

1        **Selfee: Self-supervised Features Extraction of animal behaviors**

2        **Yinjun Jia<sup>1,2,\*</sup>, Shuaishuai Li<sup>3,4</sup>, Xuan Guo<sup>1,2</sup>, Junqiang Hu<sup>1</sup>, Xiao-Hong Xu<sup>3,4</sup>, Wei**  
3        **Zhang<sup>1,2,\*</sup>**

4        <sup>1</sup>School of Life Sciences, IDG/McGovern Institute for Brain Research, Tsinghua  
5        University, Beijing, 100084, China.

6        <sup>2</sup>Tsinghua-Peking Center for Life Sciences, Beijing, 100084, China.

7        <sup>3</sup>Institute of Neuroscience, State Key Laboratory of Neuroscience, Chinese Academy  
8        of Sciences Center for Excellence in Brain Science and Intelligence Technology,  
9        Shanghai, 200031, China.

10       <sup>4</sup>Shanghai Center for Brain Science and Brain-Inspired Intelligence Technology,  
11       Shanghai, 200031, China.

12       \*Correspondence and requests for materials should be addressed to W.Z.  
13       ([wei\\_zhang@mail.tsinghua.edu.cn](mailto:wei_zhang@mail.tsinghua.edu.cn)) and Y.J. ([jyj20@mails.tsinghua.edu.cn](mailto:jyj20@mails.tsinghua.edu.cn))

## 14 **ABSTRACT**

15 Fast and accurately characterizing animal behaviors is crucial for neuroscience research.  
16 Deep learning models are efficiently used in laboratories for behavior analysis.  
17 However, it has not been achieved to use a fully unsupervised method to extract  
18 comprehensive and discriminative features directly from raw behavior video frames for  
19 annotation and analysis purposes. Here, we report a self-supervised feature extraction  
20 (Selfee) convolutional neural network with multiple downstream applications to  
21 process video frames of animal behavior in an end-to-end way. Visualization and  
22 classification of the extracted features (Meta-representations) validate that Selfee  
23 processes animal behaviors in a comparable way of human understanding. We  
24 demonstrate that Meta-representations can be efficiently used to detect anomalous  
25 behaviors that are indiscernible to human observation and hint in-depth analysis.  
26 Furthermore, time-series analyses of Meta-representations reveal the temporal  
27 dynamics of animal behaviors. In conclusion, we present a self-supervised learning  
28 approach to extract comprehensive and discriminative features directly from raw video  
29 recordings of animal behaviors and demonstrate its potential usage for various  
30 downstream applications.

31

## 32 INTRODUCTION

33 Extracting representative features of animal behaviors has long been an important  
34 strategy to study the relationship between genes, neural circuits, and behaviors.  
35 Traditionally, human observations and descriptions are the primary solutions for animal  
36 behavior analysis<sup>1</sup>. Well-trained researchers would define a set of behavior patterns and  
37 compare their intensity or proportion between experimental and control groups. With  
38 the emergence and thrive of machine learning methodology, supervised learning has  
39 been assisting human annotations and achieved impressive results<sup>2-4</sup>. Nevertheless,  
40 supervised learning is limited by prior knowledge and manually assigned labels, thus  
41 could not identify behavioral features that are not annotated.

42 Other machine learning methods were then introduced to the field which were designed  
43 to extract representative features beyond human-defined labels. These methods can be  
44 generally divided into two major categories: one estimates animal postures with a group  
45 of pre-defined key points of the body parts, and the other directly transforms raw images.  
46 The former category marks representative key points of animal bodies, including limbs,  
47 joints, trunks, and/or other body parts of interest<sup>5-7</sup>. Those features are usually sufficient  
48 to represent animal behaviors. However, it has been demonstrated that the key points  
49 generated by pose estimation are less efficient for direct behavior classification or two-  
50 dimensional visualization<sup>8,9</sup>. Sophisticated post-processing like recurrent neural  
51 networks (RNNs)<sup>8</sup>, non-locomotor movement decomposition<sup>10</sup>, or feature  
52 engineering<sup>9</sup> can be applied to transform the key points into higher-level

53 discriminative features. Additionally, neglected body parts could be catastrophic. For  
54 example, the position of the proboscis of a fly is commonly neglected in behavior  
55 studies<sup>9,11</sup>. Still, it is crucial for feeding<sup>12</sup>, licking behavior during courtship<sup>13</sup>, and  
56 hardness detection for a substrate<sup>14</sup>. Finally, best to our knowledge, there is no  
57 demonstration of these pose-estimation methods applied to multiple animals of the  
58 same color with intensive interactions. Thus, the application of pose-estimation to  
59 mating behaviors of two black mice, a broadly adopted behavior paradigm<sup>15-17</sup>, could be  
60 limited because labeling body parts during mice mounting is challenging even for  
61 humans (see Discussion for more details). Therefore, using these feature extraction  
62 methods requires rigorously controlled experimental settings, additional feature  
63 engineering, and considerable prior knowledge of particular behaviors.

64 In contrast, the other category transforms pixel-level information, thus retaining more  
65 details and requiring less prior knowledge. Feature extraction of images could be  
66 achieved by wavelet transforms<sup>18</sup> or Radon transforms<sup>19</sup> followed by principal  
67 component analysis (PCA), and these transforms can be applied to either 2D images or  
68 depth images. However, preprocessing such as segmentation and/or registration of the  
69 images is usually required to achieve spatial invariance, a task that is particularly  
70 difficult for multi-agent videos. Additionally, these methods usually use fixed  
71 transforms and could not be adapted to different behaviors. Flourished deep learning  
72 methods, especially convolutional neural networks<sup>20</sup> (CNNs), could be adaptive to  
73 extract features from diversified datasets. Also, they have been proven more potent than

74 classic computer vision algorithms like wavelet transforms<sup>21</sup> and Radon transforms<sup>22</sup>  
75 on a famous grayscale dataset MNIST, even without supervising<sup>23</sup>. Therefore, we  
76 attempt to adopt CNNs to achieve end-to-end feature extractions animal behaviors that  
77 are comprehensive and discriminative.

78 The cutting-edge self-supervised deep learning methods aim to extract representative  
79 features for downstream missions by comparing different augmentations of the same  
80 image and/or different images<sup>24-28</sup>. Compared with previous techniques, these methods  
81 have three major advantages. Firstly, self-supervised or unsupervised methods could  
82 completely avoid human biases. Secondly, the augmentations used to create positive  
83 samples promise invariance of the neural networks to object sizes, spatial orientations,  
84 and ambient laminations so that registration or other preprocessing is not required.  
85 Finally, the networks are optimized to export similar results for positive samples and  
86 separate negative ones, such that the extracted features are inherently discriminative.  
87 Even without negative samples, the networks can utilize differential information within  
88 batches to obtain remarkable results on downstream missions like classification or  
89 image segmentation<sup>27,29,30</sup>. These advances in self-supervised learning provide a  
90 promising way to analyze animal behaviors.

91 In this work, we develop Selfee (**S**elf-supervised **F**eatures **E**xtraction) that adopts  
92 cutting-edge self-supervised learning algorithms and CNNs to analyze animal  
93 behaviors. Selfee is trained on massive unlabeled behavior video frames to avoid human  
94 bias on annotating animal behaviors, and it could capture a global character of animal

95 behaviors even when detailed postures are hard to see, just like human observation.  
96 During the training process, Selfee learns to project images to a low-dimensional space  
97 without being affected by shooting conditions, image translation, and rotation, where  
98 cosine distance is proper to measure the similarities of original pictures. Selfee also  
99 provides potentials for various downstream analyses. We demonstrate that the extracted  
100 features are suitable for t-SNE visualization, k-NN-based classification, k-NN-based  
101 anomaly detection, and dynamic time warping (DTW). We also show that further  
102 integrated modeling, like the autoregressive hidden Markov model (AR-HMM), is  
103 compatible with Selfee extracted Meta-representations. After downstream analyses,  
104 Selfee provides comparable results with manual annotations on fly behavior like  
105 courtship index. We apply Selfee to fruit flies, mice, and rats, three widely used model  
106 animals, and validate our results with manual annotations. Discoveries of behavioral  
107 phenotypes in mutant flies by Selfee are proven to have biological significance. The  
108 performance of Selfee on these model species indicates its potential usage for  
109 behavioral studies of non-model animals as well as other tasks. We also provide an  
110 open-source Python package and pre-trained models of flies and mice to the community  
111 (see more in Code Availability).

112

## 113 **RESULTS**

### 114 **Workflow of Selfee and its downstream analyses**

115 Selfee is trained to generate Meta-representations at the frame level, which are then  
116 analyzed at different time scales. First, grayscale videos are decomposed into single  
117 frames, and three tandem frames are stacked into a live-frame to generate a motion-  
118 colored RGB picture (Figure 1A). These live-frames preserve not only spatial  
119 information (e.g., postures of each individual or relative distances and angles between  
120 individuals) within each channel but also temporal information across different  
121 channels. Live-frames are used to train Selfee to produce comprehensive and  
122 discriminative representations at the frame level (Figure 1B). These representations can  
123 be later used in numerous applications. For example, anomalous detection on mutant  
124 animals can discover new phenotypes compared with their genetic controls (Figure 1C).  
125 Also, the AR-HMM could be applied to model the micro-dynamics of behaviors, such  
126 as the duration of states or the probabilities of state transitions<sup>18</sup>. The AR-HMM splits  
127 videos into modules and yields behavioral state usages that visualize differences  
128 between genotypes (Figure 1D). In contrast, DTW could compare the long-term  
129 dynamics of animal behaviors and capture global differences at the video level<sup>31</sup> by  
130 aligning pairs of time series and calculating their similarities (Figure 1E). These three  
131 demonstrations cover different time scales from frame to video level, and other  
132 downstream analyses could also be incorporated into the workflow of Selfee.  
133 Compared with previous machine learning frameworks for animal behavior analysis,

134 Selfee has three major advantages. First, Selfee and the Meta-representations could be  
135 used for various tasks. The contrastive learning process of Selfee would allow output  
136 features to be appropriately compared by cosine similarity. Therefore, distance-based  
137 applications, including classification, clustering, and anomaly detection, would be  
138 easily realized. It was also reported that with some adjustment of backbones, self-  
139 supervised learning would facilitate tasks such as pose estimation<sup>32</sup> and object  
140 segmentation<sup>28,33</sup>. Those findings indicate that Selfee could be generalized, modified,  
141 and finetuned for animal pose estimation or segmentation tasks. Second, Selfee is a  
142 fully unsupervised method developed to annotate animal behaviors. Although some  
143 other techniques also adopt semi-supervised or unsupervised learning, they usually  
144 require manually labeled pre-defined key points of the images<sup>8,10</sup>; some methods also  
145 require expert-defined programs for better performance<sup>9</sup>. Key point selection and  
146 program incorporation require a significant amount of prior knowledge and are subject  
147 to human bias. In contrast, Selfee does not need any prior knowledge. Finally, Selfee is  
148 relatively hardware-inexpensive. Training Selfee only takes eight hours on a single  
149 RTX 3090, and the inference speed could reach 800 frames per second. Selfee could  
150 accept top-view 2D greyscale video frames as inputs so that neither depth cameras<sup>18</sup>  
151 nor fine-calibrated multi-view camera arrays<sup>10</sup> is required. Therefore, Selfee can be  
152 trained and used with routinely collected behavior videos on ordinary desktop  
153 workstations, warranting its accessibility by biology laboratories.

154



155 **Siamese convolutional neural networks capture discriminative representations of**  
156 **animal posture.**

157 Selfee contains a pair of Siamese CNNs trained to generate discriminative  
158 representations for live-frames. ResNet-50<sup>34</sup> is chosen as the backbone whose classifier  
159 layer is replaced by a three-layer multi-layer perceptron (MLP). These MLPs are called  
160 projectors which yield final representations during the inference stage. There are two  
161 branches in Selfee. The main branch is equipped with an additional predictor, while the  
162 reference branch is a copy of the main branch (the SimSiam style<sup>29</sup>). Both branches  
163 contain group discriminators after projectors and perform dimension reduction on  
164 extracted features for online clustering (Figure 2B).

165 During the training stage, batches of live-frames are randomly transformed twice and  
166 fed into the main branch and reference branch, respectively. Augmentations applied to  
167 live-frames include crop, rotation, flip, and application of the Turbo lookup table<sup>35</sup>  
168 followed by color jitters (Figure 2A, Figure 2—figure supplement 1). The reference  
169 branch yields a representation of received frames, while the main branch predicts the  
170 outcome of the reference branch. At the same time, they both produce clustering results  
171 of the current batch. The main branch is optimized for similar predictions and clustering  
172 results as the reference branch, and the reference branch will not receive gradient  
173 information to prevent mode collapse<sup>27,29</sup> (Figure 2C). In this way, Selfee is trained to  
174 be invariant to those transforms and focus on critical information to yield discriminative  
175 representations.

176 After the training stage, we evaluated the performance of Selfee with t-SNE  
177 visualization and k-NN classification. To investigate whether our model captures  
178 human-interpretable features, we manually labeled one clip of *Drosophila* courtship  
179 video and visualized those representations with t-SNE dimension reduction. On the t-  
180 SNE map, human-annotated courtship behaviors, including chasing, wing extension,  
181 copulation attempt, copulation, and non-interactive behaviors (“others”), separated  
182 from each other distinctively (Figure 2D).

183 Meta-representations can also be used for behavior classification. We manually labeled  
184 seven 10,000-frame videos (around five minutes each) as a pilot dataset. A weighed k-  
185 NN classifier was then constructed as previously reported<sup>24</sup>. Seven-fold cross-  
186 validation was performed on the dataset with the k-NN classifier, which achieved a  
187 mean F<sub>1</sub> score of 72.4% and achieved a similar classification result as human  
188 annotations (Figure 2E, F). The classifier had the worst recall score on wing extension  
189 behaviors (67% recall), likely because of the ambiguous intermediate states between  
190 chasing and wing extension (Figure 2—figure supplement 2A). The precisions also  
191 showed that this k-NN classifier tended to have strict criteria with wing extension and  
192 copulation and relatively loose criteria with chasing and copulation attempts (Figure  
193 2—figure supplement 2B). It was reported that independent human experts could only  
194 reach agreements on around 70% of wing extension frames<sup>36</sup>, comparable to the  
195 performance of our k-NN classifier.

196 We then asked whether Selfee can be generalized to analyze behaviors of other species.

197 We finetuned fly video pre-trained Selfee with mice mating behavior data. The mating  
198 behavior of mice can be defined mainly into five categories<sup>37</sup>, including social interest,  
199 mounting, intromission, ejaculation, and others (see Methods for detailed definitions).  
200 With t-SNE visualization, we found that five types of behaviors could be separated by  
201 Selfee, although mounting behaviors were rare and not concentrated (Figure 2G). We  
202 then used eight human-annotated videos to test the k-NN classification performance of  
203 Selfee-generated features. We achieved an F<sub>1</sub> score of 59.0% (Figure 2—figure  
204 supplement 3). Mounting, intromission, and ejaculation share similar static  
205 characteristics but are different in temporal dynamics. Therefore, we asked if more  
206 temporal information would assist the classification. Using the LightGBM classifier,  
207 we achieved a much higher classification performance by incorporating slide moving  
208 average and standard division of 81-frame time windows, the main frequencies, and  
209 their energy within 81-frame time windows. The average F<sub>1</sub> score of eight-fold cross-  
210 validation could reach 67.4%, and the classification results of the ensembled classifier  
211 (see Methods) were closed to human observations (Figure 2H, I). Nevertheless, it was  
212 still difficult to distinguish between mounting, intromission, and ejaculation because  
213 mounting and ejaculation are much rarer than social body contact or intromission.  
214 Selfee is more robust than the vanilla SimSiam networks when applied to the behavioral  
215 data. Behavioral data often suffer from catastrophic imbalance. For example, copulation  
216 attempts are around six-fold rarer than wing extension during fly courtship (Figure 2—  
217 figure supplement 5A). Therefore, we added group discriminators to vanilla SimSiam

218 networks which were reported to fight against the long-tail effect proficiently<sup>38</sup>. Aside  
219 from overcoming the long-tail effect, we also found group discriminators helpful for  
220 preventing mode collapse during ablation studies (Figure 2—figure supplement 5B, C,  
221 D, and Supplementary Table 1). Additionally, the convergence can be easily reached on  
222 grayscale images of similar objects (two flies), by which CNNs may not be well trained  
223 to extract good representations. Applying the Turbo lookup table on grayscale frames  
224 brought more complexity and made color jitters more powerful on grayscale images.  
225 Selfee would capture more useful features with this Turbo augmentation (Figure 2—  
226 figure supplement 5E, F, and Figure 2—figure supplement 6).

227

228 **Anomaly detection at the frame level identifies rare behaviors at the sub-second**  
229 **time scale.**

230 The representations produced by Selfee could be directly used for anomaly detection  
231 without further post-processing. During the training step, Selfee learns to compare  
232 Meta-representations of frames with cosine distance which is also used for anomaly  
233 detection. When given two groups of videos, namely the query group and the reference  
234 group, the anomaly score of each live-frame in the query group is calculated by two  
235 steps (Figure 3A). First, distances between the query live-frame and all reference live-  
236 frames are measured, and the k-nearest distance is referred to as its inter-group score  
237 (IES). Without further specification, k equals 1 in all anomaly detections in this work.  
238 Some false positives occurred when only the IES was used as the anomaly score (Figure

239 3—figure supplement 1A). The reason could be that two flies in a chamber could be in  
240 mathematically infinite relative positions and form a vast event space. However, each  
241 group usually only contains several videos, and each video is only recorded for several  
242 minutes. For some rare postures, even though the probability of observing them is  
243 similar in both the query and reference group, they might only occur in the query group  
244 but not in the reference group. Therefore, an intra-group score (IAS) is introduced in  
245 the second step to eliminate these false-positive effects. We assume that those rare  
246 events should not be sampled frequently in the query groups either. Thus, the IAS is  
247 defined as the k-nearest distance of the query frame against all other frames within its  
248 group, except those within the time window of  $\pm 50$  frames (Figure 3—figure  
249 supplement 1B). The final anomaly score is defined as the IES minus the IAS.

250 To test whether our methods could detect anomalous behavior in real-world data, we  
251 performed genetic screenings within fifteen neurotransmitter-related mutant alleles or  
252 neuron-silenced lines (with UAS-Kir2.1<sup>39</sup>) (Figure 3B). Their male-male interaction  
253 videos were inferred by Selfee trained on male-female courtship videos. Since we  
254 aimed to find interactions distinct from male-male courtship behaviors, a baseline of  
255 *ppk23>Kir2.1* flies was established because this line exhibit strong male-male courtship  
256 behaviors<sup>40</sup>. We compared the top-100 anomaly scores from sets of videos from  
257 experimental groups and wild-type control flies. The results revealed that one line,  
258 *CCHa2-R-RB>Kir2.1*, showed a significantly high anomaly score. By manually going  
259 through all anomalous live-frames, we further identified its phenotype as a brief tussle

260 behavior mixed with copulation attempts (Figure 3C, Video 1, 0.2x play speed). This  
261 behavior was ultra-fast and lasts for less than a quarter second (Figure 3D), making it  
262 difficult to be detected by human observers. Up to this point, we have demonstrated  
263 that the frame-level anomaly detection could capture sub-second behavior episode that  
264 human observers tend to neglect.

265 Selfee also revealed that *Trh* knock-out flies had an unusual close body contact during  
266 the screening. *Trh* is the crucial enzyme for serotonin biosynthesis, and its mutant flies  
267 showed a statistically significantly higher anomaly score (Figure 3B) than the wild-type  
268 control. Selfee identified 60 frames of abnormal behaviors within 42,000 input frames,  
269 occupying less than 0.15% of the total recording time. By manually going through all  
270 these frames, we concluded most of them as short-range body interactions (Figure 3E  
271 and Video 2, 0.2x play speed), and these social interactions could last for around half  
272 to one second on average (Figure 3F). Despite that serotonin signals were well-studied  
273 for controlling aggression behavior in flies<sup>41</sup>, to the best of our knowledge, the close  
274 body contact of flies and serotonergic neurons' role in this behavior has not been  
275 reported yet. A possible reason is that this behavior has no unique posture compared  
276 with other behaviors, like wing extension, and this behavior is too scarce to be noticed  
277 by human experts.

278 To further ask whether these close body contacts have biological significance, we  
279 performed corresponded behavior assays on mutant flies. Based on the fact that the *Trh*  
280 mutant male flies have a higher tolerance to body touch, we hypothesized that they

281 would have a decreased defensive behavior. As previously reported, fruit flies show  
282 robust defensive behavior to mechanical stimuli on their wings<sup>42,43</sup>. Decapitated flies  
283 would kick with their hind legs when a thin probe stimulates their wings. This  
284 stimulation mimics the invasion of parasitic mites and could be used to test its defensive  
285 behavior. Our results showed that *Trh* knock-out flies had a significantly lower kicking  
286 rate than control flies (Figure 3G), indicating a reduction of self-defensive intensity.  
287 Next, we performed social behavior assay<sup>44,45</sup> on the mutant flies because the close  
288 body contact can also be explained by reduced social repulsion. We measured the  
289 nearest distance, median distance, and average distance of each male flies in a forty-  
290 individual group placed in a vertical triangular chamber. By comparing median values  
291 of these distances of each replication, *Trh* knock-out flies kept significantly shorter  
292 distances from others than the control group (Figure 3H, I). The probability density  
293 function of their median distances also showed that knock-out flies had a closer social  
294 distance than control flies (Figure 3J). Therefore, we concluded that *Trh* knock-out flies  
295 had reduced social repulsion. Taken together, Selfee is capable of discovering novel  
296 features of animal behaviors with biological relevance when a proper baseline is  
297 defined.

298

### 299 **Modeling motion structure of *Drosophila* courtship behaviors.**

300 Animal behaviors have long-term structures beyond single-frame postures. The  
301 duration and proportions of each bout and transition probabilities of different behaviors

302 have been proven to have biological significance<sup>18,46</sup>. To better understand those long-  
303 term characteristics, we introduce AR-HMM and DTW analyses to model the temporal  
304 structure of *Drosophila* courtship behavior. AR-HMM is a powerful method to analyze  
305 stereotyped behavioral data<sup>18,47,48</sup>. It discovers modules of behaviors and describes the  
306 modules with auto-regressive matrixes. The transition probability of each state is  
307 defined by the transition matrix of the HMM (Figure 4A). In this way, AR-HMM could  
308 capture local structures of animal behaviors as well as syntaxes.

309 We asked if we could detect the dynamic changes of courtship behaviors of male flies  
310 by disturbing their chemosensation. *Ir76b* is an extensively studied (co)receptor that is  
311 known to mediate female pheromones detection<sup>49-52</sup>. We used an AR-HMM model with  
312 ten modules (No.1 to 10) to analyze the courtship of *Ir76b* mutant flies and their control  
313 group and focused on state usages. PCA of state usages revealed an apparent difference  
314 between mutant flies and control flies (Figure 4B). Module No.6 showed a statistically  
315 significant difference among ten discovered modules (Figure 4C). By manually going  
316 through all the frames of module No.6, we found that it mainly contained non-  
317 interactive behaviors with minor contaminations of courtship behaviors (Video 3, 1x  
318 play speed). To validate this result, we compared it with human annotations. Although  
319 this module did not cover all non-interactive behaviors that human experts would label,  
320 they showed a similar trend between the experimental and control group (Figure 4D).

321 We also performed AR-HMM analysis with a much larger module number. The PCA  
322 result was also distinct, and the previous module No.6 was split into five smaller



323 modules (No.2, 15, 24, 32, 34) containing non-interactive behaviors (Figure 4—figure  
324 supplement 1, Video 4-8, 1x play speed). This tuning indicated that AR-HMM analysis  
325 is robust regardless of the number of modules, same as a previous report<sup>18</sup>. Our results  
326 indicated that *Ir76b* mutation might affect male flies' detection of female pheromones  
327 and consequentially the temporal structure of their courtship behaviors. These findings  
328 prove that Selfee with AR-HMM could discover the differences in proportions of  
329 behaviors, similar to what was achieved with classic manual analysis such as the  
330 courtship index.

331 The AR-HMM modeling does not necessarily capture the difference of long-term  
332 dynamics intuitively, such as the latency of certain behaviors. To solve this problem,  
333 we introduce DTW analysis. DTW is a well-known algorithm to align time series,  
334 which returns the best-matched path and the matching similarity (Figure 4E). The  
335 alignment can be simplified as follows. When given the same start state and end state,  
336 it optimally maps all indices from the query series to the reference series monotonically.  
337 Pairs of mapped indices form a path to visualize the dynamic difference. The points  
338 upper than the diagonal line indicate that the current time point in the query group is  
339 matched to a future time point in the reference group so that the query group has faster  
340 dynamics and vice versa. In our experiments, cosine similarities of Selfee extracted  
341 representations are used to calculate warping paths.

342 Previously, DTW was widely applied to numerical measures of animal behaviors,  
343 including trajectory<sup>53</sup>, audios<sup>54</sup>, and acceleration<sup>55</sup>. For the first time, we applied DTW

344 to image data, with the aid of Selfee, to study the prolonged dynamic of animal  
345 behaviors. We analyzed whether the vision is essential for a male fly's copulation  
346 completion. Visual cues are essential for male flies to locate female flies during  
347 courtship<sup>56</sup>, and mutant flies of *NorpA*, which have defective visual transduction<sup>57</sup>, have  
348 a prolonged courtship latency in our experiments (Figure 4F), similar to previously  
349 findings<sup>58</sup>. When wild-type flies were used as the reference for the DTW, the group of  
350 *NorpA* mutant flies yielded a curve lower than the diagonal line, indicating a delay of  
351 their courtship behaviors (Figure 4G). In this way, our experiments confirm that Selfee  
352 and DTW could capture differences in long-term dynamics such as behavior latency. In  
353 conclusion, DTW and AR-HMM could capture temporal differences between control  
354 and experimental groups beyond single-frame postures, making Selfee a competent  
355 unsupervised method for traditional analyses like courtship index or copulation latency.

356

## 357 **DISCUSSION**

358 Here we use cutting-edge self-supervised learning methods and convolutional neural  
359 networks to extract Meta-representations from animal behavior videos. Siamese CNNs  
360 have proven their capability to learn comprehensive representations<sup>29</sup>. The cosine  
361 similarity, part of its loss function used for training, is rational and well-suited to  
362 measure similarities between the raw images. Besides, convolutional neural networks  
363 are trained end-to-end so that preprocessing steps like segmentation or key points  
364 extraction is unnecessary. By incorporating Selfee with different post-processing

365 methods, we can identify phenotypes of animal behaviors at different time scales. In  
366 the current work, we demonstrate that the extracted representations could be used not  
367 only for straightforward distance-based analyses such as t-SNE visualization or k-NN  
368 anomaly detection but also for sophisticated post-processing methods like AR-HMM.  
369 These validations confirm that the extracted Meta-representations are meaningful and  
370 valuable.

371 By applying our method to mice mating behavior, we show that our Selfee out-  
372 performed some of the widely used pose-estimation methods in multi-animal behavior  
373 analysis. The famous DeepLabCut and similar methods could identify human-defined  
374 key points on animals. However, when animals of the same color are recorded at a  
375 compromised resolution and their body contacts are intensive, the current version of  
376 DeepLabCut could hardly extract useful features (Figure 1—figure supplement 1,  
377 Video 9). The reason is that it is extremely difficult to unambiguously label body parts  
378 like nose, ears and hips when two mice are close enough, a task challenging even for  
379 human experts. By contrast, Selfee could readily identify the frame as “intromission”  
380 (Figure 1—figure supplement 1) as human experts would do. These results show that  
381 our methods could capture global characteristics of behaviors like human experts,  
382 making it well-suited for processing multi-animal behavior videos, compared with  
383 pose-estimation methods.

384 We also demonstrate that the cutting-edge self-supervised learning model is accessible  
385 to biology labs. Our model can be trained on only one RTX 3090 GPU with a batch size

386 of 256 within only 8 hours with the help of the newly proposed CLD loss function<sup>38</sup>  
387 and other improvements (see Methods for further details). Furthermore, when the model  
388 pre-trained with mice videos was applied to rat behaviors, we were able to achieve a  
389 zero-shot classification of five major types of social behaviors (Figure 2—figure  
390 supplement 4). Although the F<sub>1</sub> score was only 49.6%, it still captured the major  
391 differences between similar behaviors, such as allogrooming and social nose contact.  
392 Thus, we have demonstrated that self-supervised learning could be easily achieved with  
393 limited computation resources and a much shorter time and could be transferred to  
394 datasets that share similar visual characteristics.

395 Despite those advantages, there are some limitations of Selfee. First, because each live-  
396 frame only contains three raw frames, our model could not capture much information  
397 on the animal motion. It becomes more evident when Selfee is applied to highly  
398 dynamic behaviors such as mice mating behaviors. This can be overcome by increasing  
399 the computation because commonly used 3D convolution<sup>59</sup> or spatial-temporal  
400 attention<sup>60</sup> is good at dynamic information extraction but requires much more  
401 computational resources. Second, as previously reported, CNNs are highly vulnerable  
402 to image texture<sup>61</sup>. We observed that certain types of beddings of the behavior chamber  
403 could profoundly affect the performance of our neural networks (Figure 1—figure  
404 supplement 2), so in some cases, background removal is necessary (see Methods for  
405 further details). Lastly, Selfee could only use discriminative features within each batch,  
406 without any negative samples provided, so minor irrelevant differences could be

407 amplified and cause inconsistent results (named mode-split). This mode split may  
408 increase variations of downstream analyses.

409 We can envision at least two possible future directions for Selfee. One is to optimize  
410 the backbone of neural networks to extract better features. Advanced self-supervised  
411 learning methods like DINO<sup>33</sup> (with visual transformers, ViTs) could separate objects  
412 from the background and extract more explainable representations. Besides, by using  
413 ViTs, the neural network could be more robust against distractive textures<sup>62</sup>. At the same  
414 time, more temporal information can also be incorporated for a better understanding of  
415 motions. Combining these two, equipping ViTs with spatial-temporal attention could  
416 capture more temporal information.

417 Another direction will be explainable behavior forecasting for a deeper understanding  
418 of animal behaviors. For a long time, behavior forecasting has been a field with  
419 extensive investigations in which RNNs, LSTMs, or transformers are usually applied  
420 <sup>9,60,63</sup>. However, most of these works use coordinates of key points as inputs. Therefore,  
421 the trained model might predominantly focus on spatial movement information and  
422 discover fewer behavioral syntaxes. By representation learning, spatial information is  
423 essentially condensed so that more syntaxes might be highlighted. Transformer models  
424 for forecasting could capture correlations between sub-series as well as long-term  
425 trends like seasonality<sup>64</sup>. These deep learning methods would provide behavioral  
426 neuroscientists powerful tools to identify behavior motifs and syntaxes that organize  
427 stereotyped motifs beyond the Markov property.

428 **ACKNOWLEDGEMENTS**

429 We thank members of the Zhang lab for discussions. This work was supported by grants  
430 31871059 and 32022029 from the National Natural Science Foundation of China, grant  
431 Z181100001518001 from the Beijing Municipal Science & Technology Commission,  
432 and a ‘Brain+X’ seed grant from the IDG/McGovern Institute for Brain Research at  
433 Tsinghua to W.Z.. W.Z. is supported by Chinese Institute for Brain Research, Beijing.  
434 W.Z. is an awardee of the Young Thousand Talent Program of China.

435 **AUTHOR CONTRIBUTIONS**

436 Y.J. and J.H. coded the Selfee neural network and other accessory parts. Y.J., X.G. and  
437 S.L. performed animal experiments and analyzed data. W.Z. and X.X. supervised the  
438 project. Y.J. and W.Z. wrote the manuscript. All authors discussed and commented on  
439 the manuscript.

440 **DECLARATION OF INTERESTS**

441 The authors declare no competing interests.

442 **DATA AND CODE AVAILABILITY STATEMENT**

443 Major data used in this study were uploaded to Dryad. Data could be accessed via:

444 <https://datadryad.org/stash/share/BnIoOnaweOn2fc-sllSO0FhJJqduXQYaNu->

445 [KuPgZ394](https://datadryad.org/stash/share/BnIoOnaweOn2fc-sllSO0FhJJqduXQYaNu-KuPgZ394) or its DOI 10.5061/dryad.brV15dvb8. We also shared our pretrained weights

446 with Google Drive:

447 <https://drive.google.com/file/d/1A3U5guNEKA3Bi9H3QnfstZDEZ6aesqcR/view?usp>

448 [=sharing](#). With the uploaded dataset and pretrained weights, our experiments could be

449 replicated. However, due to its huge size and the limited internet service resources, we

450 are currently not able to share our full training dataset. The full dataset is as large as

451 400GB, which is hard to upload to a public server and will be difficult for others users

452 to download.

453 For training dataset, it would be available from the corresponding author upon

454 reasonable request ([wei\\_zhang@mail.tsinghua.edu.cn](mailto:wei_zhang@mail.tsinghua.edu.cn)), and then we can discuss how

455 to transfer such a big dataset. No project proposal is needed as long as the dataset is not

456 used for any commercial purpose.

457 Our Python scripts could be accessed on GitHub: <https://github.com/EBGU/Selfee>

458 Other software used in our project include ImageJ(<https://imagej.net/software/fiji/>) and

459 GraphPad Prism(<https://www.graphpad.com/>).

460 All data used to plot graphs and charts in the manuscript can be fully accessed on Dryad

461 (DOI 10.5061/dryad.brV15dvb8).

462 **REFERENCES**

- 463 1 Hall, J. The mating of a fly. *Science* **264**, 1702-1714,  
464 doi:10.1126/science.8209251 (1994).
- 465 2 Kabra, M., Robie, A. A., Rivera-Alba, M., Branson, S. & Branson, K. JAABA:  
466 interactive machine learning for automatic annotation of animal behavior.  
467 *Nature Methods* **10**, 64-67, doi:10.1038/nmeth.2281 (2013).
- 468 3 Jiang, Z., Chazot, P. L., Celebi, M. E., Crookes, D. & Jiang, R. Social Behavioral  
469 Phenotyping of *Drosophila* With a 2D–3D Hybrid CNN Framework. *IEEE*  
470 *Access* **7**, 67972-67982, doi:10.1109/ACCESS.2019.2917000 (2019).
- 471 4 Segalin, C. *et al.* The Mouse Action Recognition System (MARS): a software  
472 pipeline for automated analysis of social behaviors in mice. *bioRxiv*,  
473 2020.2007.2026.222299, doi:10.1101/2020.07.26.222299 (2020).
- 474 5 Mathis, A. *et al.* DeepLabCut: markerless pose estimation of user-defined body  
475 parts with deep learning. *Nature Neuroscience* **21**, 1281-1289,  
476 doi:10.1038/s41593-018-0209-y (2018).
- 477 6 Günel, S. *et al.* DeepFly3D, a deep learning-based approach for 3D limb and  
478 appendage tracking in tethered, adult *Drosophila*. *eLife* **8**, e48571,  
479 doi:10.7554/eLife.48571 (2019).
- 480 7 Graving, J. M. *et al.* DeepPoseKit, a software toolkit for fast and robust animal  
481 pose estimation using deep learning. *eLife* **8**, e47994, doi:10.7554/eLife.47994  
482 (2019).
- 483 8 Luxem, K., Fuhrmann, F., Kürsch, J., Remy, S. & Bauer, P. Identifying  
484 Behavioral Structure from Deep Variational Embeddings of Animal Motion.  
485 *bioRxiv*, 2020.2005.2014.095430, doi:10.1101/2020.05.14.095430 (2020).
- 486 9 Sun, J. J. *et al.* in *2021 IEEE/CVF Conference on Computer Vision and Pattern*  
487 *Recognition (CVPR)*. 2876-2885 (2021).
- 488 10 Huang, K. *et al.* A hierarchical 3D-motion learning framework for animal  
489 spontaneous behavior mapping. *Nature Communications* **12**, 2784,  
490 doi:10.1038/s41467-021-22970-y (2021).
- 491 11 Calhoun, A. J., Pillow, J. W. & Murthy, M. Unsupervised identification of the  
492 internal states that shape natural behavior. *Nature Neuroscience* **22**, 2040-2049,  
493 doi:10.1038/s41593-019-0533-x (2019).
- 494 12 Zhou, Y., Cao, L.-H., Sui, X.-W., Guo, X.-Q. & Luo, D.-G. Mechanosensory  
495 circuits coordinate two opposing motor actions in *Drosophila* feeding. *Science*  
496 *Advances* **5**, eaaw5141, doi:10.1126/sciadv.aaw5141 (2019).
- 497 13 Mezzera, C. *et al.* Ovipositor Extrusion Promotes the Transition from Courtship  
498 to Copulation and Signals Female Acceptance in *Drosophila melanogaster*.  
499 *Current Biology* **30**, 3736-3748.e3735,  
500 doi:<https://doi.org/10.1016/j.cub.2020.06.071> (2020).
- 501 14 Zhang, L. *et al.* Parallel Mechanosensory Pathways Direct Oviposition  
502 Decision-Making in *Drosophila*. *Current Biology* **30**, 3075-3088.e3074,



- 503 doi:<https://doi.org/10.1016/j.cub.2020.05.076> (2020).
- 504 15 Bayless, D. W. *et al.* Limbic Neurons Shape Sex Recognition and Social  
505 Behavior in Sexually Naive Males. *Cell* **176**, 1190-1205.e1120,  
506 doi:10.1016/j.cell.2018.12.041 (2019).
- 507 16 Zhang, S. X. *et al.* Hypothalamic dopamine neurons motivate mating through  
508 persistent cAMP signalling. *Nature* **597**, 245-249, doi:10.1038/s41586-021-  
509 03845-0 (2021).
- 510 17 Wei, Y.-C. *et al.* Medial preoptic area in mice is capable of mediating sexually  
511 dimorphic behaviors regardless of gender. *Nature Communications* **9**, 279,  
512 doi:10.1038/s41467-017-02648-0 (2018).
- 513 18 Wiltschko, Alexander B. *et al.* Mapping Sub-Second Structure in Mouse  
514 Behavior. *Neuron* **88**, 1121-1135,  
515 doi:<https://doi.org/10.1016/j.neuron.2015.11.031> (2015).
- 516 19 Berman, G. J., Choi, D. M., Bialek, W. & Shaevitz, J. Mapping the stereotyped  
517 behaviour of freely moving fruit flies. *Journal of the Royal Society Interface* **11**  
518 (2014).
- 519 20 Lecun, Y., Bottou, L., Bengio, Y. & Haffner, P. Gradient-based learning applied  
520 to document recognition. *Proceedings of the IEEE* **86**, 2278-2324,  
521 doi:10.1109/5.726791 (1998).
- 522 21 Romero, D., Ruedin, A. M. C. & Seijas, L. in *ICIAP*.
- 523 22 Aradhya, V. N. M., Kumar, G. H. & Noushath, S. in *2007 International*  
524 *Conference on Signal Processing, Communications and Networking*. 626-629.
- 525 23 Ji, X., Vedaldi, A. & Henriques, J. F. Invariant Information Clustering for  
526 Unsupervised Image Classification and Segmentation. *2019 IEEE/CVF*  
527 *International Conference on Computer Vision (ICCV)*, 9864-9873 (2019).
- 528 24 Wu, Z., Xiong, Y., Yu, S. X. & Lin, D. Unsupervised Feature Learning via Non-  
529 Parametric Instance-level Discrimination. *ArXiv abs/1805.01978* (2018).
- 530 25 Caron, M. *et al.* Unsupervised Learning of Visual Features by Contrasting  
531 Cluster Assignments. *ArXiv abs/2006.09882* (2020).
- 532 26 Chen, T., Kornblith, S., Norouzi, M. & Hinton, G. E. A Simple Framework for  
533 Contrastive Learning of Visual Representations. *ArXiv abs/2002.05709* (2020).
- 534 27 Grill, J.-B. *et al.* Bootstrap Your Own Latent: A New Approach to Self-  
535 Supervised Learning. *ArXiv abs/2006.07733* (2020).
- 536 28 He, K., Fan, H., Wu, Y., Xie, S. & Girshick, R. B. Momentum Contrast for  
537 Unsupervised Visual Representation Learning. *2020 IEEE/CVF Conference on*  
538 *Computer Vision and Pattern Recognition (CVPR)*, 9726-9735 (2020).
- 539 29 Chen, X. & He, K. Exploring Simple Siamese Representation Learning. *ArXiv*  
540 *abs/2011.10566* (2020).
- 541 30 Zbontar, J., Jing, L., Misra, I., LeCun, Y. & Deny, S. in *ICML*.
- 542 31 Myers, C., Rabiner, L. & Rosenberg, A. Performance tradeoffs in dynamic time  
543 warping algorithms for isolated word recognition. *IEEE Transactions on*  
544 *Acoustics, Speech, and Signal Processing* **28**, 623-635,

- 545 doi:10.1109/TASSP.1980.1163491 (1980).
- 546 32 Dahiya, A., Spurr, A. & Hilliges, O. in *NeurIPS 2020 Workshop on Pre-*  
547 *registration in Machine Learning* Vol. 148 (eds Bertinetto Luca *et al.*) 255--  
548 271 (PMLR, Proceedings of Machine Learning Research, 2021).
- 549 33 Caron, M. *et al.* Emerging Properties in Self-Supervised Vision Transformers.  
550 *ArXiv abs/2104.14294* (2021).
- 551 34 He, K., Zhang, X., Ren, S. & Sun, J. Deep Residual Learning for Image  
552 Recognition. *2016 IEEE Conference on Computer Vision and Pattern*  
553 *Recognition (CVPR)*, 770-778 (2016).
- 554 35 Mikhailov, A. *Turbo, An Improved Rainbow Colormap for Visualization*,  
555 [https://ai.googleblog.com/2019/08/turbo-improved-rainbow-colormap-](https://ai.googleblog.com/2019/08/turbo-improved-rainbow-colormap-for.html)  
556 [for.html](https://ai.googleblog.com/2019/08/turbo-improved-rainbow-colormap-for.html)> (2019).
- 557 36 Leng, X., Wohl, M., Ishii, K., Nayak, P. & Asahina, K. Quantifying influence of  
558 human choice on the automated detection of *Drosophila* behavior by a  
559 supervised machine learning algorithm. *PLoS One* **15**, e0241696,  
560 doi:10.1371/journal.pone.0241696 (2020).
- 561 37 McGill, T. E. Sexual Behavior in Three Inbred Strains of Mice. *Behaviour* **19**,  
562 341-350 (1962).
- 563 38 Wang, X., Liu, Z. & Yu, S. X. Unsupervised Feature Learning by Cross-Level  
564 Instance-Group Discrimination. *arXiv: Computer Vision and Pattern*  
565 *Recognition* (2020).
- 566 39 Paradis, S., Sweeney, S. T. & Davis, G. W. Homeostatic Control of Presynaptic  
567 Release Is Triggered by Postsynaptic Membrane Depolarization. *Neuron* **30**,  
568 737-749, doi:[https://doi.org/10.1016/S0896-6273\(01\)00326-9](https://doi.org/10.1016/S0896-6273(01)00326-9) (2001).
- 569 40 Thistle, R., Cameron, P., Ghorayshi, A., Dennison, L. & Scott, K. Contact  
570 Chemoreceptors Mediate Male-Male Repulsion and Male-Female Attraction  
571 during *Drosophila* Courtship. *Cell* **149**, 1140-1151,  
572 doi:<https://doi.org/10.1016/j.cell.2012.03.045> (2012).
- 573 41 Alekseyenko, Olga V. *et al.* Single Serotonergic Neurons that Modulate  
574 Aggression in *Drosophila*. *Current Biology* **24**, 2700-2707,  
575 doi:<https://doi.org/10.1016/j.cub.2014.09.051> (2014).
- 576 42 Li, J. *et al.* A Defensive Kicking Behavior in Response to Mechanical Stimuli  
577 Mediated by *Drosophila* Wing Margin Bristles. *The Journal of Neuroscience* **36**,  
578 11275-11282, doi:10.1523/jneurosci.1416-16.2016 (2016).
- 579 43 Liu, C. *et al.* A neural circuit encoding mating states tunes defensive behavior  
580 in *Drosophila*. *Nature Communications* **11**, 3962, doi:10.1038/s41467-020-  
581 17771-8 (2020).
- 582 44 Simon, A. F. *et al.* A simple assay to study social behavior in *Drosophila*:  
583 measurement of social space within a group1. *Genes, Brain and Behavior* **11**,  
584 243-252, doi:<https://doi.org/10.1111/j.1601-183X.2011.00740.x> (2012).
- 585 45 Au - McNeil, A. R. *et al.* Conditions Affecting Social Space in *Drosophila*  
586 *melanogaster*. *JoVE*, e53242, doi:doi:10.3791/53242 (2015).

- 587 46 Mueller, J. M., Ravbar, P., Simpson, J. H. & Carlson, J. M. *Drosophila*  
588 *melanogaster* grooming possesses syntax with distinct rules at different  
589 temporal scales. *PLOS Computational Biology* **15**, e1007105,  
590 doi:10.1371/journal.pcbi.1007105 (2019).
- 591 47 Wiltchko, A. B. *et al.* Revealing the structure of pharmacobehavioral space  
592 through motion sequencing. *Nature Neuroscience* **23**, 1433-1443,  
593 doi:10.1038/s41593-020-00706-3 (2020).
- 594 48 Rudolph, S. *et al.* Cerebellum-Specific Deletion of the GABAA Receptor  $\delta$   
595 Subunit Leads to Sex-Specific Disruption of Behavior. *Cell Reports* **33**, 108338,  
596 doi:<https://doi.org/10.1016/j.celrep.2020.108338> (2020).
- 597 49 Luo, R. *et al.* Molecular basis and homeostatic regulation of Zinc taste.  
598 *Protein&Cell*, doi:10.1007/s13238-021-00845-8 (2021).
- 599 50 Sánchez-Alcañiz, J. A. *et al.* An expression atlas of variant ionotropic glutamate  
600 receptors identifies a molecular basis of carbonation sensing. *Nature*  
601 *Communications* **9**, 4252, doi:10.1038/s41467-018-06453-1 (2018).
- 602 51 Ganguly, A. *et al.* A Molecular and Cellular Context-Dependent Role for Ir76b  
603 in Detection of Amino Acid Taste. *Cell Reports* **18**, 737-750,  
604 doi:<https://doi.org/10.1016/j.celrep.2016.12.071> (2017).
- 605 52 Zhang, Y. V., Ni, J. & Montell, C. The Molecular Basis for Attractive Salt-Taste  
606 Coding in *Drosophila*. *Science* **340**, 1334-1338,  
607 doi:10.1126/science.1234133 (2013).
- 608 53 Cleasby, I. R. *et al.* Using time-series similarity measures to compare animal  
609 movement trajectories in ecology. *Behavioral Ecology and Sociobiology* **73**,  
610 151, doi:10.1007/s00265-019-2761-1 (2019).
- 611 54 Kohlsdorf, D., Herzing, D. & Starner, T. Methods for discovering models of  
612 behavior: A case study with wild Atlantic spotted dolphins. *Animal Behavior*  
613 *and Cognition* **3**, 265-287 (2016).
- 614 55 Aurasopon, A. Dynamic Time Warping for classifying cattle behaviors and  
615 reducing acceleration data size. *Agricultural Engineering International: The*  
616 *CIGR Journal* **18**, 293-300 (2016).
- 617 56 Ribeiro, I. M. A. *et al.* Visual Projection Neurons Mediating Directed Courtship  
618 in *Drosophila*. *Cell* **174**, 607-621.e618, doi:10.1016/j.cell.2018.06.020 (2018).
- 619 57 Bloomquist, B. T. *et al.* Isolation of a putative phospholipase c gene of  
620 *drosophila*, *norpA*, and its role in phototransduction. *Cell* **54**, 723-733,  
621 doi:[https://doi.org/10.1016/S0092-8674\(88\)80017-5](https://doi.org/10.1016/S0092-8674(88)80017-5) (1988).
- 622 58 Markow, T. A. & Manning, M. Mating success of photoreceptor mutants of  
623 *Drosophila melanogaster*. *Behavioral and Neural Biology* **29**, 276-280,  
624 doi:[https://doi.org/10.1016/S0163-1047\(80\)90612-3](https://doi.org/10.1016/S0163-1047(80)90612-3) (1980).
- 625 59 Ji, S., Xu, W., Yang, M. & Yu, K. 3D Convolutional Neural Networks for Human  
626 Action Recognition. *IEEE Transactions on Pattern Analysis and Machine*  
627 *Intelligence* **35**, 221-231, doi:10.1109/TPAMI.2012.59 (2013).
- 628 60 Aksan, E., Cao, P., Kaufmann, M. & Hilliges, O. Attention, please: A Spatio-

- 629 temporal Transformer for 3D Human Motion Prediction. *ArXiv* **abs/2004.08692**  
630 (2020).
- 631 61 Geirhos, R. *et al.* ImageNet-trained CNNs are biased towards texture;  
632 increasing shape bias improves accuracy and robustness. *ArXiv* **abs/1811.12231**  
633 (2019).
- 634 62 Naseer, M. *et al.* Intriguing Properties of Vision Transformers. *ArXiv*  
635 **abs/2105.10497** (2021).
- 636 63 Fragkiadaki, K., Levine, S. & Malik, J. Recurrent Network Models for  
637 Kinematic Tracking. *ArXiv* **abs/1508.00271** (2015).
- 638 64 Wu, H., Xu, J., Wang, J. & Long, M. Autoformer: Decomposition Transformers  
639 with Auto-Correlation for Long-Term Series Forecasting. *ArXiv*  
640 **abs/2106.13008** (2021).

641

## 642 **FIGURE LEGENDS**

### 643 **Figure 1 | The framework of Selfee and its downstream applications.**

644 (A) One live-frame is composed of three tandem frames in R, G, and B channels,  
645 respectively.

646 (B) Live-frames are used to train Selfee, which adopts a backbone of ResNet-50.

647 (C, D, and E) Representations produced by Selfee could be used for anomaly detection  
648 that could identify unusual animal postures in the query video compared with the  
649 reference videos (C) AR-HMM (autoregressive hidden Markov model) that models the  
650 local temporal characteristics of behaviors and clusters frames into modules (states) and  
651 calculates stages usages of different genotypes (D) DTW (dynamic time warping) that  
652 aligns behavior videos to reveal differences of long-term dynamics (E) and other  
653 potential tasks including behavior classification, forecasting, or even image  
654 segmentation and pose estimation after properly modifying and finetuning of the neural  
655 networks.

656

### 657 **Figure 2 | The network structure of Selfee and validation of Selfee with human** 658 **annotations.**

659 (A) The architecture of Selfee networks. Each live-frame is randomly transformed  
660 twice before being fed into Selfee.

661 (B) Selfee adopts a SimSiam-style network structure with additional group  
662 discriminators. Loss 1 is canonical negative cosine loss, and loss 2 is the newly  
663 proposed CLD loss.

664 (C) The reference branch of Selfee is not involved in backward propagation.

665 (D) Visualization of fly courtship live-frames with t-SNE dimension reduction. Each  
666 dot was colored based on human annotations.

667 (E) The confusion matrix of the k-NN classifier for fly courtship behavior, normalized  
668 by the numbers of each behavior in the ground truth. The average F<sub>1</sub> score of the nine-  
669 fold cross-validation was 72.4%, and mAP was 75.8%.

670 (F) A visualized comparison of labels produced by the k-NN classifier and human  
671 annotations of fly courtship behaviors.

672 (G) Visualization of live-frames of mice mating behaviors with t-SNE dimension  
673 reduction. Each dot is colored based on human annotations.

674 (H) The confusion matrix of the LightGBM classifier for mice mating behaviors,  
675 normalized by the numbers of each behavior in the ground truth. For the LightGBM  
676 classifier, the average F<sub>1</sub> score of the nine-fold cross-validation was 67.4%, and mAP  
677 was 69.1%.

678 (I) A visualized comparison of labels produced by the LightGBM classifier and human  
679 annotations of mice mating behaviors.  
680

681 **Figure 3 | Anomalous posture detection using Selfee-produced features.**

682 (A) The calculation process of anomaly scores. Each query frame is compared with  
683 every reference frame, and the nearest distance was named IES (the thickness of lines  
684 indicates distances). Each query frame is also compared with every query frame, and  
685 the nearest distance is called IAS. The final anomaly score of each frame equals IES  
686 minus IAS.

687 (B) Genetic screening of fifteen fly lines with mutations in neurotransmitter genes or  
688 with specific neurons silenced ( $n \geq 6$  for each genotype). RA is short for CCHa2-R-RA,  
689 and RB is short for CCHa2-R-RB.  $CCHa2-R-RB^{Gal4} > Kir2.1$ ,  $q < 0.0001$ ;  $Trh^{Gal4}$ ,  $q =$   
690  $0.0432$ ; one-way ANOVA with Benjamini and Hochberg correction.

691 (C) Examples of mixed tussles and copulation attempts identified in  $CCHa2-R-RB^{Gal4} >$   
692  $Kir2.1$  flies.

693 (D) The temporal dynamic of anomaly scores during the mixed behavior, centralized at  
694 1.67 s. SEM is indicated with the light color region.

695 (E) Examples of close body contact behaviors identified in  $Trh^{Gal4}$  flies.

696 (F) The cosine similarity between the center frame of the close body contact behaviors  
697 (1.67s) and their local frames. SEM is indicated with the light color region.

698 (G) The kicking index of  $Trh^{Gal4}$  flies ( $n=30$ ) was significantly lower than  $w^{1118}$  flies  
699 ( $n=27$ ),  $p = 0.0034$ , Mann Whitney test.

700 (H) Examples of social aggregation behaviors of  $Trh^{Gal4}$  flies and  $w^{1118}$  flies.

701 (I) Social distances of  $Trh^{Gal4}$  flies ( $n=6$ ) and  $w^{1118}$  flies ( $n=6$ ).  $Trh^{Gal4}$  flies had much  
702 closer social distances with each other compared with  $w^{1118}$  flies; nearest,  $p = 0.0043$ ;  
703 median,  $p = 0.002$ ; average,  $p = 0.0087$ ; all Mann Whitney test.

704 (J) Distributions of the median social distance of  $Trh^{Gal4}$  flies and  $w^{1118}$  flies.  
705 Distributions were calculated within each replication. Average distributions are  
706 indicated with solid lines, and SEMs are indicated with light color regions.

707

708 **Figure 4 | Time-series analyses using Selfee produced features.**

709 (A) A brief illustration of the AR-HMM model. The local autoregressive property is  
710 determined by  $\beta_t$ , the autoregressive matrix, which is yield based on the current hidden  
711 state of the HMM. The transition between each hidden state is described by the  
712 transition matrix ( $p_{ij}$ ).

713 (B) PCA visualization of state usages of wild-type flies ( $n=7$ ) and  $Ir76b^l$  mutant flies  
714 ( $n=7$ ).

715 (C) State usages of ten modules. Module No. 6 showed significantly different usages in

716 wild-type and mutant flies;  $q = 0.022$ , Mann Whitney test with Benjamini and Hochberg  
717 correction.

718 (D) Module No.6 showed similar statistic results as manually labeled non-interactive  
719 behaviors. Module No.6,  $p = 0.0006$ ; manually labeled non-interactive behaviors,  $p =$   
720  $0.0006$ ; Both were Mann Whitney test.

721 (E) A brief illustration of the DTW model. The transformation from a rounded rectangle  
722 to an ellipse could contain six steps (grey reference shapes). The query transformation  
723 lags at step 2 but surpasses at step 4. The dynamic is visualized on the right panel.

724 (F) *NorpA*<sup>36</sup> flies ( $n=6$ ) showed a significantly longer copulation latency than wild-type  
725 flies ( $n=7$ ),  $p = 0.0495$ , Mann Whitney test.

726 (G) *NorpA*<sup>36</sup> flies had delayed courtship dynamics than wild-type flies with DTW  
727 visualization.

728

729 **Figure 1—figure supplement 1 | A comparison between Selfee and DeepLabCut on**  
730 **animals of the same color.**

731 (A) Selfee is more robust when applied to intensive social behaviors of animals of the  
732 same color, like mating behaviors of two black mice. The image is from the  
733 intromission behavior which could be identified by Selfee equipped with the trained  
734 LightGBM classifier. In contrast, a trained DeepLabCut model identified it as a single  
735 mouse.

736

737 **Figure 1—figure supplement 2 | Beddings and backgrounds that affect training**  
738 **and inference of Selfee.**

739 (A) Textures on the damped filter paper would mislead Selfee to output features similar  
740 with copulation but not wing extension (ground truth).

741 (B) Background inconsistency would affect the training process when Selfee was  
742 applied to mice behavior data. Therefore, backgrounds were removed from all frames  
743 to avoid potential defects.

744

745 **Figure 2—figure supplement 1 | Different augmentations used for Selfee training.**

746 (A) Visualization of each augmentation. Detailed descriptions of each augmentation  
747 could be found in Methods and source codes.

748

749 **Figure 2—figure supplement 2 | Difficulties on fly courtship behavior classification.**

750 (A) Some wing extension frames are hard to distinguish from chasing behaviors. Blue  
751 indicators point at slightly extended wings.

752 (B) The confusion matrix of the k-NN classifier, normalized by the numbers of each  
753 behavior in inferred labels.

754

755 **Figure 2—figure supplement 3 | Classification of mice mating behaviors with**

756 **Selfee extracted features.**

757 (A) For the k-NN classifier, the average  $F_1$  score of the nine-fold cross-validation was  
758 59.0%, and mAP was 53.0%. The confusion matrix of the k-NN classifier, normalized  
759 by the numbers of each behavior in the ground truth.

760 (B) The confusion matrix of the k-NN classifier, normalized by the numbers of each  
761 behavior in inferred labels.

762 (C) The confusion matrix of the LightGBM classifier, normalized by the numbers of  
763 each behavior in inferred labels. The LightGBM classifier had a much better  
764 performance compared with the k-NN classifier.

765

766 **Figure 2—figure supplement 4 | k-NN classification of rat behaviors with Selfee**

767 **trained on mice datasets.**

768 (A) The average  $F_1$  score of the nine-fold cross-validation was 49.6%, and mAP was  
769 46.6%. The confusion matrix of the k-NN classifier, normalized by the numbers of each  
770 behavior in the ground truth.

771 (B) The confusion matrix of the k-NN classifier, normalized by the numbers of each  
772 behavior in inferred labels.

773

774 **Figure 2—figure supplement 5 | Ablation test of Selfee training process on fly**

775 **datasets.**

776 (A) The distribution of different behaviors in wild-type flies courtship videos.

777 (B) Visualization of the same live-frames as Figure 2D with t-SNE dimension reduction.  
778 Used representations were extracted by models trained without CLD loss. Each dot is  
779 colored based on human annotations. The legend is shared with panel A.

780 (C) The confusion matrix of the k-NN classifier, normalized by the numbers of each



781 behavior in the ground truth. Used representations were extracted by models trained  
782 without CLD loss.

783 (D) Collapse levels during the training process. Without CLD loss, Selfee suffered from  
784 catastrophic mode collapse. Details for collapse level calculation could be found in  
785 Methods.

786 (E) Visualization of the same live-frames as Figure 2D with t-SNE dimension reduction.  
787 Used representations were extracted by models trained without Turbo transformation.  
788 Each dot is colored based on human annotations. The legend is shared with panel A.

789 (F) The confusion matrix of the k-NN classifier, normalized by the numbers of each  
790 behavior in the ground truth. Used representations were extracted by models trained  
791 without Turbo transformation.

792

793 **Figure 3—figure supplement 1 | Using intra-group score (IAS) to eliminate false-**  
794 **positive results in anomaly detections.**

795 (A) Anomaly scores without IAS of wild-type male-male interactions with the same  
796 genotype as references. The blue region indicates the max anomaly score when using  
797 IAS; blue dots indicate anomaly scores without IAS that fall into the blue region; red  
798 dots indicate false-positive anomaly scores.

799 (B) The cosine similarity between the center frame of wild-type courtship behaviors  
800 (1.67s) and their local frames. SEM is indicated with the light color region. Seven  
801 videos containing 70,000 frames were split into non-overlapping 100-frame fragments  
802 for calculations. Beyond  $\pm 50$  frames, the cosine similarity dropped to a much lower  
803 level, not affecting anomaly detection.

804

805 **Figure 4—figure supplement 1 | AR-HMM produced features with 40 modules**  
806 **using Selfee.**

807 (A) PCA visualization of state usages of wild-type flies and *Ir76b<sup>1</sup>* mutant flies.

808 (B) State usages of forty modules. Module No. 2, 15, 24, 32, 34 showed significantly  
809 different usages in wild-type and mutant flies;  $q = 0.029, 0.029, 0.049, 0.049, 0.029$   
810 respectively, Mann Whitney test with Benjamini and Hochberg correction.

811 (C) The collection of modules No. 2, 15, 24, 32, 34 showed similar statistic results as  
812 manually labeled non-interactive behaviors. AR-HMM module collection,  $p = 0.0006$ ;  
813 manually labeled non-interactive behaviors,  $p = 0.0006$ ; all Mann Whitney test.

814

## 815 MATERIALS AND METHODS

### 816 Fly stocks

817 All fly strains were maintained under a 12 h/12 h light/dark cycle at 25°C and 60%  
818 humidity (PERCIVAL incubator). The following fly lines were acquired from  
819 Bloomington *Drosophila* Stock Center: *CCHa1<sup>attP</sup>* (84458), *CCHa1-R<sup>attP</sup>* (84459),  
820 *CCHa2<sup>attP</sup>* (84460), *CCHa2-R<sup>attP</sup>* (84461), *CCHa2-R-RA<sup>Gal4</sup>* (84603), *CCHa2-R-*  
821 *RB<sup>Gal4</sup>* (84604), *CNMa<sup>attP</sup>* (84485), *Oamb<sup>attP</sup>* (84555), *Dop2R<sup>KO</sup>* (84720), *DopEcR<sup>Gal4</sup>*  
822 (84717), *SerT<sup>attP</sup>* (84572), *Trh<sup>Gal4</sup>* (86146), *TK<sup>attP</sup>* (84579), *Ir76b<sup>1</sup>* (51309), *NorpA<sup>36</sup>*  
823 (9048), UAS-Kir2.1 (6596). *Tdc2<sup>RO54</sup>* was a gift from Dr. Yufeng Pan at Southeast  
824 University, China. Taotie-Gal4 was a gift from Dr. Yan Zhu at Institute of Biophysics,  
825 Chinese Academy of Sciences, China.

826

### 827 Fly courtship behavior and male-male interaction

828 Virgin female flies were raised for 4~6 days in fifteen-fly groups, and naïve male flies  
829 were kept in isolated vials for 8~12 days. All behavioral experiments were done under  
830 25 °C and 45%~50% humidity. Flies were transferred into a customized chamber of  
831 3 mm height and 10 mm diameter by a homemade aspirator. Fly behaviors were  
832 recorded using a stereoscopic microscope-mounted with a CCD camera (Basler ORBIS  
833 OY-A622f-DC) at the resolution of 1000×500 (for two chambers at the same time), or  
834 640×480 (for individual chambers) and a frame rate of 30 Hz. Five types of behaviors  
835 were annotated manually, including “chasing” (a male fly follows a female fly), “wing  
836 extension” (a male fly extends unilateral wing and orientates to the female to sing  
837 courtship son, “copulation attempt” (a male fly bends its abdomen toward the genitalia  
838 of the female or the unstable state that male fly mounts on a female with its wings open),  
839 and “copulation” (male fly mounts on a female in a stable posture for several minutes).

840

### 841 Fly defensive behavior assay

842 The kicking behavior was tested based on previously reported paradigms<sup>1,2</sup>. Briefly,  
843 flies were raised in groups for 3~5 days. Flies were anesthetized on ice, and then male  
844 flies were decapitalized and transferred to 35 mm Petri dishes with damped filter paper  
845 on the bottom to keep the moisture. Flies were allowed to recover for around 30 minutes  
846 in the dishes. The probe for stimulation was homemade from a heat-melt yellow pipette  
847 tip, and the probe's tip was 0.3 mm. Each side of flies' wing margin was gently touched  
848 5 times, and the kicking behavior was recorded manually. The statistical analysis was

849 performed with the Mann Whitney test with GraphPad Prism Software.

### 850 **Social behavior assay for flies**

851 The social distance was tested based on the previously reported method<sup>3</sup>. Briefly, flies  
852 were raised in groups for 3 days. Flies were anesthetized paralyzed on ice, and male  
853 flies were picked and transferred to new vials (around 40 flies per vial). Flies were  
854 allowed to recover for one day. The vertical triangular chambers were cleaned with 75%  
855 ethanol and dried with paper towels. After assembly, flies were transferred into the  
856 chambers by a homemade aspirator. The photos were taken after 20 min, and the  
857 positions of each fly were manually marked in ImageJ. The social distances were  
858 measured with the lateral sides of the chambers (16.72 cm) as references, and the  
859 median values of the nearest, median, and average distance of each replication are  
860 calculated. The statistical analysis was performed with the Mann Whitney test in  
861 GraphPad Prism Software.

862

### 863 **Mice mating behavior assay**

864 Wild-type mice of C57BL/6J were purchased from Slac Laboratory Animal (Shanghai).  
865 Adult (8-24 weeks old) male mice were used for sexual behavior analysis. All animals  
866 were housed under a reversed 12 h/12 h light-dark cycle with water and food *ad libitum*  
867 in the animal facility at the Institute of Neuroscience, Shanghai, China. All experiments  
868 were approved by the Animal Care and Use Committee of the Institute of Neuroscience,  
869 Chinese Academy of Sciences, Shanghai, China (IACUC No. NA-016-2016).  
870 Male mice were singly housed for at least 3 days prior to sexual behavioral tests. All  
871 tests were initiated at least 1 hr after lights were switched off. Behavioral assays were  
872 recorded using infrared cameras at the frame rate of 30 Hz. Female mice were surgically  
873 ovariectomized and supplemented with hormones to induce receptivity. Hormones were  
874 suspended in sterile sunflower Selfee oil (Sigma-Aldrich, S5007) and injected 10 mg  
875 (in 50 mL oil) and 5 mg (in 50 mL oil) of 17 $\beta$ -estradiol benzoate (Sigma-Aldrich, E8875)  
876 48 h and 24 h preceding the test, respectively. On the day of the test, 50 mg of  
877 progesterone (Sigma-Aldrich, P0130; in 50 mL oil) was injected 4–6 h prior to the test.  
878 Male animals were adapted 10 min to behavioral testing rooms where a recording  
879 chamber equipped with video acquisition systems was located. A hormonal primed  
880 ovariectomized C57BL/6J female (OVX) was introduced to the home cage of male  
881 mice and videotaped for 30 min. Mating behavior tests were repeated three times with  
882 different OVX at least three days apart. Videos were manually scored using a custom-  
883 written MATLAB program. The following criteria were used for behavioral annotation:  
884 active nose contacts initiated by male mouse towards the female's genitals, body area,  
885 faces were defined collectively as “social interest”; male mouse climbs the back of the  
886 female and moves the pelvis were defined as “mount”; Rhythmic pelvic movements

887 after mount were defined as “intromission”; a body rigidity posture after final deep  
888 thrust were defined as “ejaculation”.

## 889 **Data preprocessing, augmentation and sampling**

890 Fly behavior videos were decomposed into frames by FFmpeg, and only the first 10,000  
891 frames of each video were preserved and resized into images with a resolution of  
892 224×224. For model training of *Drosophila* courtship behavior, each video was  
893 manually checked to ensure successful copulations within 10,000 frames.

894 Mice behavior videos were decomposed into frames by FFmpeg, and only frames of  
895 the first 30 min of each video were preserved. Frames were then preprocessed with  
896 OpenCV<sup>4</sup> in Python. Behavior chambers in each video were manually marked,  
897 segmented, and resized into images of a resolution of 256×192. For background  
898 removal, the average frame of each video was subtracted from each frame, and noises  
899 were removed by a threshold of 25 and the median filter with a kernel size of 5. Finally,  
900 the contrast was adjusted with histogram equalization.

901 For data augmentations, crop, rotation, flip, Turbo, and color jitter were applied. For a  
902 given frame, it formed a live-frame with its preceding and succeeding frames. For flies’  
903 behavior video, three frames were successive, and for mice, the preceding or succeeding  
904 frame is one frame away from the current frame due to their slower dynamics<sup>5</sup>. Each  
905 live-frame was randomly cropped into a smaller version containing more than 49%  
906 (70%×70%) of the original image; then the image was randomly (clockwise or  
907 anticlockwise) rotated for an angle smaller than the acute angle formed by the diagonal  
908 line and the vertical line, then the image would be vertically flipped, horizontally  
909 flipped, and/or applied the Turbo lookup table<sup>6</sup> at the probability of 50%, respectively;  
910 and finally, the brightness, contrast, saturation, and hue were randomly adjusted within  
911 10% variation. Notably, since the Turbo transformation is designed for grayscale  
912 images, for a motion-colored RGB image, each channel was transformed individually.  
913 After Turbo transformation, their corresponded channels were composited to form a  
914 new image.

915 For fly data sampling, all images of all videos were randomly ranked, and each batch  
916 contained 256 images from different videos. For mice data sampling, all images of each  
917 video were randomly ranked, and each batch contained 256 images from the same video.  
918 This strategy was designed to eliminate the inconsistency of recording conditions of  
919 mice that was more severe than flies.

920

## 921 **Selfie neural network and its training**

922 All training and inference were accomplished on a workstation with 128GB RAM,  
923 AMD Ryzen 7 5800x, and one NVIDIA GeForce RTX 3090. Selfie neural network was  
924 constructed based on publications and source codes of BYOL<sup>7</sup>, SimSiam<sup>8</sup>, and CLD<sup>9</sup>

925 with PyTorch<sup>10</sup>. In brief, the last layer of ResNet-50 was removed, and a 3-layer 2048-  
926 dimension MLP was added as the projector. Hidden layers of the projector were  
927 followed by batch normalization (BN) and ReLU activation, and the output layer only  
928 had BN. The predictor was constructed with a 2-layer bottleneck MLP with a 512-  
929 dimension hidden layer and a 2048-dimension output layer. The hidden layer but not  
930 the output layer of the predictor had BN and ReLU. As for the group discriminator for  
931 CLD loss, it had only one normalized fully connected layer that projected 2048-  
932 dimension output to 1024 dimensions, followed by a customized normalization layer  
933 that was described in the paper of CLD<sup>9</sup>.

934 The loss function of Selfee had two major parts. The first part was the negative cosine  
935 loss<sup>7,8</sup>, and the second part was the CLD loss<sup>9</sup>. For a batch of  $n$  samples,  $Z$ ,  $P$ ,  $V$   
936 represented the output of projector, predictor, and group discriminator of the main  
937 branch, respectively;  $Z'$ ,  $P'$ ,  $V'$  represented the output of the reference branch; and  $sg$   
938 as the stop-gradient operator. After k-means clustering of  $V$ , the centroids of  $k$  classes  
939 were given by  $M$ , and labels of each sample were provided in the one-hot form as  $L$ .  
940 The hyperparameter  $\theta$  was 0.07, and  $\lambda$  was 2. The loss function was given by the  
941 following equations:

942

$$943 \quad \text{CosineDistance}(m, n) = 1 - \frac{m}{\|m\|_2} \cdot \frac{n}{\|n\|_2}$$

944

$$945 \quad \text{Loss}_1 = \frac{1}{2n} \sum_{i=1}^n \text{CosineDistance}(sg(z_i), p'_i) + \frac{1}{2n} \sum_{i=1}^n \text{CosineDistance}(sg(z'_i), p_i)$$

946

$$947 \quad \text{CrossEntropyLoss}(x, l) = -x \cdot l + \log(\|\exp(x)\|_1)$$

948

$$949 \quad \text{Loss}_2 = \frac{1}{2} \text{CrossEntropyLoss}\left(\frac{v_i}{\theta} \cdot M'^T, l_i\right) + \frac{1}{2} \text{CrossEntropyLoss}\left(\frac{v'_i}{\theta} \cdot M^T, l'_i\right)$$

950

$$951 \quad \text{Loss} = \text{Loss}_1 + \lambda \text{Loss}_2$$

952

953 For all training processes, the Selfee network was trained for 20,000 steps with the SDG  
954 optimizer with a momentum of 0.9 and a weight decay of 1e-4. The learning rate was  
955 adjusted in the one-cycle learning rate policy<sup>11</sup> with base learning rates and a pct start  
956 of 0.025. The model for *Drosophila* courtship behavior was initialized with ResNet-50  
957 pre-trained on the ImageNet, and the base learning rate was 0.025 per batch size of 256.  
958 As for the mating behaviors of mice, the model was initialized with weights trained on  
959 the fly dataset, and the base learning rate was 0.05 per batch size of 256.

960

## 961 **t-SNE visualization**

962 Video frames for t-SNE visualization were all processed by Selfee. Embeddings of three  
963 tandem frames were averaged to eliminate potential noises. All embeddings were  
964 transformed using t-SNE provided in the scikit-learn<sup>12</sup> package in Python without  
965 further turning of parameters. Results were visualized with the Matplotlib<sup>13</sup> package in  
966 Python, and their colors were assigned based on human annotations of video frames.  
967

## 968 **Classification**

969 Two kinds of classification methods were implemented, including the k-NN classifier  
970 and the LightGBM (Light Gradient Boosting Machine) classifier. The weighed k-NN  
971 classifier was constructed based on the previous reports<sup>8,14</sup>. LightGBM classifier<sup>15</sup> was  
972 provided by its official package in Python. The F<sub>1</sub> score and mAP were calculated with  
973 the scikit-learn<sup>12</sup> package in Python.

974 For fly behavior classification, seven 10,000-frame videos were annotated manually.  
975 Seven-fold cross-validation was performed using embeddings generated by Selfee and  
976 the k-NN classifier. Inferred labels were forced to be continuous through time by using  
977 inferred labels of 21 neighbor frames to determine the final result. Then, a video  
978 independent of the cross-validation was annotated and inferred by a k-NN classifier  
979 using all 70,000 samples, and the last 3,000 frames were used for the raster plot.

980 For rat behavior classification, the RatSI dataset<sup>16</sup> (a kind gift from Noldus Information  
981 Technology bv) contains nine manually annotated videos. We neglected three rarest  
982 annotated behaviors: moving away, nape attacking, and pinning, and we combined  
983 approaching and following into a larger category. Therefore, we used five kinds of  
984 behaviors, including allogrooming, approaching or following, social nose contact,  
985 solitary, and others. Nine-fold cross-validation was performed using embeddings  
986 generated by Selfee and the k-NN classifier. Inferred labels were forced to be  
987 continuous through time by using inferred labels of 81 neighbor frames to determine  
988 the final result.

989 For mice behavior classification, eight videos were annotated manually. Eight-fold  
990 cross-validation was performed using embeddings generated by Selfee and the k-NN  
991 classifier. To incorporate more temporal information, the LightGBM classifier and  
992 additional features were also used. Additional features include slide moving average  
993 and standard division of 81-frame time windows, the main frequencies, and their energy  
994 (using short-time Fourier transform in SciPy<sup>17</sup>) within 81-frame time windows. Early-  
995 stop was used to prevent over-fitting. Inferred labels were forced to be continuous  
996 through time by using inferred labels of 81 neighbor frames to determine the final result.  
997 Then, a video independent of the cross-validation was annotated and inferred by an  
998 ensemble classifier of eight previously constructed classifiers, and all frames were used  
999 for the raster plot.

1000

## 1001 **Anomaly detection**

1002 For a group of query embeddings of sequential frames  $q_1, q_2, q_3, \dots, q_n$ , and a group of  
1003 reference embeddings of sequential frames  $r_1, r_2, r_3, \dots, r_m$ , the anomaly score of each  
1004 query frame was given by the following equation:

1005

$$1006 \quad \text{AnomalyScore}(q_i) = \min_{j=1}^m (\text{CosineDistance}(q_i, r_j)) - \min_{|j-i| < 50} (\text{CosineDistance}(q_i, q_j))$$

1007

1008 A PyTorch implementation of cosine similarity<sup>18</sup> was used for accelerated calculations.

1009 The anomaly score of each video was the average anomaly score of the top 100

1010 anomalous frames. The statistical analysis of the genetic screening was performed with

1011 one-way ANOVA with Benjamini and Hochberg correction in GraphPad Prism

1012 Software.

1013 If negative controls are provided, anomalous frames are defined as frames with higher

1014 anomaly scores than the maximum anomaly score of frames in negative control videos.

1015

## 1016 **Autoregressive hidden Markov model (AR-HMM)**

1017 All AR-HMM models were built with the implementation of MoSeq<sup>5</sup>

1018 (<https://github.com/mattjj/pyhsmm-autoregressive>). A principal component analysis

1019 (PCA) model that could explain 95% of variance of the control group was built and

1020 used to transform both control and experiment groups. The max module number was

1021 set as 10 for all experiments unless indicated otherwise. Each model was sampled for

1022 1000 iterations. We kept other hyperparameters the same as the examples provided by

1023 this package. State usages of each module in control and experimental groups were

1024 analyzed by Mann Whitney test with SciPy<sup>17</sup> followed with Benjamini and Hochberg

1025 correction. The state usages were also visualized after PCA dimensional reduction with

1026 scikit-learn<sup>12</sup> and Matplotlib<sup>13</sup>.

1027

## 1028 **Dynamic time warping (DTW)**

1029 Dynamic time warping was modified from the Python implementation<sup>19</sup>

1030 (<https://dynamictimewarping.github.io/python/>). Specifically, PyTorch implementation

1031 of cosine similarity<sup>18</sup> was used for accelerated calculations.

1032

## 1033 Pose-estimation with DeepLabCut

1034 We used the official implementation of DeepLabCut<sup>20,21</sup>. For training, 120 frames of a  
1035 mating behavior video were labeled manually, and 85% of them were used as the  
1036 training set. Marked body parts included nose, ears, body center, hips, and bottom,  
1037 following previous publications<sup>22,23</sup>. The model (ResNet-50 as the backbone) was  
1038 trained for 100,000 iterations, with a batch size of 16. We kept other hyperparameters  
1039 the same as default settings.  
1040

## 1041 REFERENCES

1042

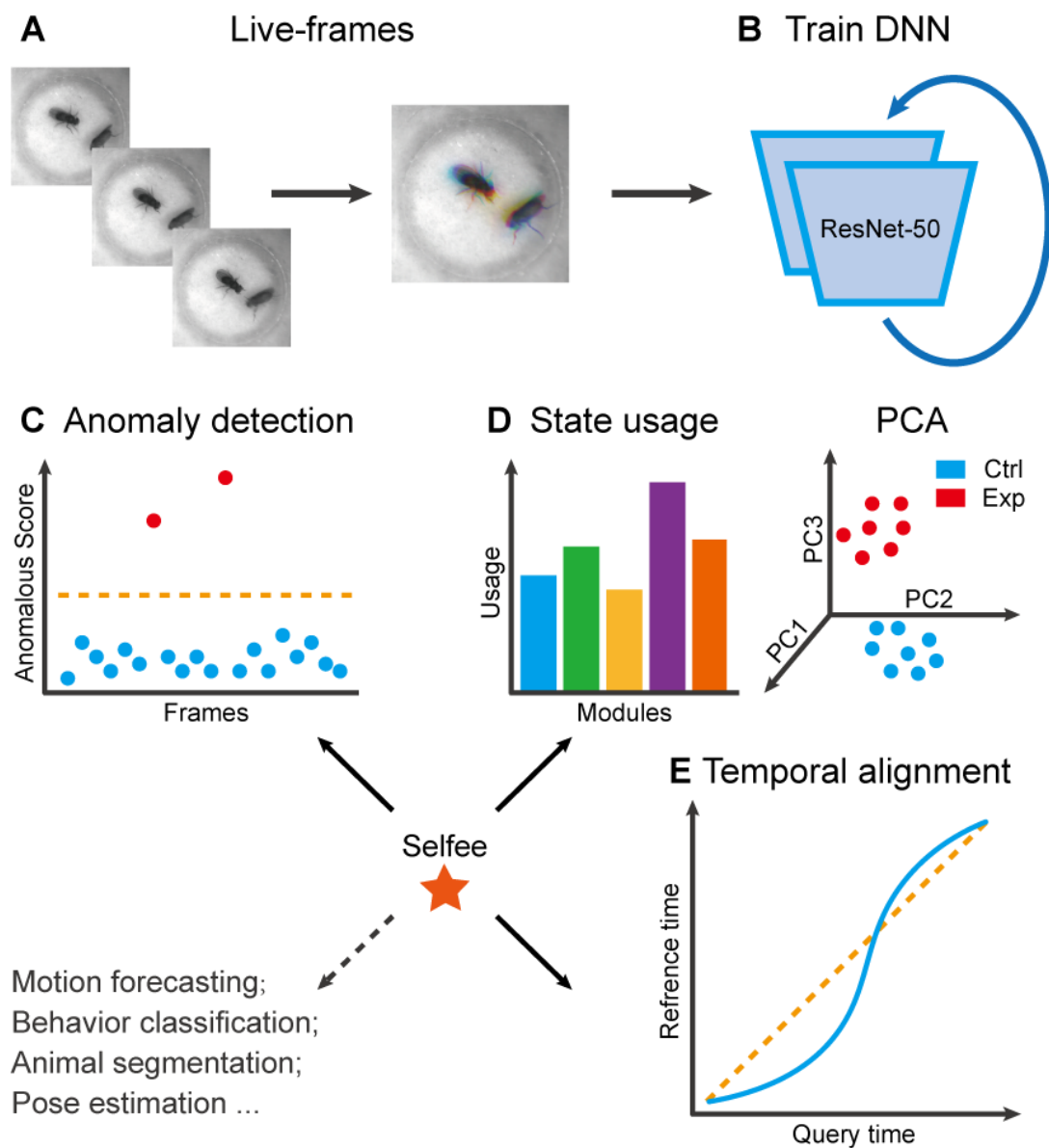
- 1043 1 Li, J. *et al.* A Defensive Kicking Behavior in Response to Mechanical Stimuli  
1044 Mediated by *Drosophila* Wing Margin Bristles. *The Journal of Neuroscience* **36**,  
1045 11275-11282, doi:10.1523/jneurosci.1416-16.2016 (2016).
- 1046 2 Liu, C. *et al.* A neural circuit encoding mating states tunes defensive behavior  
1047 in *Drosophila*. *Nature Communications* **11**, 3962, doi:10.1038/s41467-020-  
1048 17771-8 (2020).
- 1049 3 Au - McNeil, A. R. *et al.* Conditions Affecting Social Space in *Drosophila*  
1050 *melanogaster*. *JoVE*, e53242, doi:doi:10.3791/53242 (2015).
- 1051 4 Bradski, G. The OpenCV Library. *Dr. Dobb's Journal of Software Tools*.
- 1052 5 Wiltschko, Alexander B. *et al.* Mapping Sub-Second Structure in Mouse  
1053 Behavior. *Neuron* **88**, 1121-1135,  
1054 doi:<https://doi.org/10.1016/j.neuron.2015.11.031> (2015).
- 1055 6 Mikhailov, A. *Turbo, An Improved Rainbow Colormap for Visualization*,  
1056 <[https://ai.googleblog.com/2019/08/turbo-improved-rainbow-colormap-](https://ai.googleblog.com/2019/08/turbo-improved-rainbow-colormap-for.html)  
1057 [for.html](https://ai.googleblog.com/2019/08/turbo-improved-rainbow-colormap-for.html)> (2019).
- 1058 7 Grill, J.-B. *et al.* Bootstrap Your Own Latent: A New Approach to Self-  
1059 Supervised Learning. *ArXiv* **abs/2006.07733** (2020).
- 1060 8 Chen, X. & He, K. Exploring Simple Siamese Representation Learning. *ArXiv*  
1061 **abs/2011.10566** (2020).
- 1062 9 Wang, X., Liu, Z. & Yu, S. X. Unsupervised Feature Learning by Cross-Level  
1063 Instance-Group Discrimination. *arXiv: Computer Vision and Pattern*  
1064 *Recognition* (2020).
- 1065 10 Paszke, A. *et al.* in *NeurIPS*.
- 1066 11 Smith, L. N. & Topin, N. Super-Convergence: Very Fast Training of Neural  
1067 Networks Using Large Learning Rates. *arXiv:1708.07120* (2017).  
1068 <<https://ui.adsabs.harvard.edu/abs/2017arXiv170807120S>>.
- 1069 12 Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *J. Mach. Learn.*  
1070 *Res.* **12**, 2825-2830 (2011).
- 1071 13 Hunter, J. D. Matplotlib: A 2D Graphics Environment. *Computing in Science &*



- 1072            *Engineering* **9**, 90-95, doi:10.1109/MCSE.2007.55 (2007).
- 1073    14    Wu, Z., Xiong, Y., Yu, S. X. & Lin, D. Unsupervised Feature Learning via Non-  
1074            Parametric Instance-level Discrimination. *ArXiv* **abs/1805.01978** (2018).
- 1075    15    Ke, G. *et al.* in *Proceedings of the 31st International Conference on Neural*  
1076            *Information Processing Systems*    3149–3157 (Curran Associates Inc., Long  
1077            Beach, California, USA, 2017).
- 1078    16    Lorbach, M. *et al.* Learning to recognize rat social behavior: Novel dataset and  
1079            cross-dataset application. *Journal of Neuroscience Methods* **300**, 166-172,  
1080            doi:<https://doi.org/10.1016/j.jneumeth.2017.05.006> (2018).
- 1081    17    Virtanen, P. *et al.* SciPy 1.0: fundamental algorithms for scientific computing in  
1082            Python. *Nature Methods* **17**, 261-272, doi:10.1038/s41592-019-0686-2 (2020).
- 1083    18    Musgrave, K., Belongie, S. J. & Lim, S.-N. PyTorch Metric Learning. *ArXiv*  
1084            **abs/2008.09164** (2020).
- 1085    19    Toni, G. Computing and Visualizing Dynamic Time Warping Alignments in R:  
1086            The dtw Package. *Journal of Statistical Software* **31**, doi:10.18637/jss.v031.i07  
1087            (2009).
- 1088    20    Mathis, A. *et al.* DeepLabCut: markerless pose estimation of user-defined body  
1089            parts with deep learning. *Nature Neuroscience* **21**, 1281-1289,  
1090            doi:10.1038/s41593-018-0209-y (2018).
- 1091    21    Lauer, J. *et al.* Multi-animal pose estimation and tracking with DeepLabCut.  
1092            *bioRxiv*, 2021.2004.2030.442096, doi:10.1101/2021.04.30.442096 (2021).
- 1093    22    Segalin, C. *et al.* The Mouse Action Recognition System (MARS): a software  
1094            pipeline for automated analysis of social behaviors in mice. *bioRxiv*,  
1095            2020.2007.2026.222299, doi:10.1101/2020.07.26.222299 (2020).
- 1096    23    Sun, J. J. *et al.* in *2021 IEEE/CVF Conference on Computer Vision and Pattern*  
1097            *Recognition (CVPR)*.    2876-2885 (2021).
- 1098

1099 **FIGURES AND TABLES**

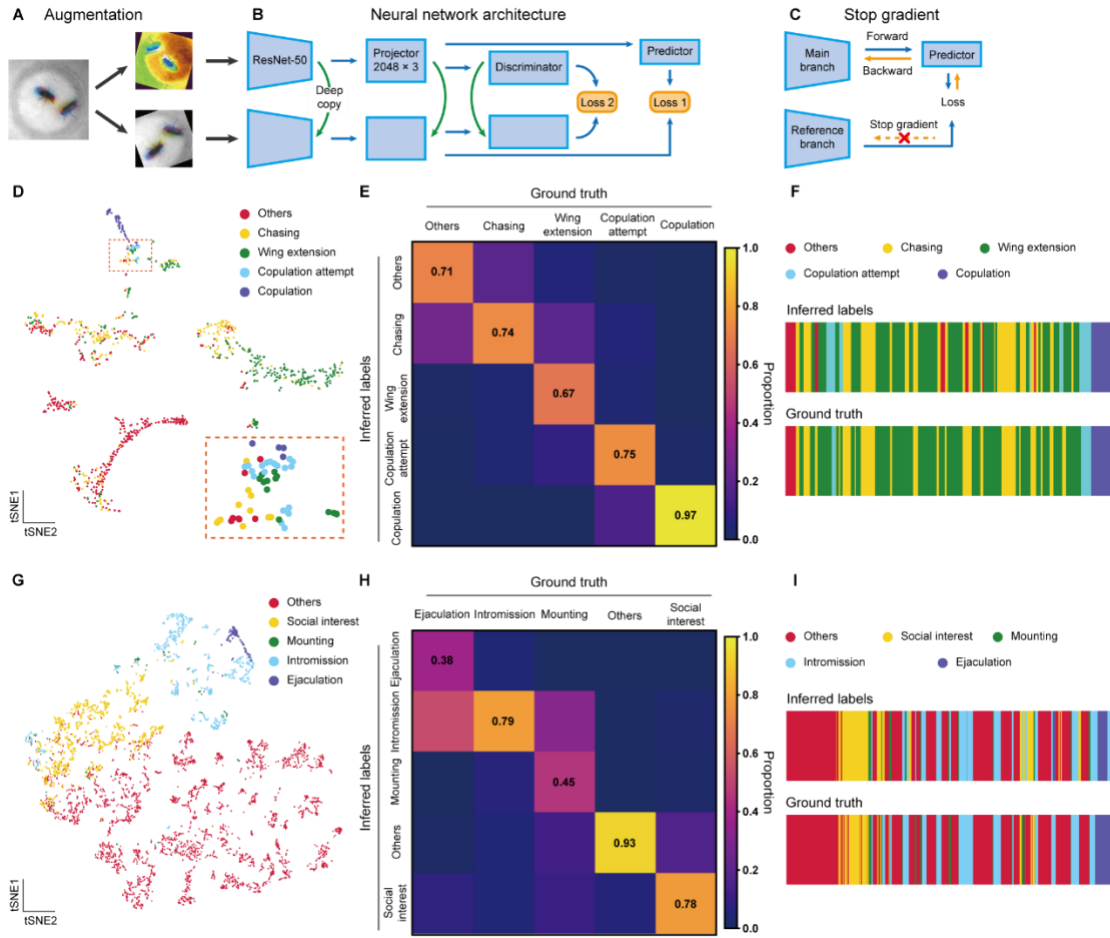
1100 **Figure 1**



1101

1102

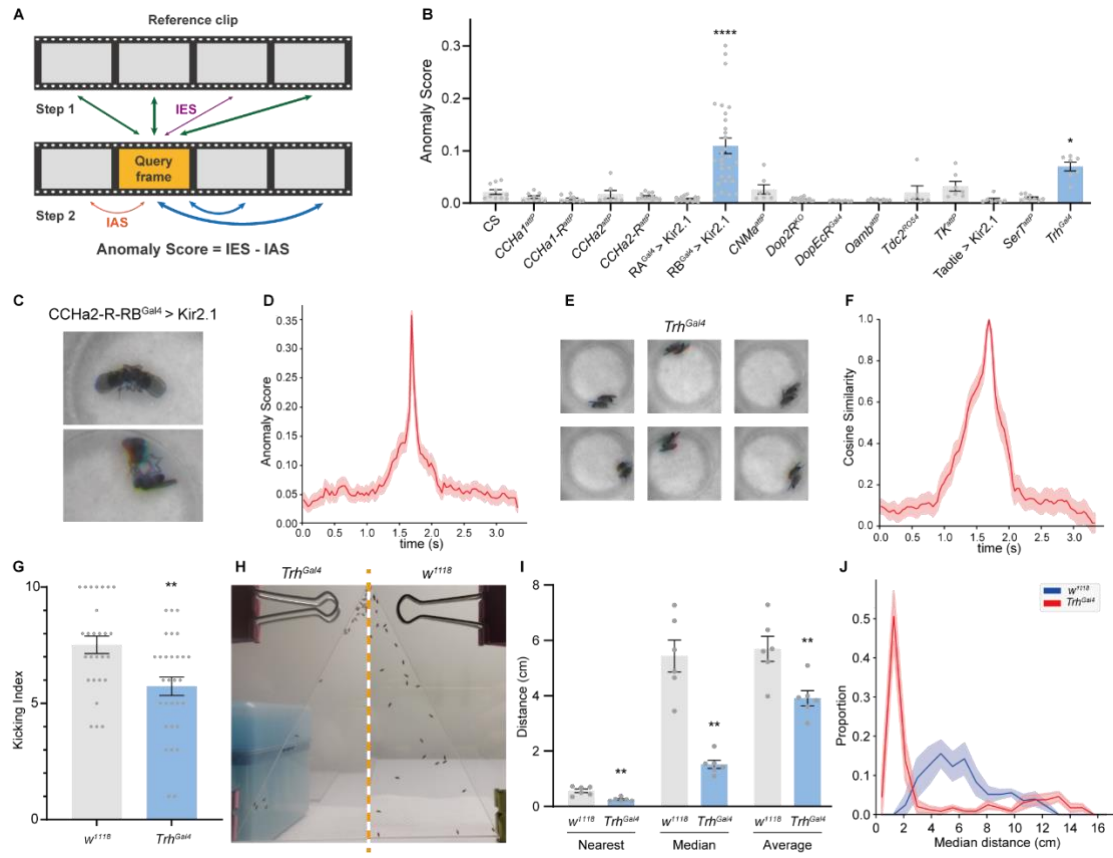
1103 **Figure 2**



1104

1105

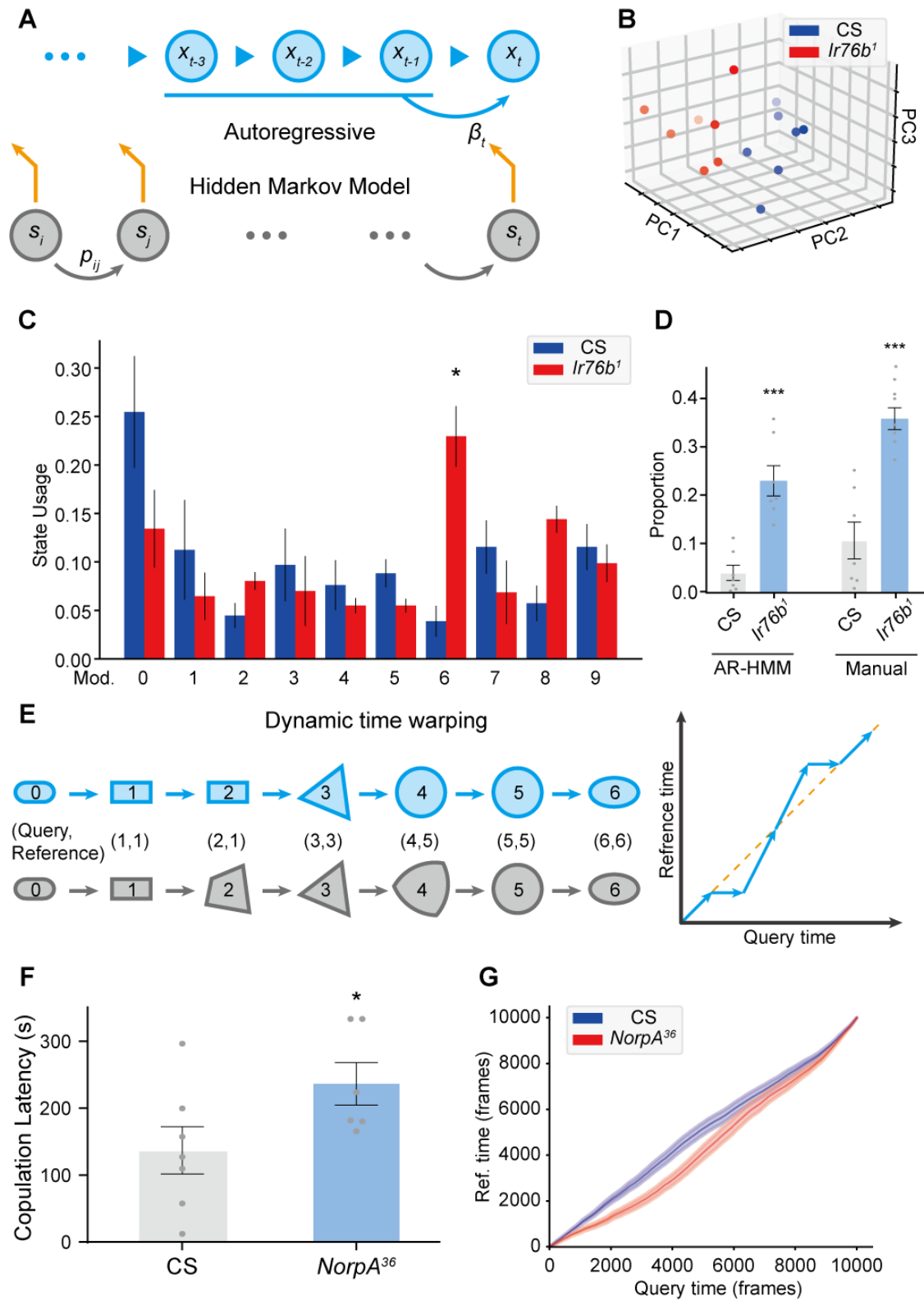
1106 **Figure 3**



1107

1108

1109 **Figure 4**



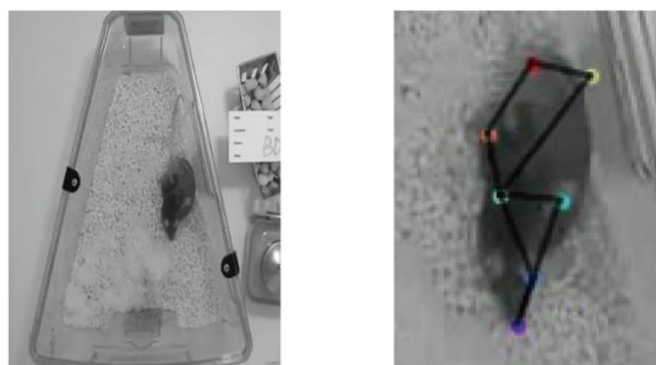
1110

1111

1112

1113 **Figure 1—figure supplement 1**

**A**



Ground Truth

Intromission

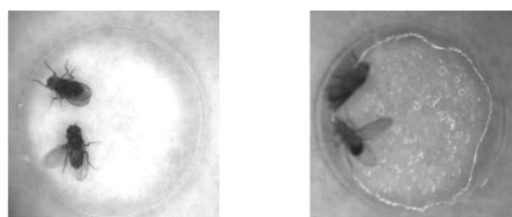
DeepLabCut

1114 Selfee with LightGBM

Intromission

1115 **Figure 1—figure supplement 2**

**A**



Ground Truth

Wing extension

Wing extension

k-NN classification

Wing extension

Copulation

**B**



Original

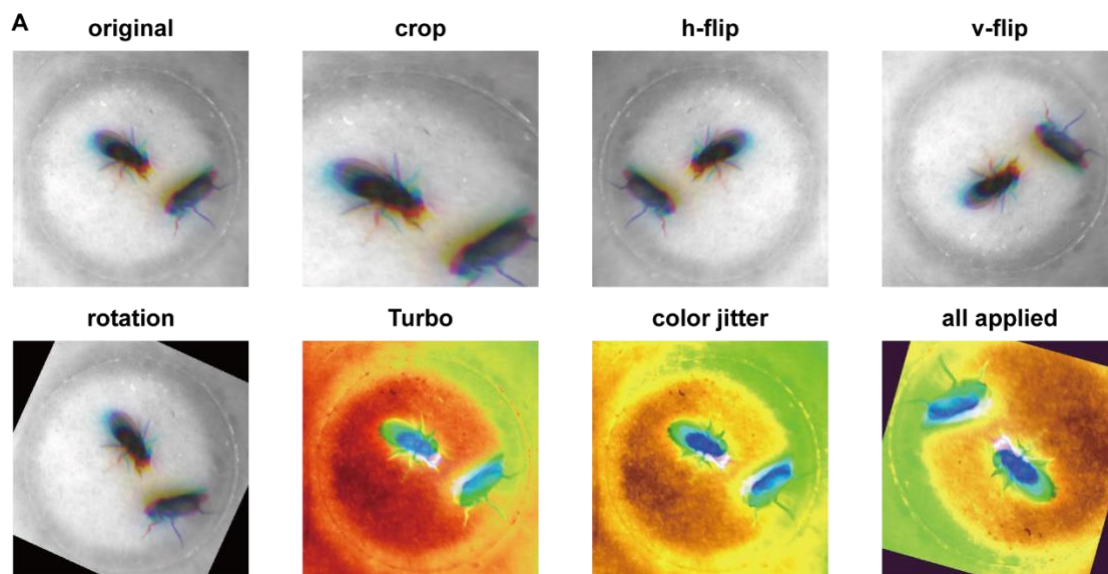
Background

Mice

1116

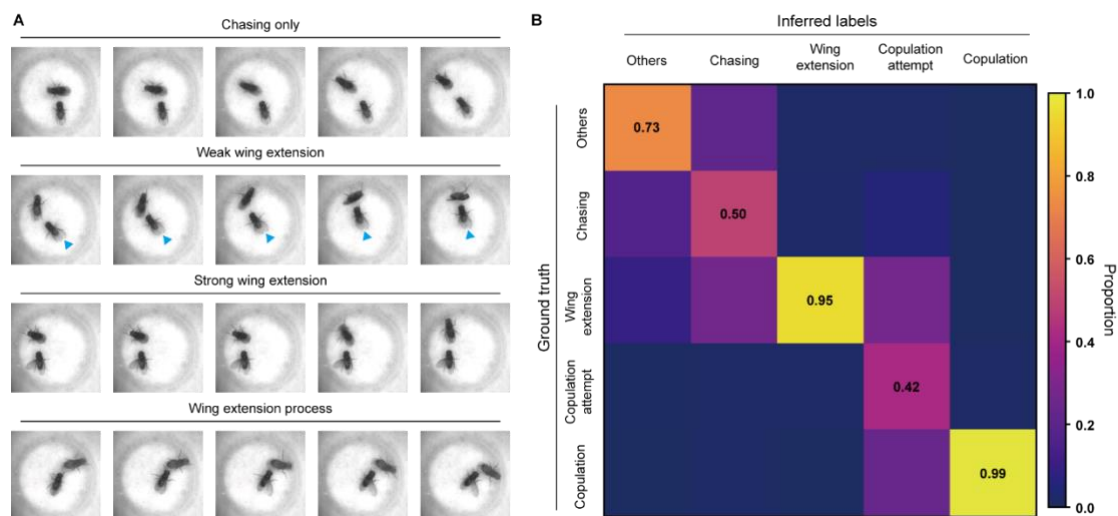
1117

1118 **Figure 2—figure supplement 1**



1119

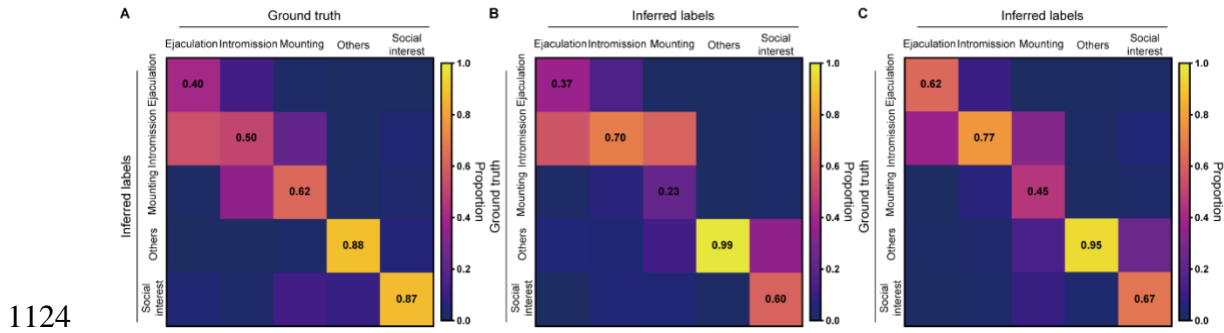
1120 **Figure 2—figure supplement 2**



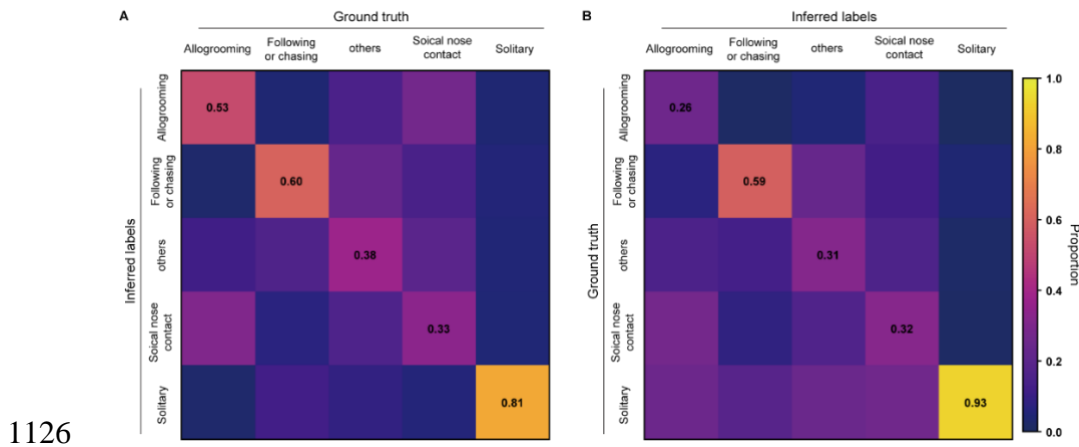
1121

1122

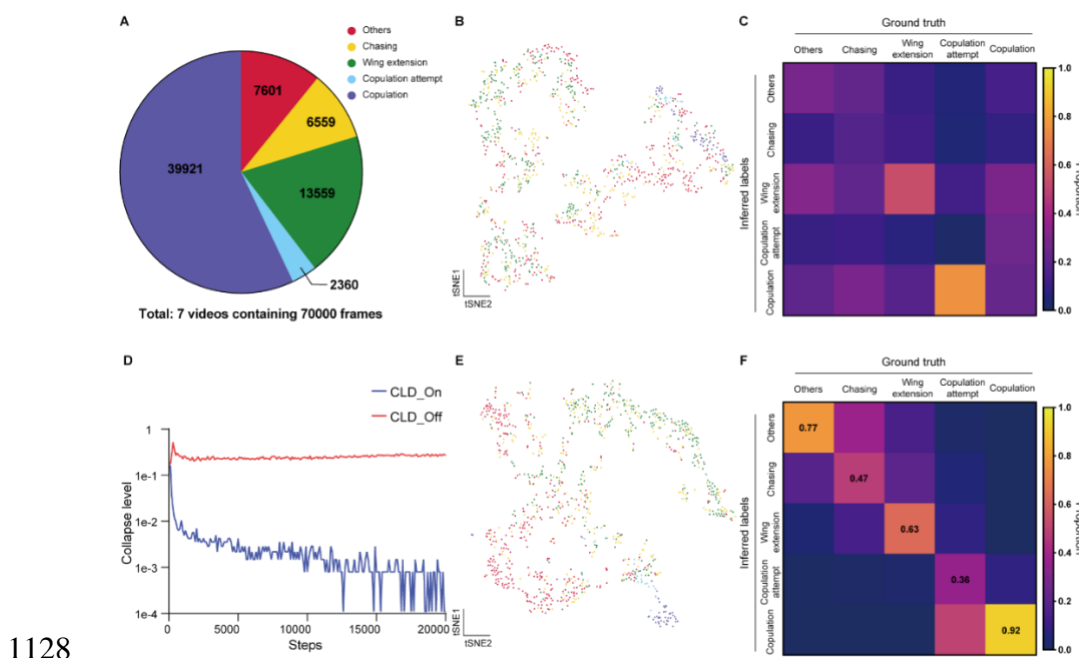
1123 **Figure 2—figure supplement 3**



1125 **Figure 2—figure supplement 4**



1127 **Figure 2—figure supplement 5**





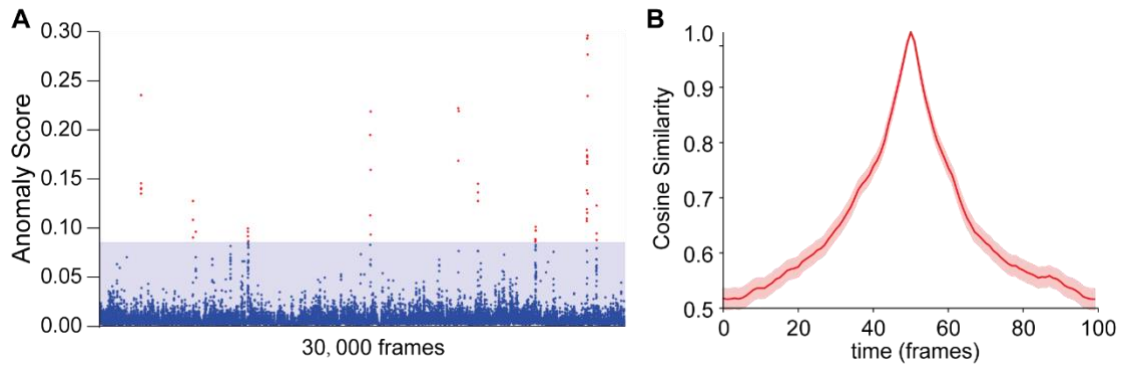
1129 **Figure 2—figure supplement 6 | Ablation test of Selfee training process on fly datasets.**

1130

Model	Pretrained ResNet-50 with random projectors		Selfee		Selfee without CLD loss		Selfee without Turbo transformation	
Evaluation	Mean F1 Score	Mean AP	Mean F1 Score	Mean AP	Mean F1 Score	Mean AP	Mean F1 Score	Mean AP
Replication 1	0.586	0.580	0.724	0.758	0.227	0.227	0.604	0.550
Replication 2	0.597	0.570	0.676	0.683	0.163	0.200	0.574	0.551
Replication 3	0.596	0.586	0.714	0.754	0.172	0.214	0.517	0.497
Best	0.597	0.586	<b>0.724</b>	<b>0.758</b>	0.227	0.227	0.604	0.551

1131

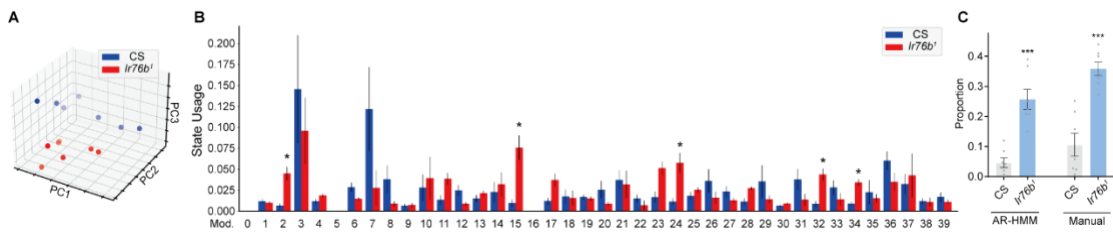
1132 **Figure 3—figure supplement 1**



1133

1134

1135 **Figure 4—figure supplement 1**



1136