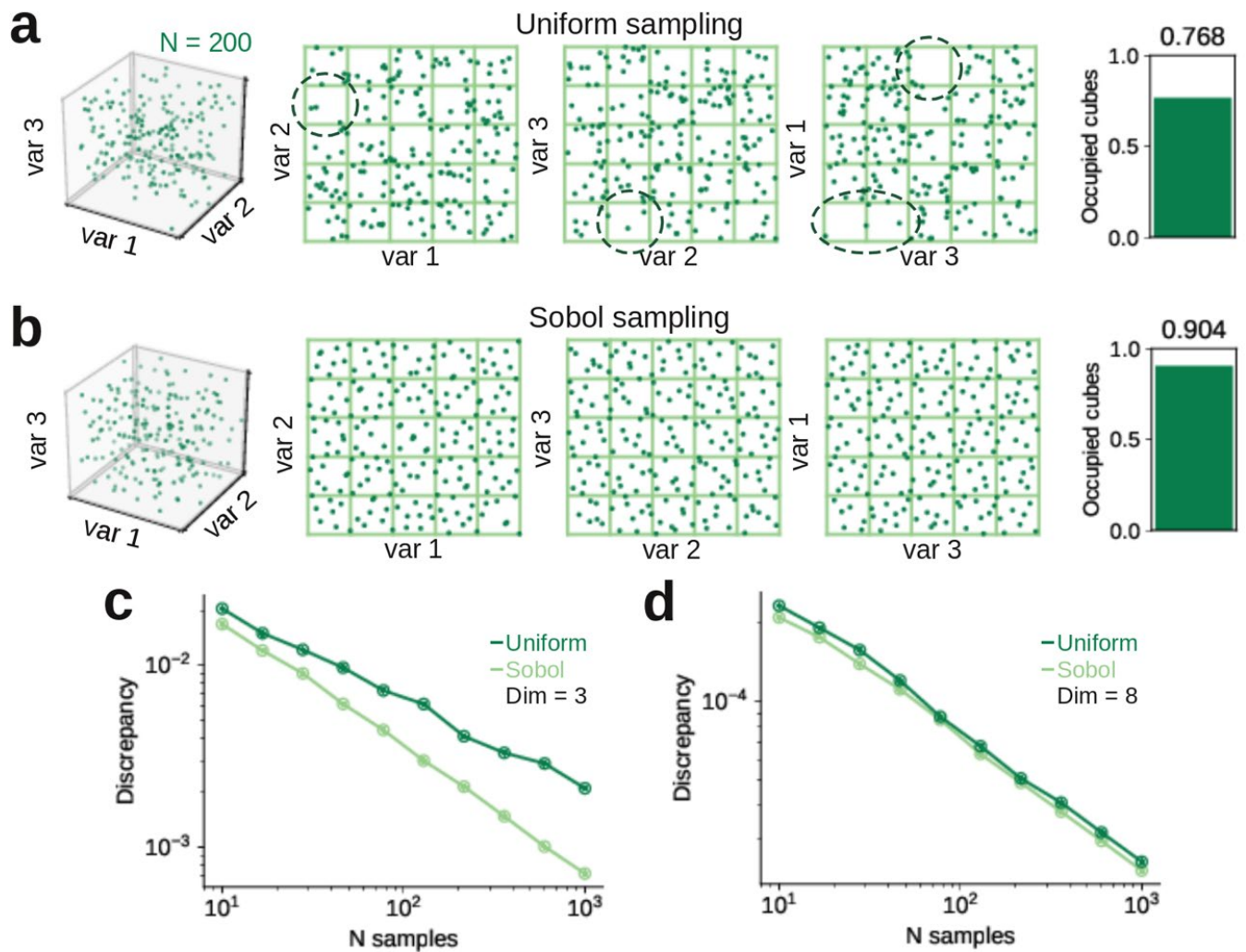
82 and a lower number of false alarm and missed detections. As the performance differences with/without forcing

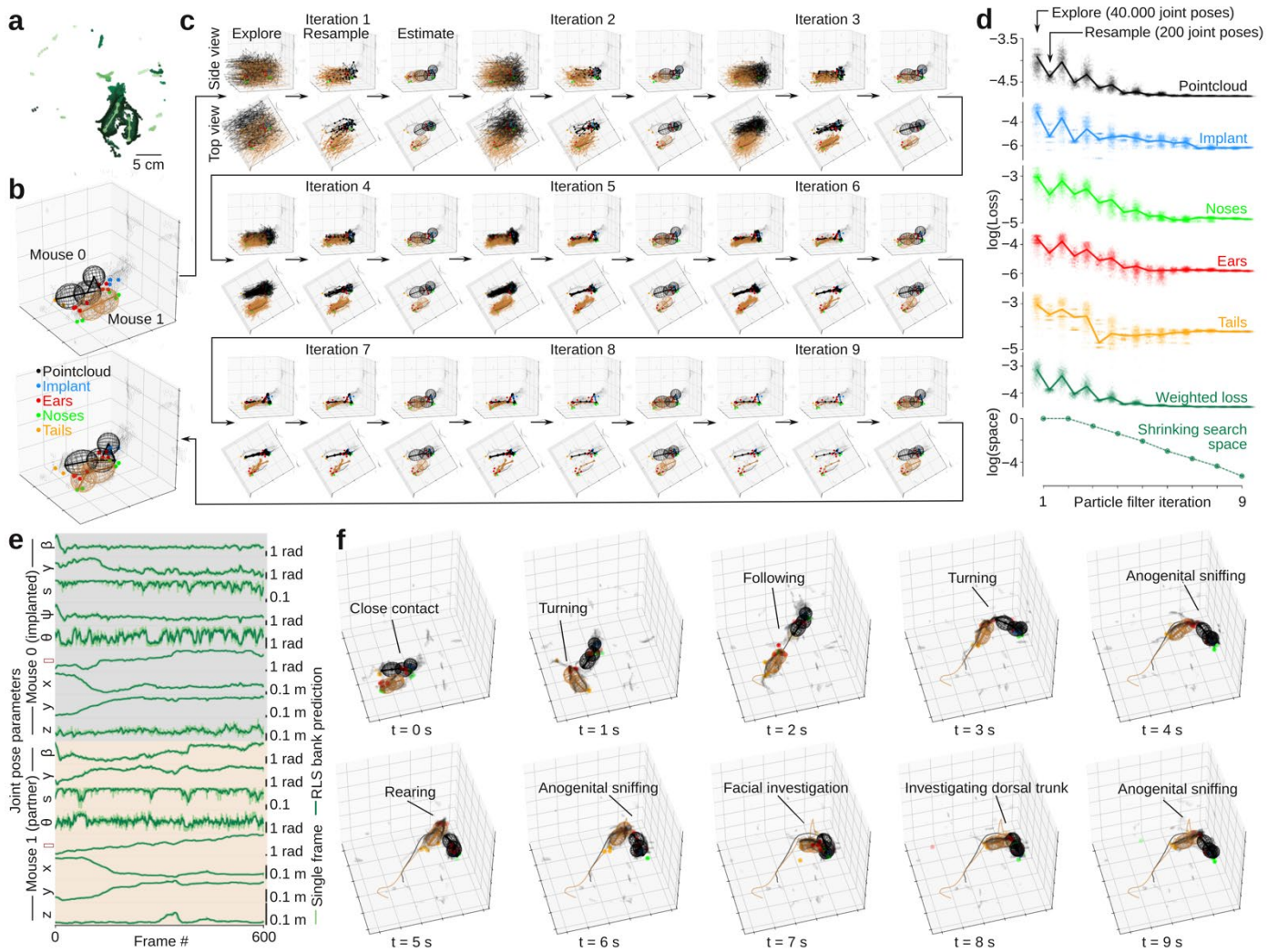83 the PAFs are minimal, our code imposes the PAFs by default.

86 **Supplementary Figure 8. Loss function calculation details. a**, Shortest distance to surface of an ellipsoid,

87 $d$, and our approximation, $\tilde{d}$. **b**, $\tilde{d}$ is a good approximation to $d$. Color map, value of $d/\tilde{d}$. White line, ellip-

88 soid surface. **c**, The loss function, $\rho$, associated with the pointcloud is the mean absolute error of the distance

89 estimate, truncated at +/- 3.0 cm **d**, Pixel density of the point-cloud depends on distance to the fixed-resolution

90 depth cameras. **e**, Pixel density is inversely proportional to the square of the distance to depth camera. **f**,

91 Overlap barrier spheres (implant sphere and spheres centered on the body ellipsoids with a radius equal to the

92 minor axis). **g**, Example of mirror symmetric body position (side-by-side, facing same direction), resulting in

93 ambiguity in animal identity if only one frame is considered. **h**, To include the context of previous frames, we

94 add an overlap loss penalty (similar to **f**) between each mouse and the position of the interaction partner in

95 the previous frame. In panel **g**, right, we would add a penalty term to the particle representing joint body pose.

96 In contrast, in panel **g**, left, this penalty is zero as there is no overlap with the position of the conspecific in

97 the previous frame.

**a** First frame in video: mice are too close

**b** Automatically scan forward until well-separated, automatically initialize from keypoints

Initialization

Initialization

Initialization

99 **Supplementary Figure 9. Automatic initialization procedure for the tracking algorithm. a**, Example

100 starting frame, where the mice are too close together for the automatic initialization. **b**, By scanning forward

101 in the 3D video, the algorithm finds a frame, where both the "head cluster" (cluster of detected ear/nose points)

102 and the "tail cluster" (cluster of detected tail points) are separated by a threshold distance. The algorithm uses

103 an average of the head, tail and implant clusters to initialize the tracking procedure (initial guess shown by

104 the black/brown lines, top row of plots shows the pointcloud data, bottom row of points shows only the initial
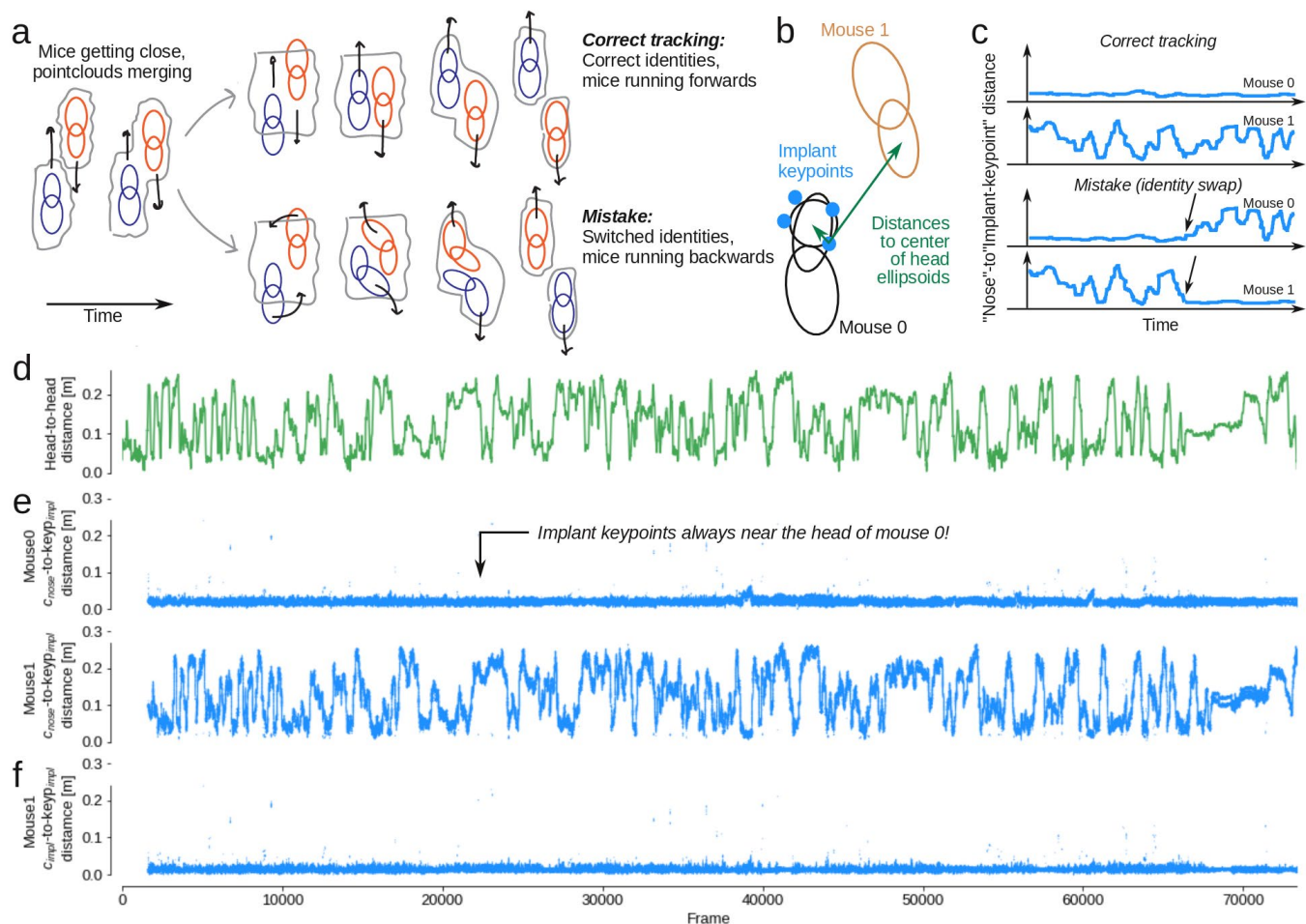
105 guesses).

107 **Supplementary Figure 10. Quasi-random particle filter exploration strategy. a**, Left, 3D plot. Middle,

108 2D projection plots of three random variables, drawn from independent uniform distributions. Points in 3D

109 space do not fill space well; in the 2D projections, there are squares (i.e., full rows, columns and pipes) of the

110 3D space not sampled at all (dashed lines). Right, partitioning space in 20%-cubes (green lines), only 76.8 %

111 of cubes are occupied. **b**, Same as **a**, but variables are drawn from quasi-random Sobol sequence (Sobol,

112 1967). Points are more evenly dispersed in space, and 90.4% of all 20%-cubes are sampled. **c**, Mean discrep-

113 ancy as function of sample number, for 3-dimensional (like panels **a,b**) uniform random sequence and a Sobol

114 sequence, calculated across 100 random sequences. The Sobol sequences have a lower discrepancy, i.e. sam-

115 ple more regions of space. **d**, Same as **c**, but for 17-dimensional variables (like our joint body posture particles).

116

117

**Supplementary Figure 11. Particle filter convergence and examples of tracked behavioral data. a,** Example of manual initialization of the tracking algorithm, by manual clicking of approximate locations of the two animals (light green dots, lines) on a top-down view of the behavioral arena (dark green dots, shade indicates z-coordinate). **b,** 3D view of initialized body model (top) and fitted body model (bottom) after running tracking algorithm on the frame. Black wireframe model, implanted mouse; brown wireframe model, partner animal. **c,** Particle filter state across 9 iterations of the fitting algorithm. After iteration 2, we shrink ('anneal') the exploration space with each step. **d,** Loss function values and size of filter search space across filter iterations. **e**, Tracked data (light green) and running adaptive estimate (dark green) across 600 frames (10 s). **f**, Data and fitted joint posture model, across 10 seconds of behavior. Trailing lines, location of hip ellipsoid center in the last 10 seconds.

**a**

Mouse 1    Mouse 0

$c_{mid}$  $c_{nose}$
$c_{tail}$  $c_{hip}$   $c_{impl}$
$c_{tip}$          $c_{tip}$
$c_{nose}$       $c_{tail}$
$c_{mid}$ $c_{hip}$

**b**

Raw tracking data          State-space filtered 3D data

$z\ y\ x\ \varphi\ \theta\ s\ \psi\ \gamma\ \beta$
$z\ y\ x\ \varphi\ \theta\ s\ \psi\ \gamma\ \beta$

$c_{tail}$    $c_{impl}$

$c_{hip}$    $c_{tip}$

$c_{mid}$    $s$

$c_{nose}$

10 s

**c**

Raw ellipsoid rotation data          Quaternion smoothing

$R_{body}$ (axis-angle)

Recompute rotation from 3D filtered skeleton

**d**

Recompute $c_{tip}$ and $c_{tail}$ from smoothed $s$ and $R_{body}/R_{nose}$

$c_{tip}$

$c_{tail}$

**e**

Raw tracking data          After filtering

**Supplementary Figure 12. State space filtering of tracked body models. a**, Estimated 3D locations of body model surfaces (wireframes, left) and skeletons (dots and lines, right) for an example frame. **b**, Fitted joint

131   pose parameters for the two mouse body models (left, 100 s snippet) and corresponding 3D coordinates of the

132   body skeleton points, and the spine scaling, $s$, for the implanted mouse (right, same 100 s). **c**, Raw 3D rotation

133   angle of the nose ellipsoid of the implanted animal (axis-angle representation), recalculated 3D rotation angles

134   from the filtered skeleton points, and final 3D rotation angles after quaternion smoothing (note the smoothing

135   out of noise, indicated by arrow). **d**, Recalculated c_nose and c_tail from the smoothed 3D rotations and

136   smoothed spine scaling. **e**, Example frame before (left) and after state space filtering of the tracked data (right).

138 **Supplementary Figure 13. Implant-to-nose distance demonstrates that there are no swapped identities.**

139 **a**, Schematic showing two common errors in tracking algorithms: Swapped identities and swapped directions.

140 When the mice approach each other, their point clouds will merge. Because resolution and frame rates are

141 limited, it can be hard to estimate body postures in this configuration. For example, if the tracking algorithm

142 is not properly spatiotemporally regularized, the algorithm might mistakenly swap mouse identities, such that

143 the mice appear to be running backwards (shown in bottom row). Direction swaps and identity swaps can also

144 happen independently. For example, when mice are allogrooming, or passing over/under each other, identities
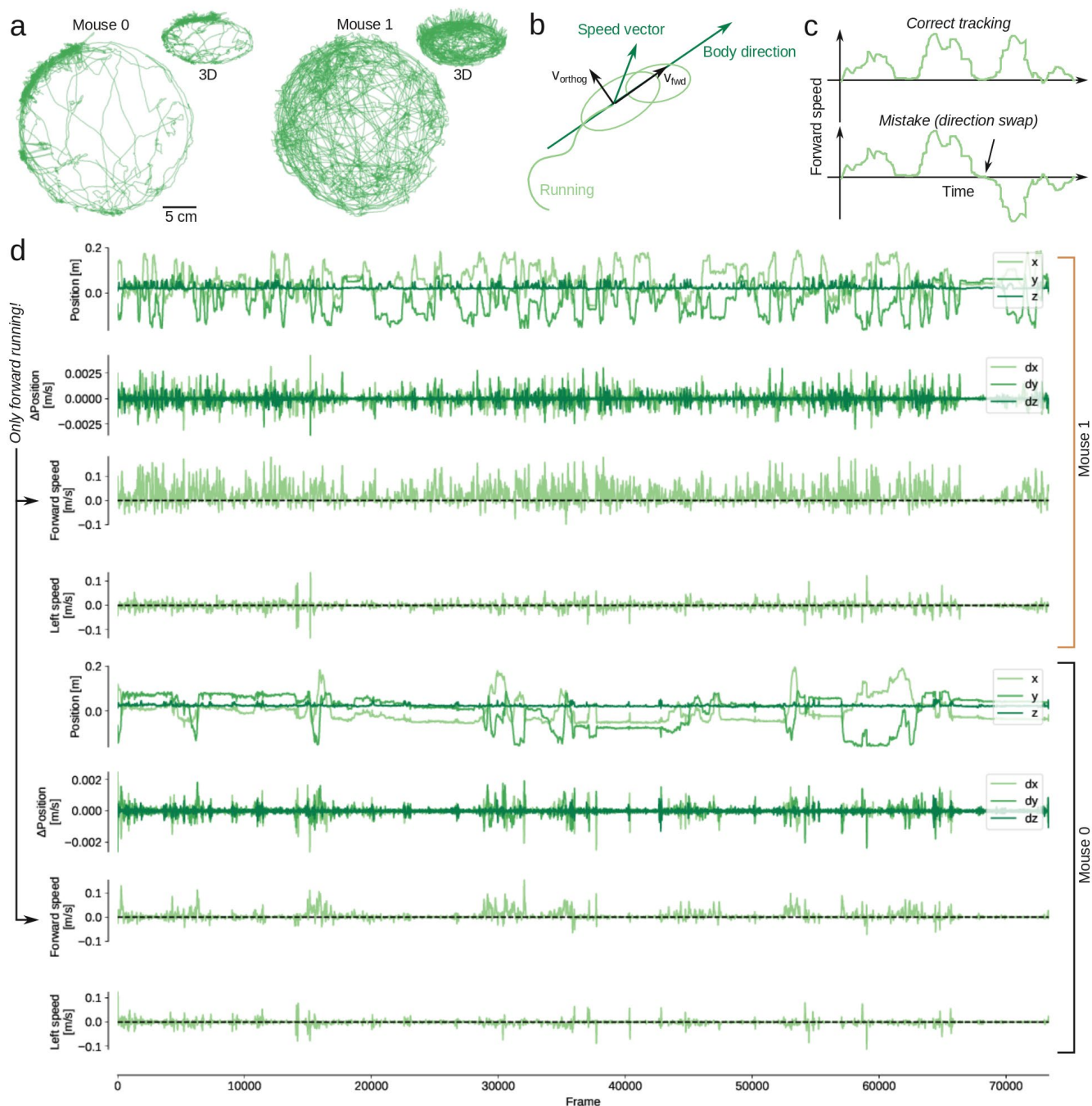
145 might swap, but both mice can still appear to run normally with no apparent errors. Conversely, when a mouse

146 is self-grooming, their point-cloud essentially resembles a ball, and when they start moving again, it may not

147 be clear if they are 'really' moving forward or backwards. **b**, For all frames, we calculated the distance be-
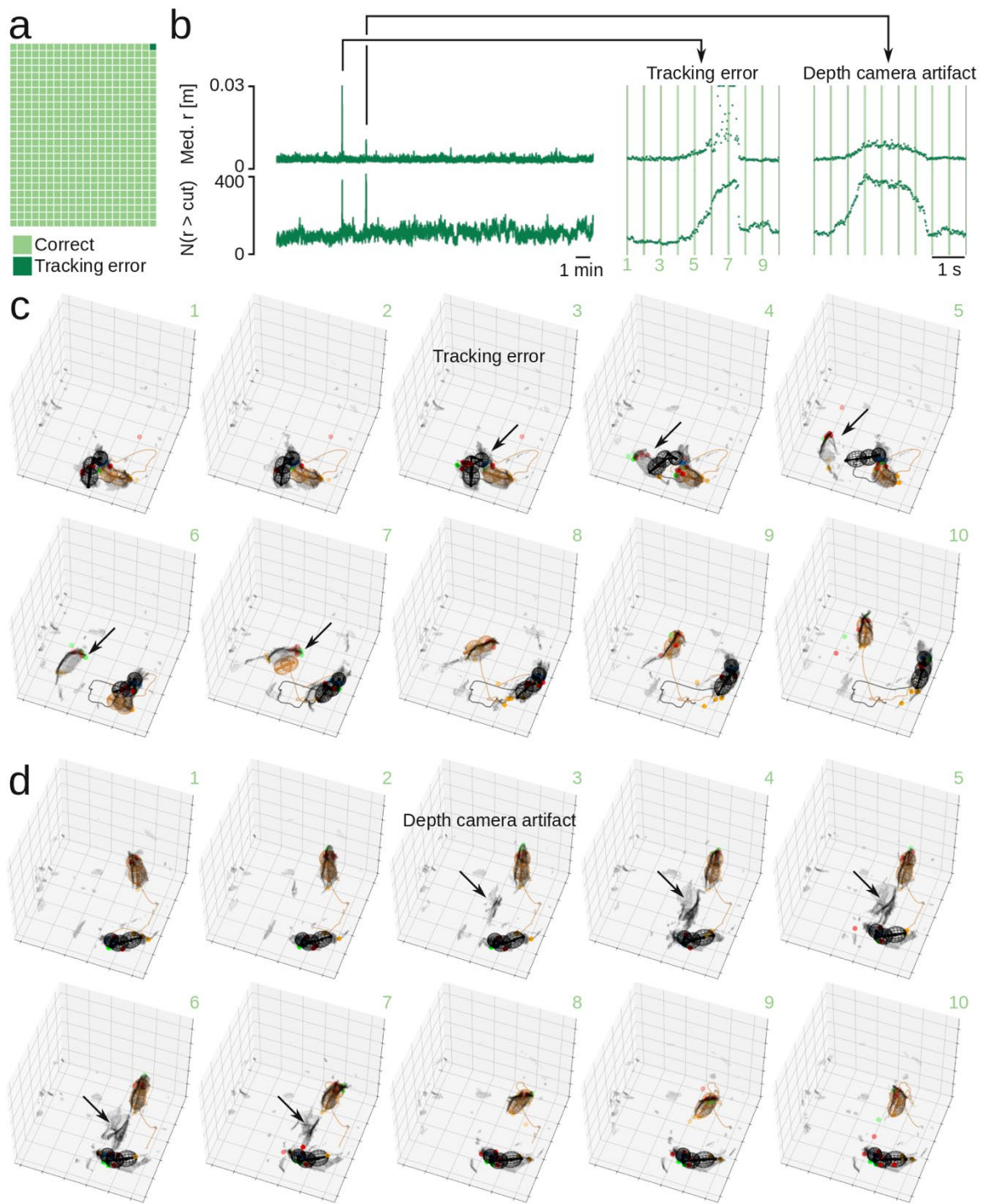
148 tween implant key-points and the centroid of both nose ellipsoids. **c**, If there is an identity swap of the mice,

149    this is will be evident in the distance between the implant key-points and the head of both mice. In correct

150    tracking (top row), implant body model always follows the same mouse. In tracking with mistakes (bottom

151    row), implant will switch from being close to one mouse, to being close to the other mouse. **d**, The head-to-

152    head (nose-centroid-to-nose-centroid) distance for the two mice, across the session. The mice often closely

153    interact (low head-to-head distance), allowing for potential identity swaps. **e**, Distance between implant key-

154    points and the nose centroid for both mice, across the session. The implant key-points are always near mouse0

155    and there are no identity swaps. **f**, The actual implant-key-point to implant-skeleton-point distance for mouse

156    0, across the session, is lower than the distance to the centroid of the nose ellipsoid.
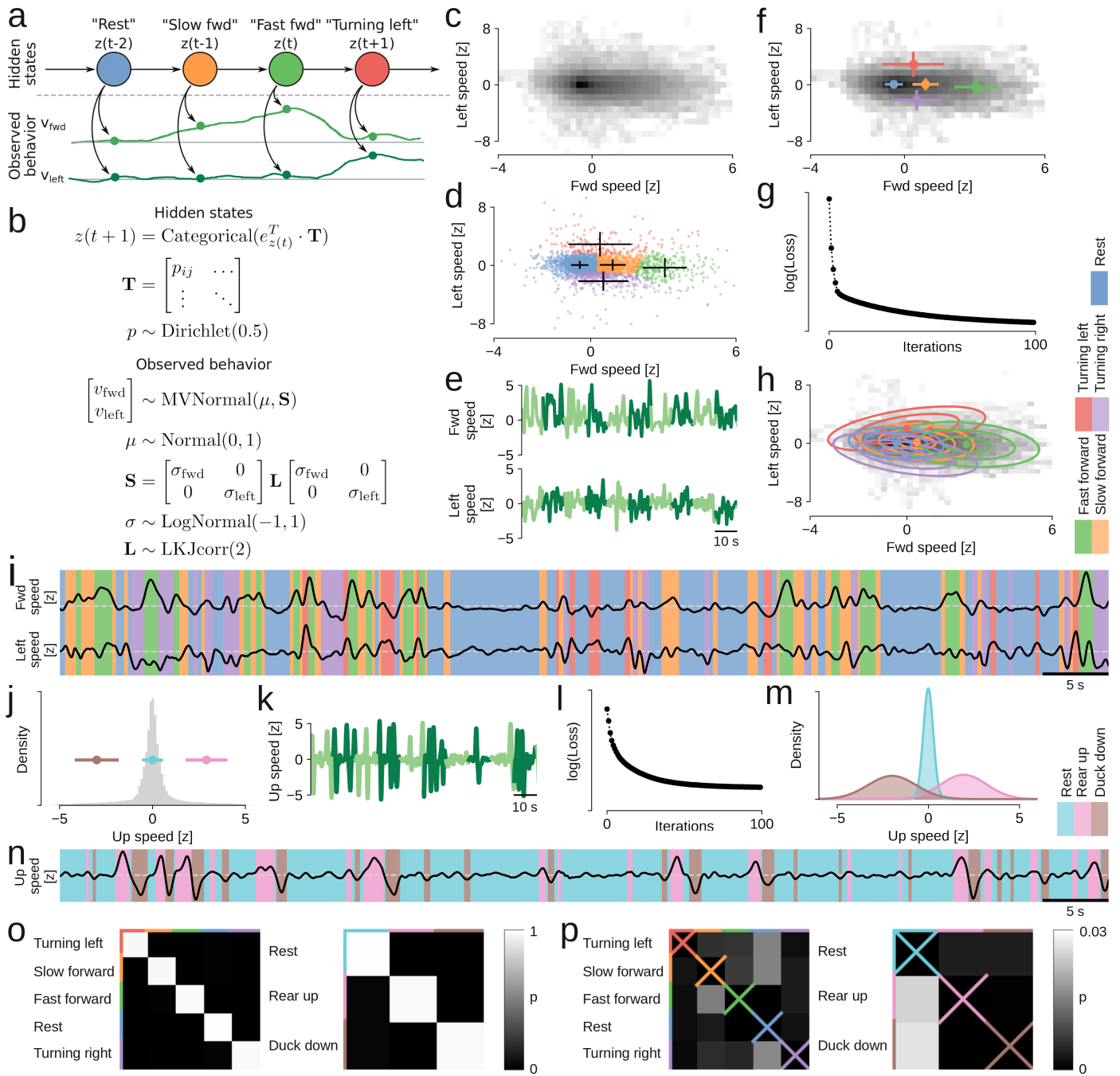
158  **Supplementary Figure 14. Calculation of movement speed in egocentric coordinates. a**, Running be-

159  havior of the two mice (centroid of the hip ellipsoid) across the behavioral session, shown in 2D (top-down

160  view) and 3D. **b**, Running speed decomposed into two components, 'forward speed' (v_fwd, projected onto

161  the orientation of the hip ellipsoid) and 'left speed' (v_orthog, the orthogonal component). **c**, In correct

162  tracking (top row), running bouts will have positive forward speed. If there is a mistake in the tracking (bot-

163  tom row), such that the mouse body model has switched direction, the mouse will appear to be 'running

164   backwards'. **d**, Top to bottom: The x,y,z-coordinates of the position (c_hip) of the mouse at each tracked

165   frame, the change in position between frames, the forward speed, and the left speed. The four rows are re-

166   peated for both mice. There are no direction swap mistakes, and across the whole session, both mice only

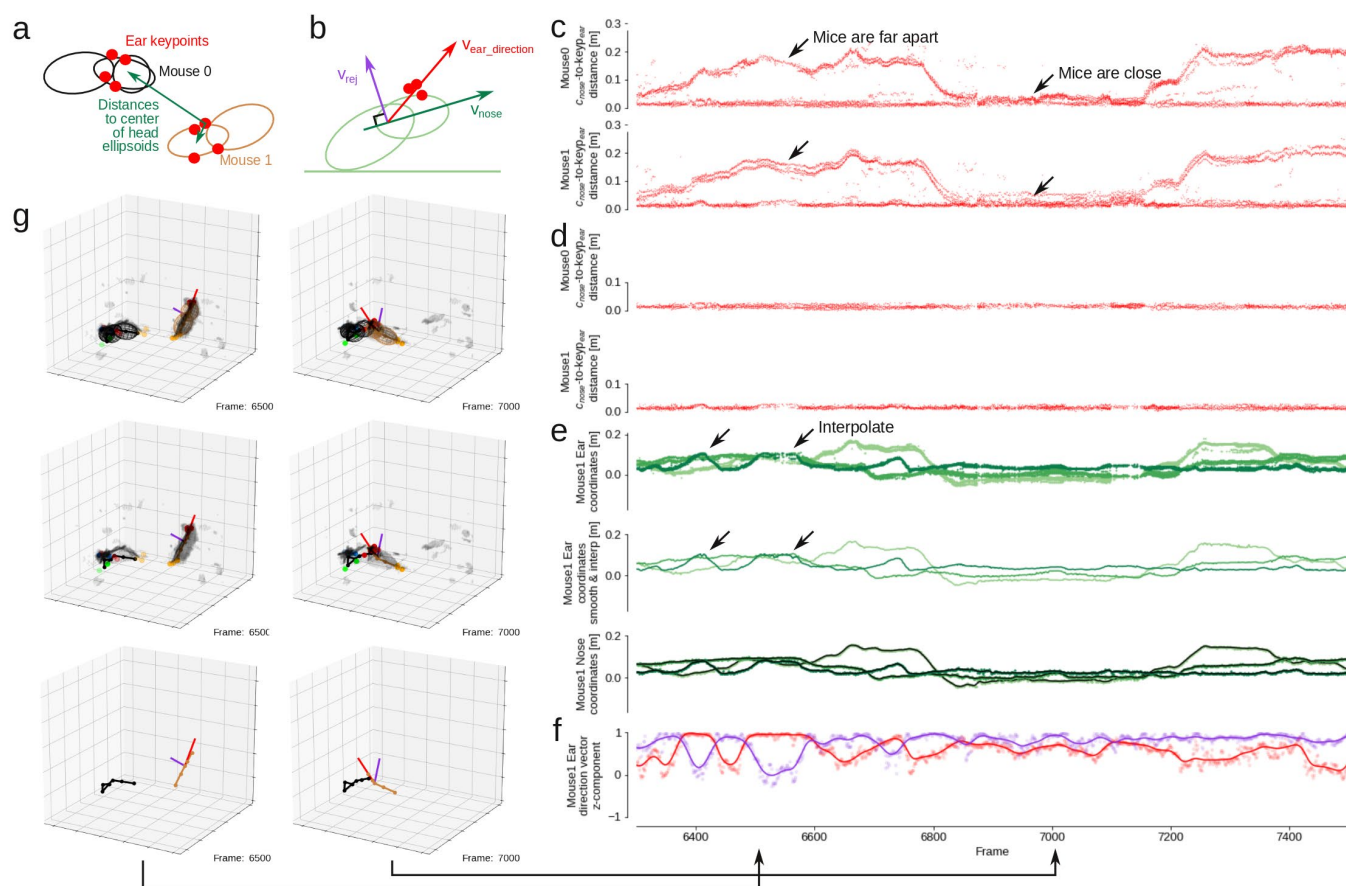167   displayed bouts of forward running (confirmed by visual inspection of raw video).

169 **Supplementary Figure 15. Manual error checking. a**, By manual inspection of 500 frames, we detected

170 one tracking error. **b,** Median point-cloud residual (top) and number of point-cloud points with a residual

171 larger than the cutoff (bottom, cut = 0.03 m) across an example 21 min recording. These traces show two

172 anomalies: One tracking error (around frame 17000, the error we also detected by manual inspection of the

173 500 frames) and one depth camera artifact (tracking was fine, but a ghostly artifact showed up in the point-

174  cloud for few a seconds. Due to the of the robust loss function, tracking was not distorted by the artifact). **c,**

175  Ten example frames showing the tracking error (0.5 s between frames, indicated by vertical lines in panel b).

176  Note that after the error, the particle filter quickly recovers to correct tracking again. **d,** Ten example frames
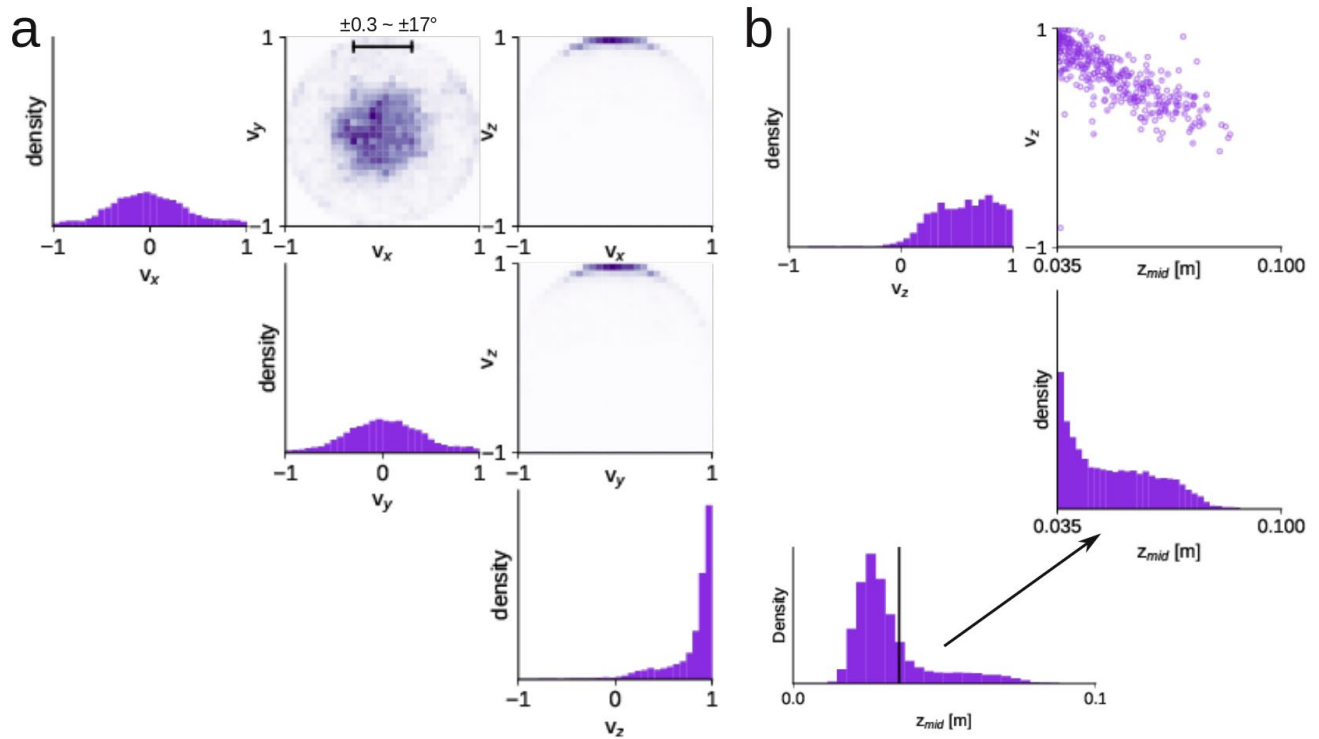
177  showing the depth camera artifact.

**Supplementary Figure 16. Bayesian modeling and automatic classification of behavioral states. a**, Generative model fit to the running behavior to automatically classify behavioral states. The model is a hidden Markov model with discrete latent states (circles), and each state emits a forward speed and a left speed, drawn from a two-dimensional gaussian distribution with a full covariance matrix. **b**, The generative model expressed as equations, showing Bayesian priors for estimating the parameters. **c**, Joint distribution (on a log-scale) of the forward speed and left speed, for both mice, across an entire behavioral session. **d**, Initial position

185  for the variation inference, for the model of forward and left speed. Crosses, cluster centers and standard

186  deviations (calculated independently for fwd/left speed) for clusters assigned by k-means clustering into 5

187  clusters. Dots, individual samples of fwd/left speed (colors indicate clusters, every $50^{th}$ sample is show). **e,**

188  Bayesian model was fitted to a subset of the data (5 mins), split and run in parallel on 10-s sequences. The

189  plot shows example 10-s sequences. **f,** Joint distribution (on a log-scale) of the forward speed and left speed,

190  for the training data, overlaid with the cluster centers and standard deviations from all data (i.e., from **d**).

191  Training data cover same velocity space as the whole session (compare with **a**). **g,** Convergence plot showing

192  the decrease in model loss (increased evidence lower bound) across iterations for the training data. **h,** Loca-

193  tions and covariance ellipsoids (indicating three standard deviations) for the gaussian emission distributions

194  associated with the five latent states, after model fitting. The five clusters are easily interpretable, and the

195  labels are shown on the right. **j,** Initial position for the variation inference for the up speed. Distribution of the

196  up speed (grey bars), as well as the center and standard deviation of three clusters (colored bars and dots),

197  assigned by k-means clustering. **i,** Automatically-assigned states (by maximum a posterior probability) to an

198  example sequence of forward and left speed. **k,l,m,n,** same as **e,g,h,i**, but for the model fitting of the emission

199  gaussians (in one dimension) of the up speed. **o,** Transition probabilities between latent states, for both for-

200  ward/left speed and up speed models. The sample rate is 60 frames/s, so – since behavioral states are longer

201  than that – the self-transition probabilities (diagonals) are very high. **p,** As **o**, but without showing the self-

202  transition probabilities (the diagonals, crossed out). These matrices have understandable structure. For exam-

203  ple, in the left matrix, the most likely transition from "rest" is to "slow forward". From "slow forward", the

204  mouse is likely to transition to "turning left", "fast forward" or "turning right". It is very unlikely to transition

205  directly from "fast forward" to "rest" or from "turning left" to "turning right". From the right matrix, we can

206  see that it is unlikely to transition directly from "rear up" to "rear down", it is more probable to have a period
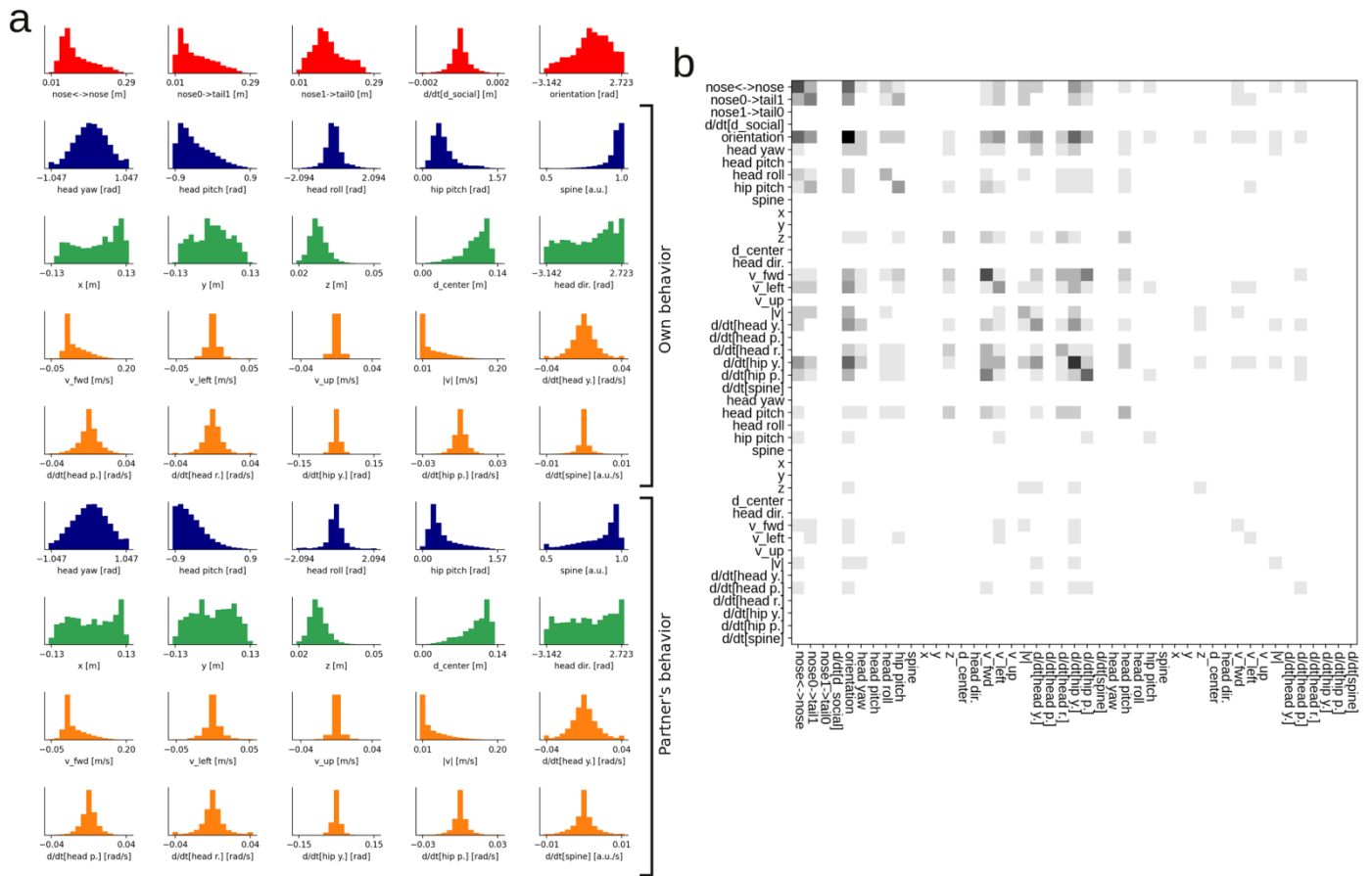
207  of "rest" in between.

**Supplementary Figure 17. Estimation of 3D heading direction in the partner animal, part I. a**, We use the 3D position of the ear keypoints to determine the 3d head direction of the partner animal. We assign the ear keypoints to a mouse body model by calculating the distance from each keypoint to the center of the nose ellipsoid of both animals. **b**, To estimate the 3D head direction, we calculate the unit rejection (v_rej) between a unit vector along the nose ellipsoid (v_nose) and a unit vector from the neck joint (c_mid) to the average 3D position of the ear keypoints that are associated with that mouse (v_ear_direction). **c**, The distance from all ear keypoints to the center of the nose ellipsoid, for both mice, for an example portion of the recording session. **d**, The distance from ear keypoints to the center of the nose ellipsoid, only showing the keypoints that we estimate to be associated with each mouse. **e**, Estimated mean 3D position of the ear keypoints associated with the partner animal ('Mouse 1'). Top to bottom: Raw 3D position of all keypoints, mean position using linear interpolation, smoothed with a Gaussian kernel ($\sigma = 3$ frames). **f**, The z-component of v_ear_direction and v_rej. The z-component is high, indicating that the ears are on the dorsal side of the head ellipsoid.

221  When the mouse is running on the ground, both v_ear_direction and v_rej have high z-components (marked

222  with rightmost arrow), but when the mouse is rearing and tilting the head backwards, v_rej will be more in

223  the xy-plane, and have a low z-component (marked with leftmost arrow). **g**, The 3D body positions, of the

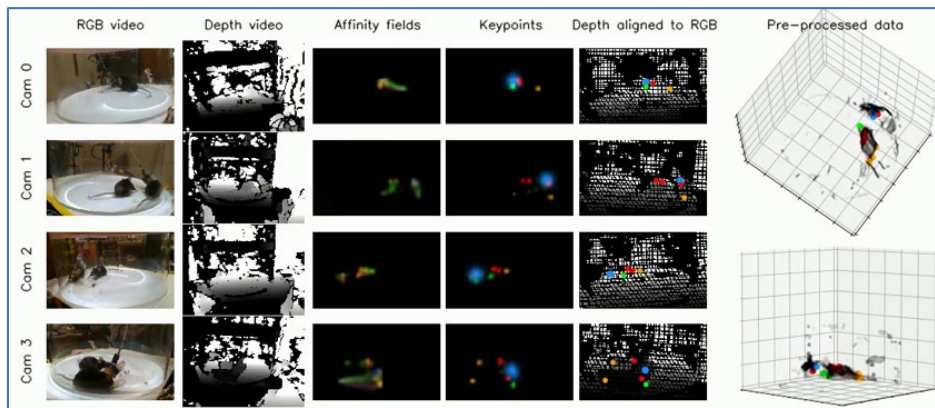224  frames indicated by arrows in panel **f**.

226   **Supplementary Figure 18. Estimation of 3D heading direction in the partner animal, part II. a**, The joint

227   distributions of the components of v_rej shows that mouse mostly keeps the ears horizontal, rarely tilting the

228   head more than 17 degrees towards the left or right. The z-component is mostly close to 1 (pointing straight

229   up), but sometimes smaller, closer to 0 (meaning that the nose is pointing up towards the sky). **b**, We can

230   examine the details of the 3D head direction behavior. For example, we can monitor the head direction, when

231   the z-coordinate of the neck (z_mid) is high (i.e., when the mouse is rearing). Here we find a clear negative

232   correlation between the z-component of v_rej and z_mid, which matches the visual inspection of the videos:

233   When the mouse rears up or climbs up against the walls of the transparent social arena, the head tilts back to

234   extend the nose upwards.

236 **Supplementary Figure 19. Details of 3D social behavior. a**, Distribution of the behavior occupancy in the

237 bins of the GLM model, for both the social features, the 'own body' features and the 'partner body' features.

238 **b,** The 'co-encoding matrix' of the neural population: The grayscale color in i'th and j'th bin in the heatmap

239 indicates the number neurons that encode both feature i and j, shown here with the full variable names on
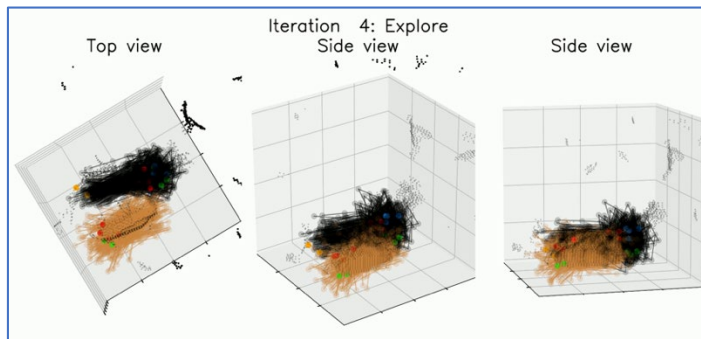
240 the matrix axes.

## Supplementary Videos and Legends
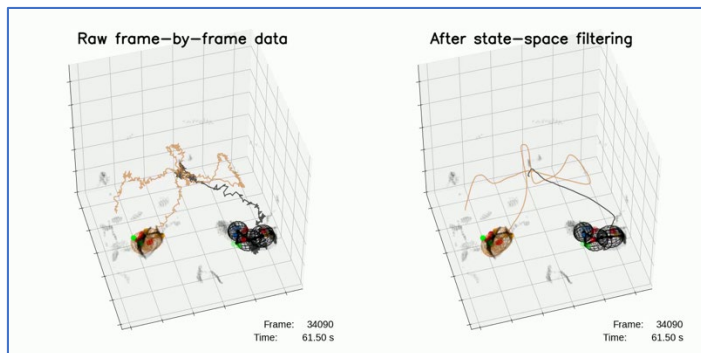


**Supplementary Video 1. Pre-processing pipeline.**



**Supplementary Video 2. Particle filter behavior.**



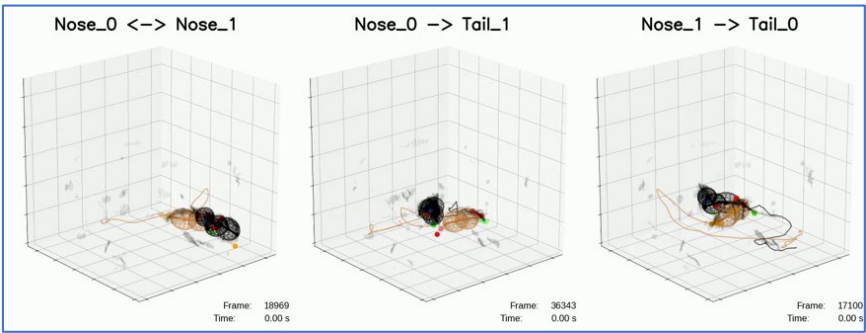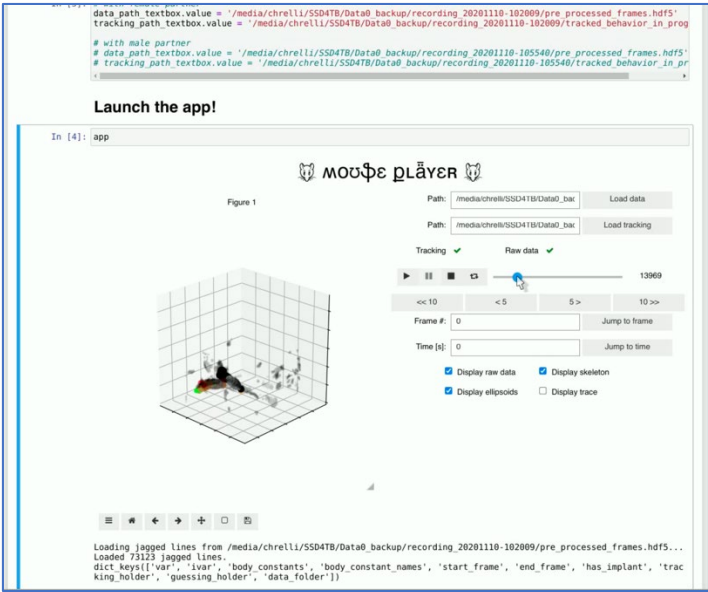**Supplementary Video 3. State-space filtering.**

255 **Supplementary Video 4. Social events.**

256

257



259 <mark>**Supplementary Video 5. MousePlayer.**</mark>