# Supplementary Figures and Movie Legends

**Title:** Fast animal pose estimation using deep neural networks

**Authors:**
Pereira, T. D.[1,#], Aldarondo, D. E.[1,#], Willmore, L.[1], Kislin, M.[1], Wang, S. S.-H.[1,2], Murthy, M.[1,2*], and Shaevitz, J. W.[1,3,4*]

1 Princeton Neuroscience Institute, Princeton University

2 Department of Molecular Biology, Princeton University

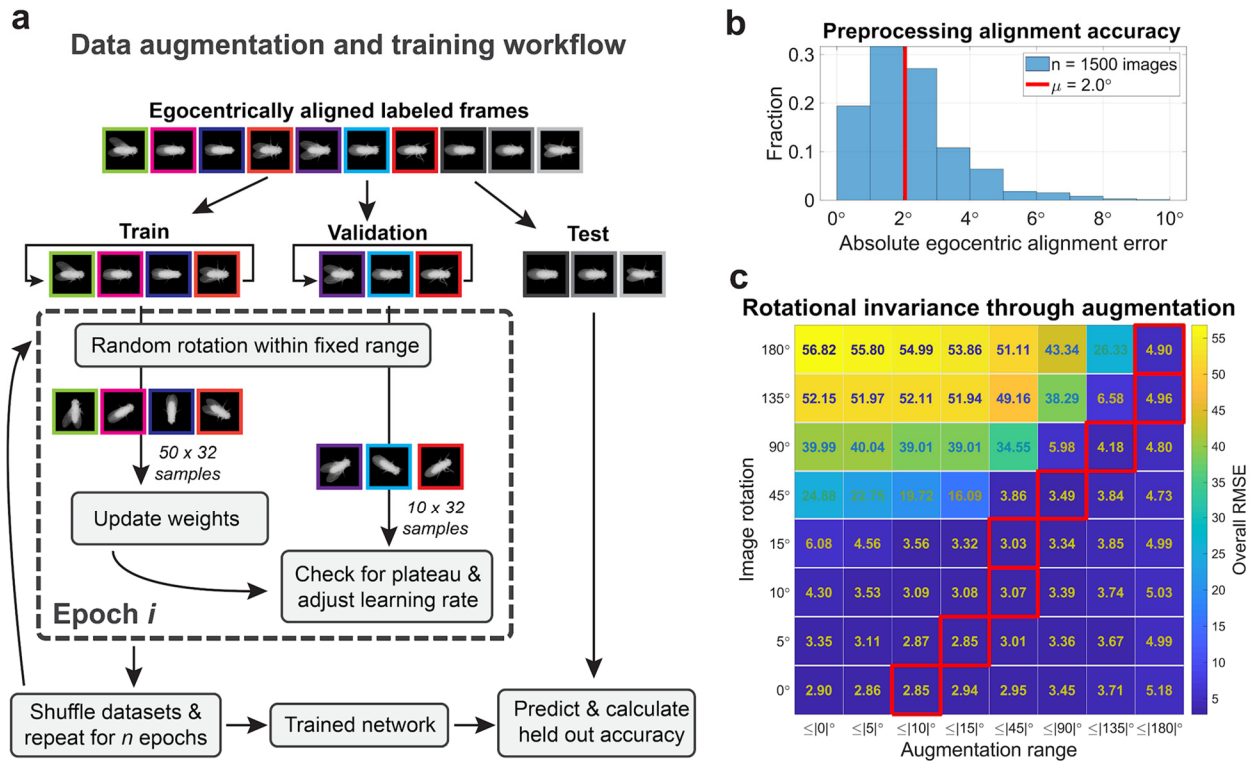3 Lewis-Sigler Institute for Integrative Genomics, Princeton University

4 Department of Physics, Princeton University

#equal authors
*lead contacts and co-corresponding authors: Mala Murthy (mmurthy@princeton.edu) and
Joshua W. Shaevitz (shaevitz@princeton.edu)

**a**

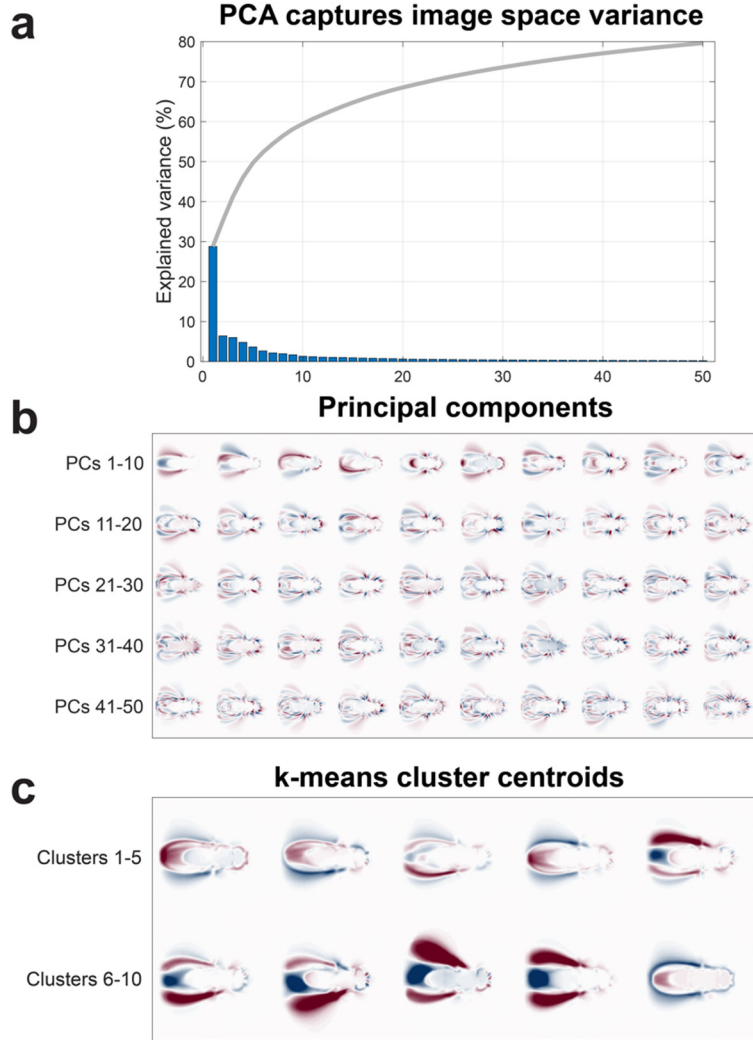**Data augmentation and training workflow**

**Egocentrically aligned labeled frames**

Train | Validation | Test

Random rotation within fixed range

*50 x 32 samples*

Update weights

*10 x 32 samples*

Check for plateau & adjust learning rate

**Epoch *i***

Shuffle datasets & repeat for *n* epochs → Trained network → Predict & calculate held out accuracy

**b**

**Preprocessing alignment accuracy**

Fraction

n = 1500 images
μ = 2.0°

0° 2° 4° 6° 8° 10°
Absolute egocentric alignment error

**c**

**Rotational invariance through augmentation**

| Image rotation | ≤\|0\|° | ≤\|5\|° | ≤\|10\|° | ≤\|15\|° | ≤\|45\|° | ≤\|90\|° | ≤\|135\|° | ≤\|180\|° |
|---|---|---|---|---|---|---|---|---|
| 180° | 56.82 | 55.80 | 54.99 | 53.86 | 51.11 | 43.34 | 26.33 | 4.90 |
| 135° | 52.15 | 51.97 | 52.11 | 51.94 | 49.16 | 38.29 | 6.58 | 4.96 |
| 90° | 39.99 | 40.04 | 39.01 | 39.01 | 34.55 | 5.98 | 4.18 | 4.80 |
| 45° | 24.88 | 22.78 | 19.72 | 16.09 | 3.86 | 3.49 | 3.84 | 4.73 |
| 15° | 6.08 | 4.56 | 3.56 | 3.32 | 3.03 | 3.34 | 3.85 | 4.99 |
| 10° | 4.30 | 3.53 | 3.09 | 3.08 | 3.07 | 3.39 | 3.74 | 5.03 |
| 5° | 3.35 | 3.11 | 2.87 | 2.85 | 3.01 | 3.36 | 3.67 | 4.99 |
| 0° | 2.90 | 2.86 | 2.85 | 2.94 | 2.95 | 3.45 | 3.71 | 5.18 |

Augmentation range

Overall RMSE

**Supplementary Figure 1: Rotational invariance is learned at the cost of prediction accuracy**

(a) The augmentation procedure consists of random rotations about the center of egocentrically aligned labeled frames. Labeled frames are split into training, validation and test sets. Colors are used to indicate unique images. Only training and validation sets are augmented and used for training.. During training, images are drawn sequentially from the training and validation sets to form batches of 32 images, cycling back to the beginning if there are less images than required, and then rotated randomly within a range of angles; confidence maps are rotated accordingly (not shown). After each epoch, the ordering of the datasets are shuffled so as to create new combinations of batches. The test set images are not augmented before computing accuracy metrics reported throughout.

(b) Egocentric alignment accuracy of the preprocessing algorithm from [1] when compared to manual labels of head/thorax. The error is the absolute deviation of the angle formed between the thorax and head from the horizontal centerline in the image. The mean of 2.0° indicates that there is little alignment error to which the network has to learn robustness.

(c) The accuracy measured as the RMSE of position estimates when evaluated on data artificially rotated at a fixed angle (rows) with networks trained on data augmented by rotations between a range of angles (columns). Red boxes denote the best accuracy for each data angle, denoting that optimal performance is achieved when the network is trained on augmented images with the minimally inclusive range of angles. Top accuracy decreases relative to the degree of rotational invariance the network must learn.
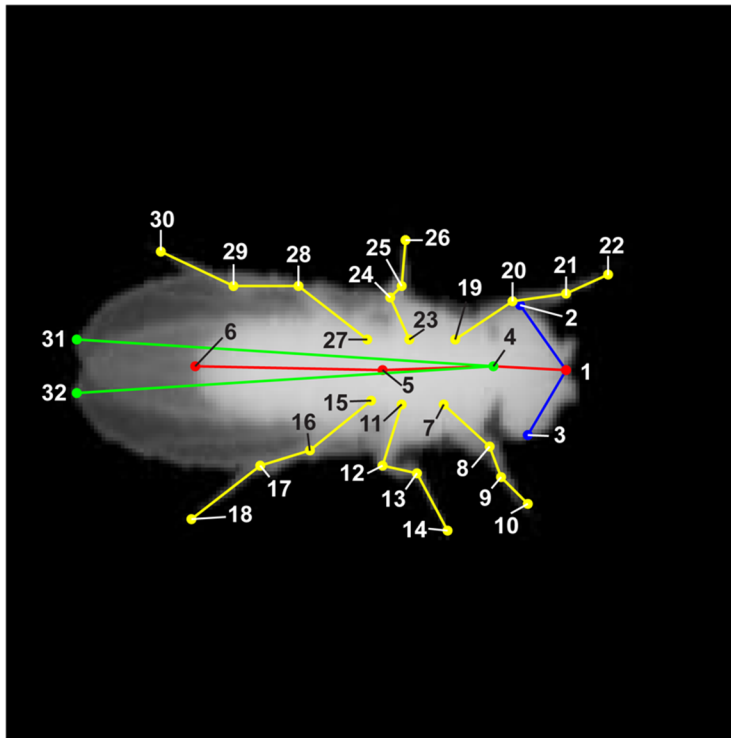
**a** PCA captures image space variance

**b** Principal components

PCs 1-10
PCs 11-20
PCs 21-30
PCs 31-40
PCs 41-50

**c** k-means cluster centroids

Clusters 1-5
Clusters 6-10

**Supplementary Figure 2: Cluster sampling to promote pose diversity in labeling dataset**
(a) Principal component analysis (PCA) of unlabeled images captures the majority of the
variance in the data within 50 components. The cumulative variance explained (line) suggests
that using PCA for dimensionality reduction does not sacrifice substantial information within the
images.
(b) Top PCA eigenmodes visualized as coefficient images. Red and blue shading denote
positive and negative coefficients at each pixel. Areas of similar colors indicate correlated pixel
intensities within a given mode. After mean subtraction, each image in the initially sampled
dataset is projected onto the first 50 eigenmodes.
(c) Cluster centroids identified by k-means after PCA. Red and blue shading denote pixels with
higher or lower intensity than the overall mean. Cluster centroids illustrate the diversity of poses
that are detected in image space by this sampling method. Samples are then drawn evenly from
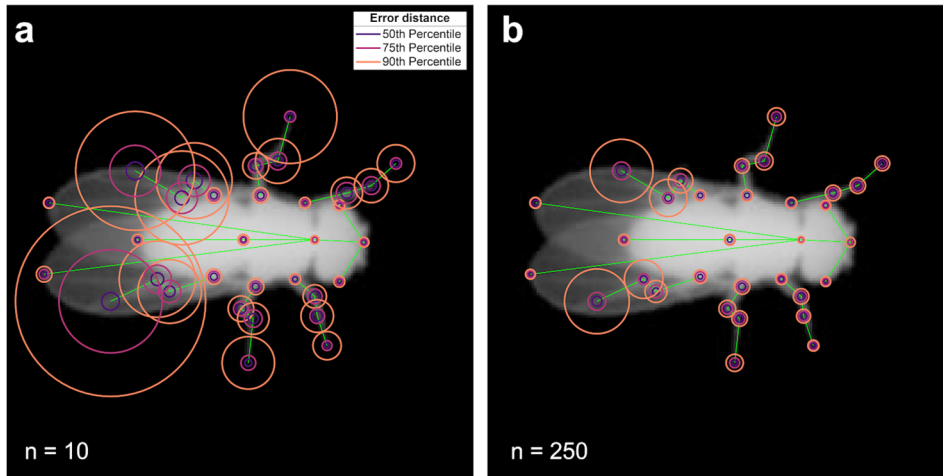each cluster to select representative images for labeling with the GUI.

(1) Tip of head (between antennae)
(2) Left eye
(3) Right eye
(4) Neck / connects head and thorax
(5) Mesothoracic phragma (connects thorax and abdomen)
(6) Tip of abdomen
(7) Right foreleg / thorax-coxa
(8) Right foreleg / coxa-femur
(9) Right foreleg / femur-tibia
(10) Right foreleg / tarsus tip
(11) Right midleg / thorax-coxa
(12) Right midleg / coxa-femur
(13) Right midleg / femur-tibia
(14) Right midleg / tarsus tip
(15) Right hindleg / thorax-coxa
(16) Right hindleg / coxa-femur
(17) Right hindleg / femur-tibia
(18) Right hindleg / tarsus tip
(19) Left foreleg / thorax-coxa
(20) Left foreleg / coxa-femur
(21) Left foreleg / femur-tibia
(22) Left foreleg / tarsus tip
(23) Left midleg / thorax-coxa
(24) Left midleg / coxa-femur
(25) Left midleg / femur-tibia
(26) Left midleg / tarsus tip
(27) Left hindleg / thorax-coxa
(28) Left hindleg / coxa-femur
(29) Left hindleg / femur-tibia
(30) Left hindleg / tarsus tip
(31) Left wing tip
(32) Right wing tip

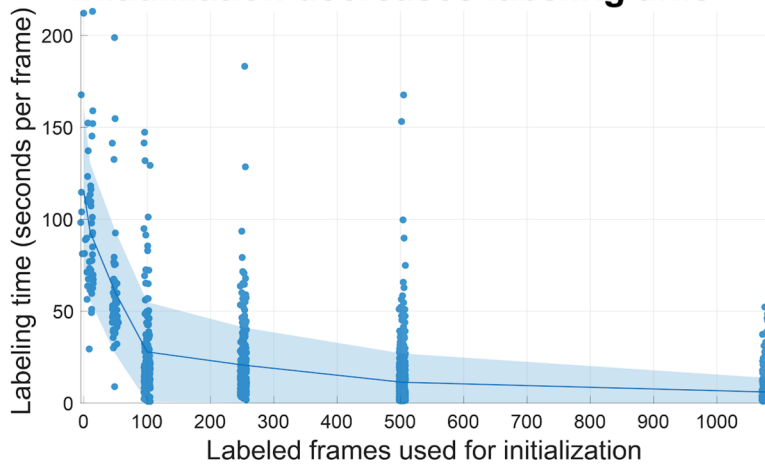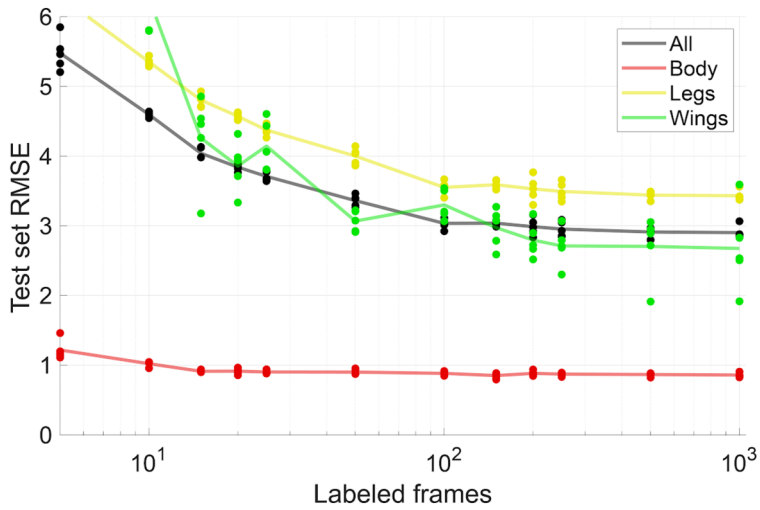**Supplementary Figure 3: User-defined skeleton**

We selected 32 points to cover the body parts of the fly; these parts were chosen to approximately match the set of visible joints and interest points in the anatomy of the animal.

**a**

Error distance
— 50th Percentile
— 75th Percentile
— 90th Percentile

n = 10

**b**

n = 250

**c** **Initialization decreases labeling time**



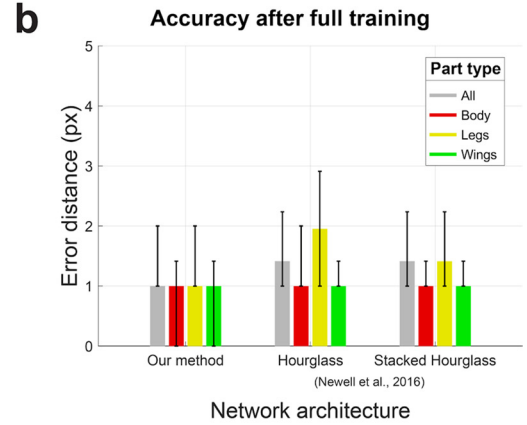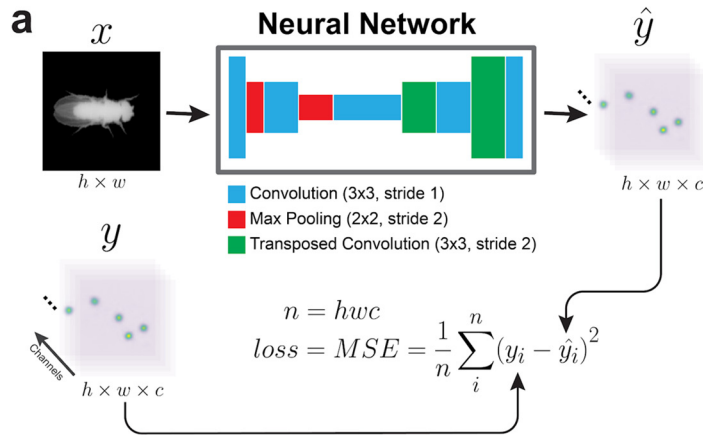**d** **Accuracy improves quickly with few samples**



65
66
67 **Supplementary Figure 4: Estimation accuracy improves with few samples**
68 (a-b) Error distance distributions per body part when estimated with networks trained for 15

5

epochs on 10 (a) or  250 (b) labeled frames. The majority of estimates fall within few pixels of the ground truth, reducing the labeling procedure to simply correcting estimates.

(c) Time spent labeling each frame decreases with the quality of initialization. Line and shaded region correspond to mean and standard deviation. Starting frames require 115.4+-45.0 (mean+-s.d.) seconds to label, decreasing to 6.1±7.7 seconds after initializing with a network trained on 1000 labeled frames.

(d) Large accuracy improvements are observed with very few labeled samples, corresponding with the decrease in time required to fix initial labels on new frames. A plateau is observed at around 150-200 frames, with marginal improvements with additional labeling. Circles denote the test set RMSE for one replicate of fast training (15 epochs) at each dataset size, lines denote mean of all replicates.
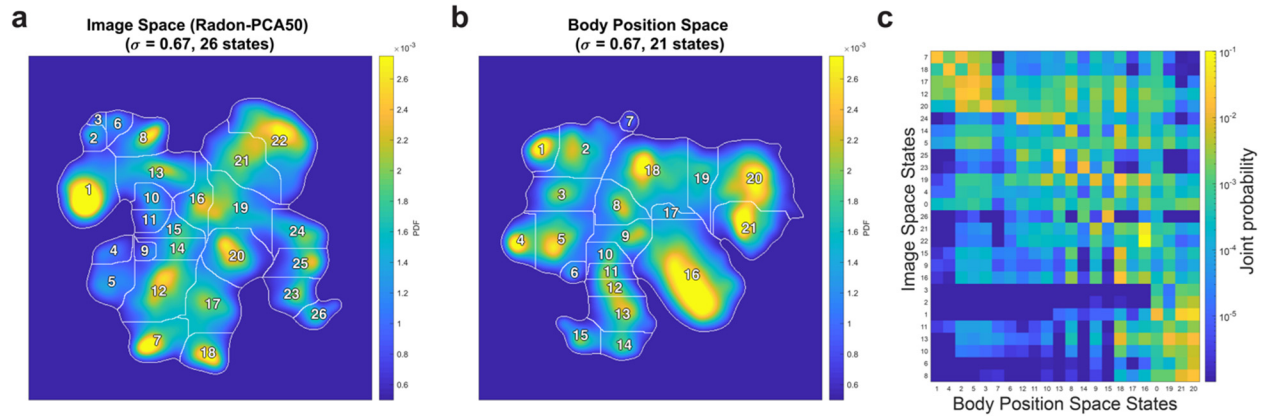
**a** $x$ Neural Network $\hat{y}$

$h \times w$

Convolution (3x3, stride 1)
Max Pooling (2x2, stride 2)
Transposed Convolution (3x3, stride 2)

$h \times w \times c$

$y$

Channels

$h \times w \times c$

$n = hwc$

$loss = MSE = \frac{1}{n}\sum_{i}^{n}(y_i - \hat{y}_i)^2$

**b** Accuracy after full training

Part type
— All
— Body
— Legs
— Wings

Error distance (px)

Our method    Hourglass    Stacked Hourglass
             (Newell et al., 2016)

Network architecture

**Supplementary Figure 5: Neural network architecture comparison**
(a) Diagram of our neural network architecture. Raw images are provided as input into the
network, which then computes a set of confidence maps of the same height and width as the
input image (top row). The network consists of a set of convolutions, max pooling and
transposed convolutions whose weights are learned during training (top middle). Estimated
confidence maps are compared to ground truth maps generated from user labels using a mean
squared error loss function, which is then minimized during training (bottom row).
(b) Accuracy comparison between architectures. We compared the accuracy of our architecture
to the hourglass and stacked hourglass versions of the network described in[2]. The accuracy of
our network is equivalent or better than those achieved when training with these reference
architectures (over all body parts, $p < 1e-10$, Wilcoxon rank sum test, 1-tailed). Bar and error
bars denote median and 25th and 75th percentiles.

97
98
99  **Supplementary Figure 6: Comparison of behavioral space distributions generated from**
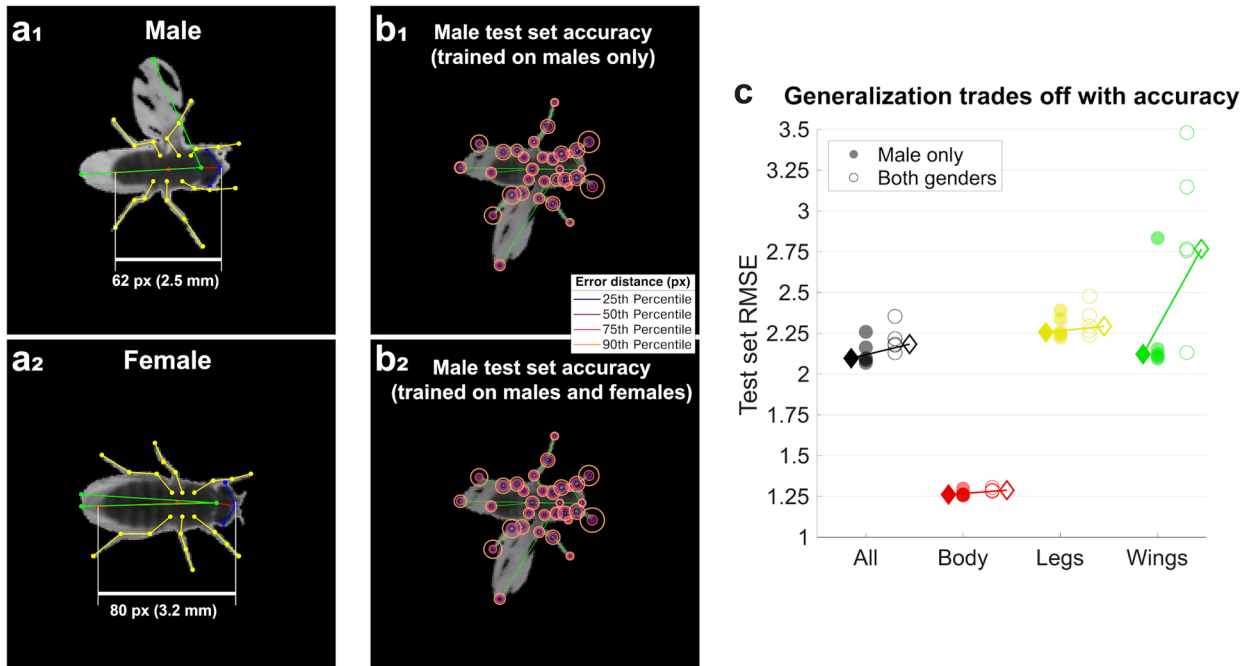100 **compressed images versus body part positions.**
101 (a) Behavioral space distribution from 59 male flies calculated using the original MotionMapper
102 pipeline (data and pipeline from [3]), including Radon-transform compression and PCA-based
103 projection onto the first 50 principal components followed by a nonlinear embedding of the
104 resultant spectrograms.
105 (b) Behavioral space distribution from 59 male flies (data and pipeline from [3]) calculated using
106 spectrograms generated from tracked body part positions rather than PCA modes (see **Online**
107 **Methods**). We note that this distribution has fewer peaks than that from (a) and a more
108 symmetric topology (e.g in the top-left clusters, **Fig. 4c-g**).
109 (c) Joint probability distribution of the cluster labels from (a) and (b); sorted by row and column
110 peaks. Many clusters identified using the pixel-based representation (rows) match up with those
111 of the position-based representation (columns), but some are distributed into newly separated
112 clusters.
113
114
115
116
117
118
119
120
121
122
123

124

**Supplementary Figure 7: Generalization to more diverse morphologies with a single network trades off with accuracy.**

(a) Male and female flies differ in anatomical morphology, in part due to differences in their body length. The males ($a_1$) more often extend their wings as they are used to produce courtship song. The females ($a_2$) rarely extend their wings in this context, resulting in different requirements for pose estimation between the two genders, despite their overall similarity in morphology.

(b) Training on labeled images of just males ($b_1$) results in similar accuracy (on male test set images) to when training on both males and females ($b_2$). This suggests that there is little discernible difference (up to the 90th percentile) of having a network trained on two different types of body morphologies.
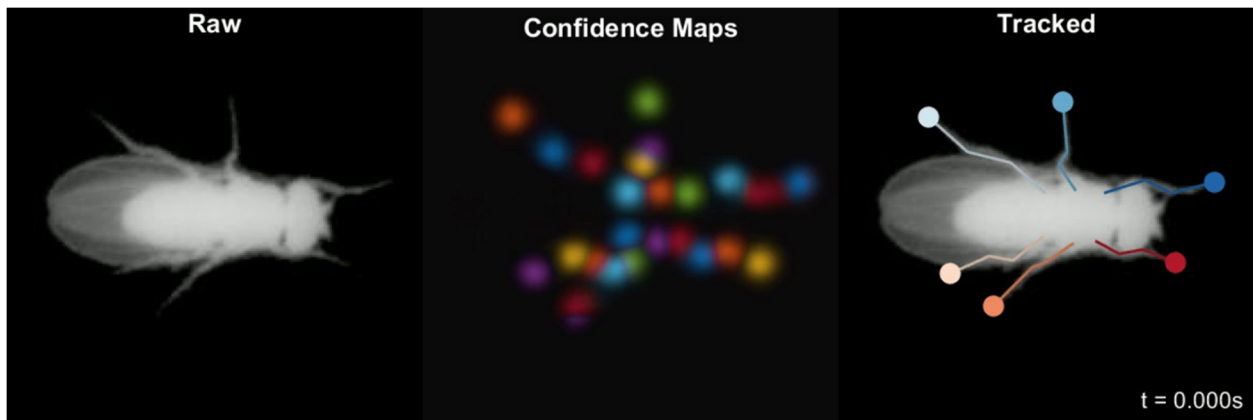
(c) Quantification of RMSE on the male test set shows that generalization to two different morphologies increases the error metric. Circles denote training replicates, diamonds denote median RMSE for all replicates, and filled and empty markers correspond to specialized versus generalized training respectively. Although the increased error rate is very small overall when generalizing, the greatest difference is observed in the body parts with greater difference in pose distributions (wings, green).

142

**Supplementary Movie 1: Body part tracking is reliable over long periods without temporal constraints.**

Raw images (left), max projection of all confidence maps (center), and tracked images (right) during a 20 second bout of free movement. Video playback at 0.2x realtime speed.
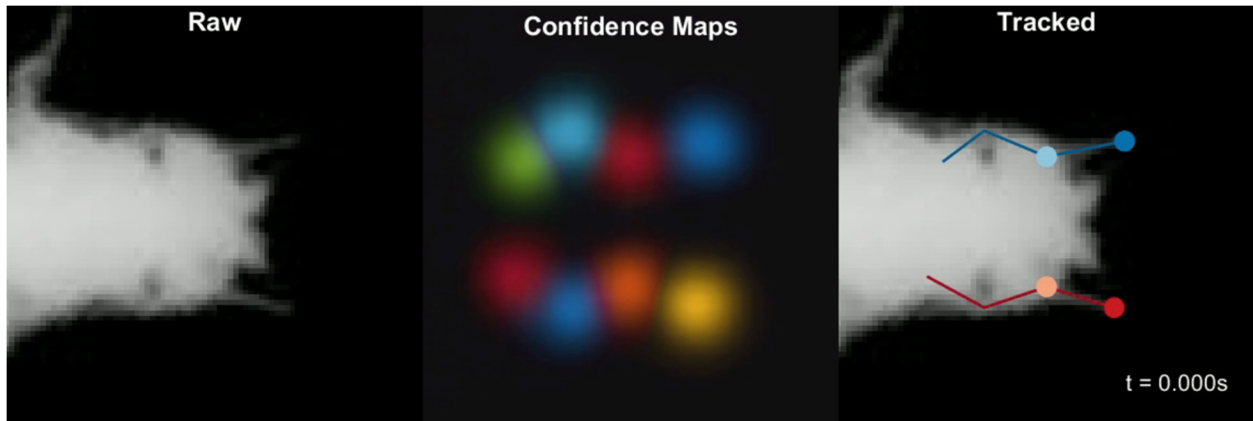


**Supplementary Movie 2: Body part tracking during freely moving locomotion.**

Raw images (left), max projection of all confidence maps (center), and tracked images (right) during a bout of locomotion. Video playback at 0.15x realtime speed. Video corresponds to Fig. 1d.

167
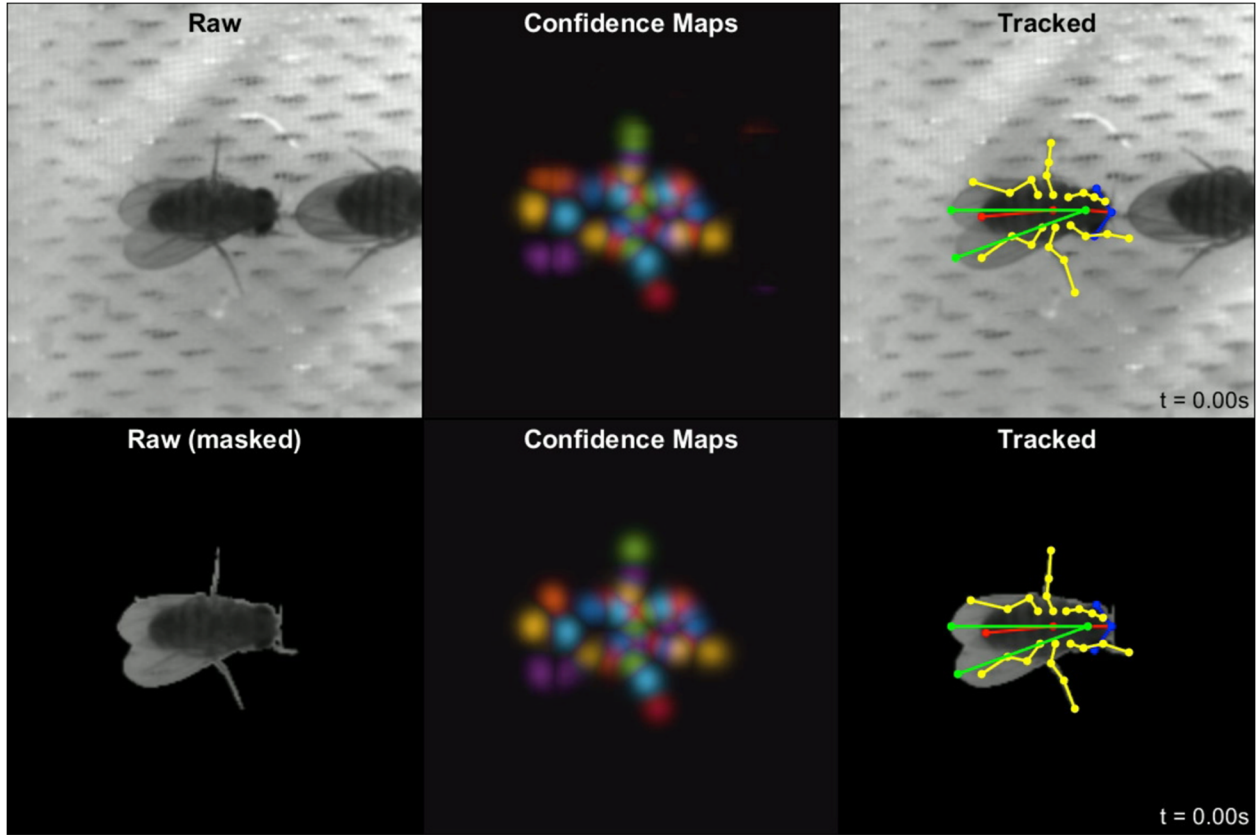168
169
170
171
172
173
174
175
176



177
**Supplementary Movie 3: Body part tracking during head grooming.**
Raw images (left), max projection of all confidence maps (center), and tracked images (right)
during a bout of head grooming. Video playback at 0.15x realtime speed. Video corresponds to
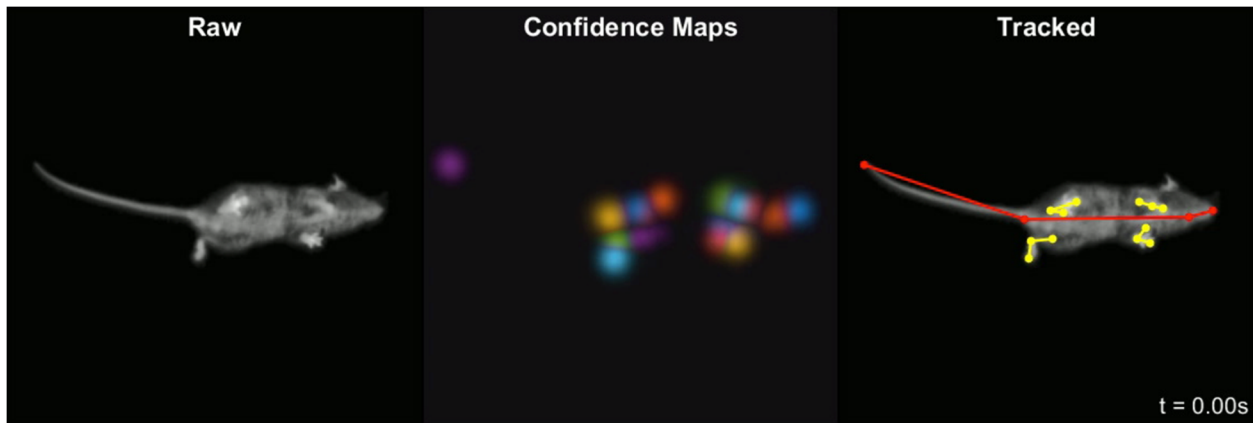Fig. 1e.

182
183
184

185
186
187
188
189
190



191
192
193 **Supplementary Movie 4: Tracking joints robustly in images with heterogeneous**
194 **background and noisy segmentation.**
195 Raw images (left), max projection of all confidence maps (center), and tracked images (right) of
196 a freely moving courting male fly. Rows correspond to results from a network trained on
197 unmasked and masked images, respectively. Video playback at 0.2x realtime speed.
198

199



200
201
202 **Supplementary Movie 5: Tracking joints in freely moving rodents.**
203 Raw images (left), max projection of all confidence maps (center), and tracked images (right) of
204 a freely moving mouse in an open field arena imaged from below through a clear acrylic floor.
205 Video playback at 0.2x realtime speed. Tracking is reliable over time but degenerate when
206 certain parts are occluded, such as when the animal rears.
207
208
209
210 <u>**References**</u>

211　1.　Berman, G. J., Bialek, W. & Shaevitz, J. W. Predictability and hierarchy in Drosophila

212　　　behavior. *Proc. Natl. Acad. Sci. U. S. A.* **113,** 11943–11948 (2016).

213　2.　Newell, A., Yang, K. & Deng, J. Stacked Hourglass Networks for Human Pose Estimation.

214　　　in *Computer Vision – ECCV 2016* 483–499 (Springer International Publishing, 2016).

215　3.　Berman, G. J., Choi, D. M., Bialek, W. & Shaevitz, J. W. Mapping the stereotyped behaviour

216　　　of freely moving fruit flies. *J. R. Soc. Interface* **11,** (2014).

217