

Cross-species Behavior Analysis with Attention-based Domain-adversarial Deep Neural Networks

Takuya Maekawa (✉ maekawa@ist.osaka-u.ac.jp)

Osaka University <https://orcid.org/0000-0002-7227-580X>

Daiki Higashide

Osaka University

Takahiro Hara

Osaka University

Kentarou Matsumura

Okayama University

Kaoru Ide

Doshisha University

Takahisa Miyatake

Okayama University <https://orcid.org/0000-0002-5476-0676>

Koutarou Kimura

Nagoya City University <https://orcid.org/0000-0002-3359-1578>

Susumu Takahashi

Doshisha University

Article

Keywords: behavior analysis, deep neural networks, cross-species features

Posted Date: December 15th, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-123107/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at Nature Communications on September 17th, 2021. See the published version at <https://doi.org/10.1038/s41467-021-25636-x>.

Cross-species Behavior Analysis with Attention-based Domain-adversarial Deep Neural Networks

Takuya Maekawa^{1,*}, Daiki Higashide¹, Takahiro Hara¹, Kentarou Matsumura², Kaoru Ide³, Takahisa Miyatake², Koutarou D. Kimura⁴ & Susumu Takahashi³

¹*Graduate School of Information Science and Technology, Osaka University, Osaka 565-0871, Japan*

²*Graduate School of Environmental and Life Science, Okayama University, Okayama 700-8530, Japan*

³*Graduate School of Brain Science, Doshisha University, Kyoto 610-0321, Japan*

⁴*Graduate School of Natural Sciences, Nagoya City University, Aichi 467-8501 Japan*

**Correspondence to Takuya Maekawa.*

Since the variables inherent to various diseases cannot be controlled directly in humans, behavioral dysfunctions have been examined in model organisms, leading to better understanding their underlying mechanisms. However, because the spatial and temporal scales of animal locomotion vary widely among species, conventional statistical analyses cannot be used to discover knowledge from the locomotion data. We propose a new procedure to automatically discover locomotion features shared among animal species by means of domain-adversarial deep neural networks. Our neural network is equipped with a function which explains the meaning of segments of locomotion where the cross-species features are hidden by incorporating an attention mechanism into the neural network, regarded as a black box.

21 **It enables us to formulate a human-interpretable rule about the cross-species locomotion fea-**
22 **ture and validate it using statistical tests. We demonstrate the versatility of this procedure**
23 **by identifying locomotion features shared across different species with dopamine-deficiency,**
24 **namely humans, mice, and worms, despite their evolutionary differences.**

25 **INTRODUCTION**

26 Neurodegenerative diseases including Parkinson's disease (PD), Alzheimer's disease, and schiz-
27 ophrenia are disorders characterized by motor dysfunctions. Since the variables inherent to such
28 diseases cannot be controlled directly in humans, behavioral dysfunctions and their neural under-
29 pinnings have been examined in model organisms^{1,2}. Assuming that fundamental aspects of the
30 behavior of humans are evolutionarily conserved among other animal species, studies in model
31 organisms gained insights into understanding the underlying mechanisms of those diseases³⁻⁶. In
32 contrast to cognitive abnormalities, motor dysfunctions can be externally assessed by comparative
33 behavioral analyses. A central concept of comparative behavioral analysis is to identify human-
34 like behavioral repertoires, called behavioral phenotype, in animals, as animal models of PD have
35 gained insight into understanding the behavioral and neural underpinning of symptoms underlying
36 a specific disease⁷ and can potentially provide clues for the development of therapeutics. How-
37 ever, behavioral analysis across animal species could not be realized using conventional statistical
38 analyses because the body scale and locomotion methods vary among species (Fig. 1a), resulting
39 in wide variations in the spatial and temporal scales of locomotion among species as shown in Fig.
40 1b in which worms, beetles, and mice show similar locomotion trajectory, but differ in the velocity
41 and spatial area occupied.

42 Deep learning is a novel technique of automated feature extraction, reducing the cost of
43 manual feature design. We employ it for automatically discovering scale-invariant locomotion
44 features shared among animal species. In the procedure (see Fig. 1c), a human operator first feeds
45 locomotion data from animals of different species with different properties into a neural network
46 designed to extract cross-species locomotion features using domain-adversarial training⁸. This is
47 trained to extract features that classify a trajectory into an appropriate class (e.g., healthy or PD) but
48 cannot label the trajectory into an appropriate domain (i.e., species), using a gradient reversal layer
49 as shown in Fig. 1d. Because these features are incapable of distinguishing between domains,
50 we can regard them as species-independent. In contrast, because we can distinguish between
51 trajectories belonging to different classes using the extracted features, these can be regarded as
52 cross-species hallmarks of the diseases. By designing a deep neural network so that it extracts
53 features based on the above idea, we can obtain locomotion features shared across different species
54 independent of their body scales and locomotion methods.

55 Despite the human-level outstanding performances of the neural network, the human op-
56 erator cannot understand the meaning of the extracted cross-species locomotion feature by the
57 deep neural network containing a huge amount of hidden parameters. To address this issue, we
58 design an explainable architecture by incorporating an attention mechanism^{9,10} into the domain-
59 adversarial neural network, which identifies segments in the trajectories where the cross-species
60 features are hidden and provides visualized trajectories and time-series of basic locomotion fea-
61 tures (e.g., speed) by highlighting the identified segments (Fig. 1e,f). From these highlighted
62 graphs, the operator can understand cross-species locomotion features extracted by the neural net-

63 work. For instance, short-duration peaks in speed are characteristic to PD mice (Fig. 1f). To
64 explain why the neural network pays attention to a certain segment and how it labels an input
65 time-series using attended segments within the time-series, we employ decision trees. They help
66 formulate a human-interpretable rule about the cross-species locomotion feature. Then, the op-
67 erator proposes a hypothesis related to the locomotion features and performs a statistical test to
68 validate it (see Fig. 1c).

69 To demonstrate the performance of our proposed cross-species behavioral analysis, we iden-
70 tified locomotion features shared across different species with dopamine-deficiency, namely hu-
71 mans, mice, and worms.

72 RESULTS

73 **Attention-based Domain-adversarial Neural Network** This study assumes that locomotion data
74 from two different species that belong to two different classes, for example, PD individuals and
75 healthy individuals of humans and mice are given. The locomotion data are used to train the neural
76 network (see Fig. 1d). The model inputs are time-series of primitive locomotion features such as
77 speed. The model has two types of outputs: estimated domain and class of an input time-series.

78 The convolutional layers in the feature extraction block are used to extract features, which
79 are used to output the two estimates. To make the neural network incapable of distinguishing be-
80 tween the two domains, we introduce the gradient reversal layer⁸ before the 1st domain predictor.
81 When we train the network using the backpropagation algorithm¹¹, the gradient reversal layer mul-
82 tiplies the gradient with a negative constant value, making the convolutional layers in the feature

83 extraction block incapable of estimating domains but classes.

84 In addition, to enable the user to understand which segments within input time-series the
85 neural network focuses on in order to discriminate between two classes, we introduce an attention
86 mechanism⁹ into the model. The attention of a data point at each time slice is regarded as impor-
87 tance of the data point when the input time-series is classified, indicating that attended data points
88 are characteristic to a class to which the input time-series belongs. The convolutional layers in
89 the attention computation block are used to compute attention for each time slice. The attention is
90 multiplied by the extracted features to contrast the data points to which the network pays attention.
91 Furthermore, to make the way of attention computation domain-independent, we introduce the gra-
92 dient reversal layer after the 2nd convolutional layer in the attention computation block. The 2nd
93 domain predictor outputs domain estimate for each time slice using the output of the 2nd convolu-
94 tional layer in the attention computation block and we make the network incapable of classifying
95 the output of the 2nd convolutional layer into domains using the gradient reversal layer.

96 The network is trained to minimize the error of the class estimates as well as maximize the
97 errors of the two types of the domain estimates. However, achieving these conflicting goals at the
98 same time makes it difficult for the neural network to converge. Therefore, we iterate the following
99 procedures to train the network.

- 100 1. We first train the network to minimize the error of the class estimates as well as maximize
101 the error of the domain estimates by the 2nd domain predictor, which employs outputs of the

102 attention computation block to predict the domain. Thus we extract important segments in
 103 the input time-series for class estimation in the domain-independent manner. This procedure
 104 consists in minimizing

$$E(\theta_f, \theta_a, \theta_c, \theta_{d2}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_c^i(\theta_f, \theta_a, \theta_c) - \lambda_1 \frac{1}{n} \sum_{i=1}^n \mathcal{L}_{d2}^i(\theta_a, \theta_{d2}), \quad (1)$$

105 where $\theta_f, \theta_a, \theta_c, \theta_{d2}$ represent network parameters of the feature extraction block, attention
 106 computation block, class predictor, and 2nd domain predictor, respectively, n is the num-
 107 ber of training instances (time-series), $\mathcal{L}_c^i(\theta_f, \theta_a, \theta_c)$ shows the loss of class prediction, and
 108 $\mathcal{L}_{d2}^i(\theta_a, \theta_{d2})$ is the loss of domain prediction by the 2nd domain predictor. The hyperpareme-
 109 ter λ_1 is used to trade-off between the two loss functions. For more details, see **Methods**.

110 2. Then we update the network to maximize the error of domain estimates by the 1st domain
 111 predictor to accelerate training in the next procedure. This procedure consists in minimizing

$$E(\theta_f, \theta_a, \theta_{d1}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_{d1}^i(\theta_f, \theta_a, \theta_{d1}), \quad (2)$$

112 where θ_{d1} is network parameter of the 1st domain predictor and $\mathcal{L}_{d1}^i(\theta_f, \theta_a, \theta_{d1})$ is the loss of
 113 domain prediction by the 1st domain predictor.

114 3. Finally, we update the network to minimize the error of class estimates as well as maxi-
 115 mize the error of domain estimates by the 1st domain predictor. This procedure consists in
 116 minimizing

$$E(\theta_f, \theta_a, \theta_c, \theta_{d1}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_c^i(\theta_f, \theta_a, \theta_c) - \lambda_2 \frac{1}{n} \sum_{i=1}^n \mathcal{L}_{d1}^i(\theta_f, \theta_a, \theta_{d1}), \quad (3)$$

117 where λ_2 is a trade-off hyperparameter of the two loss functions. With this procedure, we
 118 can train the network to predict class labels correctly but domain labels incorrectly.

119 In addition, to help understand the meaning of the attention, we employ machine learning
120 tools to explain it. We construct a decision tree that is trained to detect the attended segments by
121 the neural network using exhaustively-handcrafted features listed in Supplementary Table 1.

122 Furthermore, to interpret how the neural network predicts classes using attended segments,
123 we build a decision tree that is trained to classify an input time-series using handcrafted features
124 extracted only from attended segments within the time-series. See **Methods** for more details.

125 As above, an interpretable rule is extracted from our neural network with the explainable
126 architecture and/or by our methods based on decision trees. After that, the interpretable rule is
127 investigated via a statistical test. To validate our procedure, we prepared locomotion data from
128 four different dopamine-deficient species (humans, mice, beetles, and worms). We train our neural
129 network on data from two species to obtain an interpretable rule and then validate this using data
130 from all the species. We only use data from two species because doing otherwise would increase
131 the complexity of the classification task, making it difficult for the network to converge.

132 **Network Training on Worm and Mouse Data** We use locomotion data from worms with/without
133 DOP-3, one of the D2 dopamine receptors¹² and from healthy and PD mice. A PD mouse is a uni-
134 lateral 6-hydroxydopamine (OHDA) lesioned model of PD whose dopaminergic neurons in the
135 compact part of substantia nigra of a left or right hemisphere are lesioned with a neurotoxin 6-
136 OHDA. Fig. 2a, b show the data collection protocols. Prior studies discovered relevant locomotion
137 features of PD mice^{1,2} as well as of worms during odor avoidance behavior^{13,14}. However, locomotion
138 features across these species have not been investigated because the scales of their movements

139 are significantly different. We convert 2D locomotion data from the worms and mice into time-
140 series of locomotion speed and then feed them to the neural network after standardization. More
141 details about the data are available in **Methods** and Supplementary Table 2.

142 We randomly split the speed time-series data into training and test sets, which are fed into
143 the neural network. We deal with two classes, DA(+) class (healthy mice and worms) vs. DA(-)
144 class (PD mice and worms lacking D2 dopamine receptors), and two domains (mouse and worm).
145 The classification accuracy is shown in Fig. 2c, d, indicating that our network cannot distinguish
146 the domains but the classes. These results indicate that cross-species features of DA(-) among
147 mice and worms exist. Fig. 2e shows an example of time-series of speed for worms and mice
148 highlighted by the trained model. From these highlighted time-series, we can speculate that the
149 neural network focuses on segments corresponding to high-speed for the DA(+) worms and mice.

150 Fig. 2f shows a decision tree trained to detect the attended segments. The root node and
151 its right child node indicate that attended segments correspond to segments with long-lasting high
152 speed.¹ Fig. 2g shows a decision tree trained to classify an input time-series into an appropriate
153 class using features extracted from attended segments within the time-series. As shown in the
154 root node and Fig. 2f, the neural network seems to pay attention to long-lasting high speed of the
155 worms/mice using the attention mechanism, and then distinguishes between DA(+) and DA(-) by
156 the skewness of speed within the attended segments. The attended segments with positive skewness
157 (≥ 0.360) are classified into the DA(-) class. In contrast, the attended segments with skewness

¹High minimum speed in a sliding window, i.e., high moving minimum of speed, indicates keeping high speed.

158 smaller than 0.360 are classified into the DA(+) class (between -0.397 and 0.360). Therefore, these
159 results indicate that the DA(+) worm/mouse keeps stable speed (small skewness) when moving in
160 high speed.

161 As above, we could extract an interpretable rule from the trained neural network, i.e., DA(-)
162 worms/mice cannot keep high speed. To validate the hypothesis, we perform a statistical test by
163 computing the minimum speed within a time window when the speed is high (see Supplemen-
164 tary Information for details). As shown in Fig. 2h, we observe significant differences between
165 DA(+) and DA(-) for both the mice and worms. Surprisingly, we could also observe significant
166 differences between PD and healthy humans when we performed the same test (see Supplemen-
167 tary Information), even though the neural network was trained only on the data from the mice and
168 worms.

169 As above, our method revealed that humans, mice, and worms with disabilities in the dopamin-
170 ergic system cannot keep high speed even though their body scales and locomotion methods are
171 completely different. Our study employs worms with a lack of D2 dopamine receptors. Although
172 the PD symptoms are considered to be induced by various impairment of neural circuits such as
173 dopamine transmission impairment, morphological alterations of the basal ganglia circuitry, and
174 lack of dopamine receptors¹⁵⁻¹⁸, a major source of the discovered locomotion feature shared by
175 PD mice and humans might ascribe the lack of D2 dopamine receptors. The verification of this
176 hypothesis is beyond the scope of this study.

177 **Network Training on Worm and Human Data** Next, we show results obtained when the neural
178 network is trained on data from worms and humans. As for the human data (see **Methods** and
179 Supplementary Table 2), we convert time-series of foot-mounted pressure sensors collected during
180 walking into time-series of locomotion speed.

181 We deal with two classes, i.e., DA(+) class (healthy worms and humans) vs. DA(-) class
182 (PD humans and worms lacking D2 dopamine receptors). Fig. 3a, b show the classification accu-
183 racy, indicating that our network seems to extract a cross-species feature between the humans and
184 worms. Fig. 3c shows the time-series of speed for DA(+)/DA(-) worms and humans highlighted
185 by the trained model. From these highlighted time-series, the neural network seems to focus on
186 segments corresponding to smooth acceleration for the DA(+) worms and humans. (A decision
187 tree shown below also takes an acceleration feature as a root node.) Fig. 3d shows the time-series
188 of acceleration for worms and humans, as segments of DA(+) worms/humans corresponding to
189 acceleration are attended.

190 Fig. 3e shows a decision tree that is trained to detect the attended segments. The root
191 node and its right child node indicate that attended segments correspond to segments with high
192 acceleration values as well as high minimum speed values within a time window. This result
193 indicates that the neural network focuses on a moment of long-lasting acceleration. Fig. 3f shows
194 a decision tree that is trained to classify an input time-series into an appropriate class using features
195 extracted from attended segments within the time-series. As shown in the root node and Fig. 3e, the
196 neural network seems to pay attention to long-lasting accelerations of the DA(+) worms/humans

197 using the attention mechanism, and then distinguishes between DA(+) and DA(-) simply by the
198 minimum of speed within the attended segments.

199 From the above results, we can say that the DA(-) worms/humans present unstable speed
200 while accelerating (lower minimum speed during high acceleration). To validate the hypothesis,
201 we perform a statistical test. Based on the decision trees, we compute the minimum speed within
202 a time window when the average acceleration is high (see Supplementary Information for details).
203 As shown in Fig. 3g, we observe significant differences between DA(+) and DA(-) for both humans
204 and worms. Interestingly, we could also observe significant differences between the two classes of
205 mice (see Supplementary Information).

206 As above, we could notice that the speed of these animals is unstable when accelerating. PD
207 mice used in this study (6-OHDA mouse model of PD) are considered to exhibit motor symptoms
208 of akinesia and bradykinesia¹⁹. The discovered locomotion feature is similar to the symptom of
209 akinesia. Our study reveals that the similar feature is also found in worms lacking D2 dopamine
210 receptors, indicating that the receptors play an important role in the symptom.

211 **Network Training on Worm and Beetle Data** Tonic immobility (sometimes called death-feigning
212 behavior or thanatosis) has been observed in many species and it is thought to have evolved as an
213 anti-predator strategy²⁰. The selection regimes for short or long duration of tonic immobility have
214 been established in the red flour beetle, *Tribolium castaneum*²¹. We use trajectory data of beetles
215 collected from the short and long selection regimes on a treadmill (Fig. 4a; short: 20 beetles; long:
216 20 beetles). The long selection regime showed significantly lower levels of brain dopamine expres-

217 sion and a lower locomotor activity than those of the short selection regime²². Further, Uchiyama
218 et al.²³ showed 518 differentially expressed genes between the selection regimes. As expected from
219 physiological studies described above, genes associated with the metabolic pathways of tyrosine, a
220 precursor of dopamine, were differentially expressed between the short and long selection regime;
221 these enzyme-encoding genes were expressed at higher levels in the long selection regime than in
222 the short selection regime²³. Therefore, we train the neural network with the long selection regime
223 corresponding to the DA(-) class.

224 Fig. 4b, c show the classification accuracy. Fig. 4d shows trajectories of DA(+)/DA(-)
225 worms and beetles highlighted by the trained model, indicating that the network focuses on seg-
226 ments before the movement direction changes in many cases. Fig. 4e shows the time-series of
227 speed for DA(+)/DA(-) worms and beetles highlighted by the trained model. As shown in these
228 figures, the network focuses on segments before the local minima of speed (corresponding to di-
229 rection changes, e.g., $t=17, 23$ for DA(+) beetle). Before the turns, DA(+) worms and beetles seem
230 to quickly decrease their speed. In contrast, DA(-) worms and beetles do not seem to smoothly
231 decrease their speed (e.g., $t=15, 20$ for DA(-) beetle). Based on the observation, we compare the
232 difference in acceleration before turns between the DA(+) and DA(-) classes (see Supplementary
233 Information for details). As shown in Fig. 4f, we could observe significant differences between
234 DA(+) and DA(-) for both beetles and worms. We could also observe significant differences be-
235 tween the two classes of mice (see Supplementary Information for details). Because we only have
236 speed time-series for the human data, we could not perform this test for humans.

237 As above, we observed that DA(-) animals exhibit significantly high acceleration before they
238 change the moving direction. This indicates that the animals with dopamine-deficiency cannot
239 smoothly decrease the speed before turns. Specifically, this hypothesis found by our method fo-
240 cuses on the transition from the “running” mode to the “turn” mode. The disability in locomotion
241 mode transition caused by PD has been studied on mice and humans^{19,24}. These lines of evidence
242 suggest that the disability is caused by combined factors constituting morphological abnormalities
243 of the neural circuitry and changes of the DA transporter and in the DA receptor densities induced
244 by the lack of DA. On the other hand, our cross-species comparative analysis proposes a hypothesis
245 that the disability can be simply explained by the deficiency of dopamine expression level.

246 **DISCUSSION**

247 We propose an attention-based domain-adversarial neural network to study cross-domain behav-
248 ior by analyzing locomotion data from different species. Comparative behavioral analysis be-
249 tween two classes has been performed by using classic classification methods and manual feature
250 design^{14,25,26} as well as studies on locomotion features of PD mice using statistical analysis^{1,2}.
251 However, these features are not necessarily observed in other species. In fact, we could not find
252 significant differences between DA(+) and DA(-) for humans (and worms) by calculating the am-
253 bulation period, widely used to evaluate PD symptoms of mice¹ (see Supplementary Figure 1).
254 DeepHL, which is our prior work, is a pioneering study on deep learning-assisted animal behavior
255 analysis using attention mechanisms²⁷. However, DeepHL focuses only on behavioral data from a
256 single species and cannot be used to conduct cross-species behavior analysis.

257 To demonstrate the usefulness of the proposed neural network, we found cross-species lo-
258 comotion features among species which are far away from each other in the evolutionary lineage.
259 The study reveals that the DA(-) humans, mice, and worms cannot keep high speed. In addition,
260 speed of the DA(-) humans, mice, and worms is unstable when accelerating. Moreover, the DA(-)
261 worms, mice, and beetles present significantly high acceleration before they changing direction.
262 While training the neural network using data from two species, the discovered locomotion features
263 from one of them were observed in other species as well, indicating that the neural network cap-
264 tures latent locomotion features across different species. As a result, we propose the hypothesis
265 that various aspects of motor dysfunction across animal species can be explained by the deficiency
266 of dopamine expression levels. Our method could thus be useful in identifying animal models
267 for a variety of Parkinson’s disease symptoms such as akinesia, bradykinesia, tremor, rigidity, and
268 postural instability.

269 Domain-adversarial neural networks have been originally proposed for transfer learning⁸. To
270 the best of our knowledge, this is the first study that employs them for highlighting cross-species
271 behavioral features. Moreover, this study introduces the attention mechanism to the domain-
272 adversarial network in order to interpret the discovered cross-species behavioral features by the
273 neural network.

274 The proposed method can be potentially applied to evaluate animal models with other dis-
275 eases, accelerating therapeutic drug development. The proposed method can also be used to ana-
276 lyze various time-series such as those for metabolic change and gene expression change.

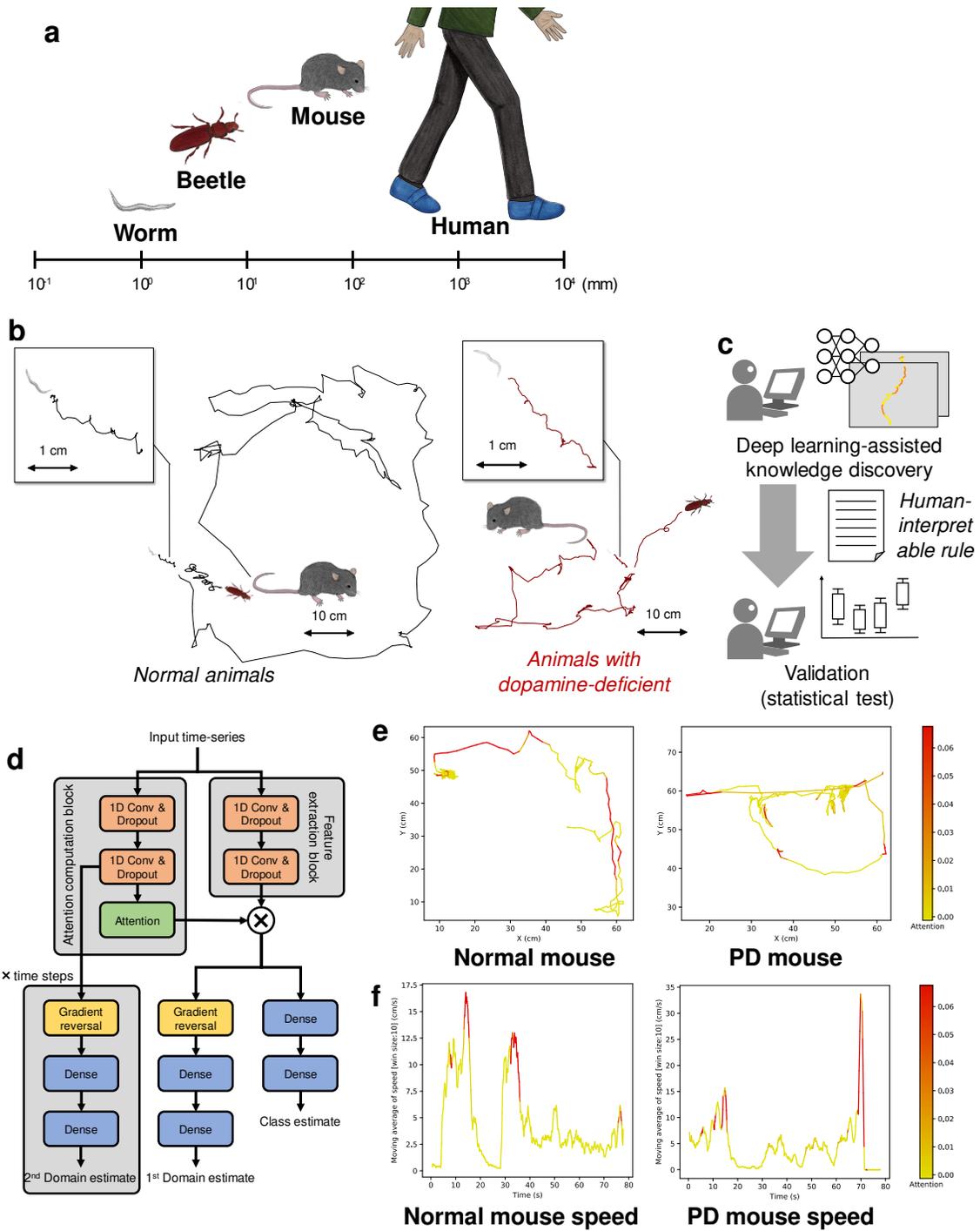


Figure 1: Cross-species behavior analysis using a domain-adversarial neural network with attention mechanism: **a**) differences in body scales among different species; **b**) locomotion trajectories of different animals (worm, beetle, and mouse) differ at the spatial scale, but show similar patterns,

left, black lines depict locomotion trajectories of normal animals, right, red lines show trajectories of dopamine-deficient individuals, inset, expanded trajectories of worms; **c)** our proposed procedure automatically finds a locomotion feature shared by different animals using deep learning, and exhibits the learned locomotion features, enabling the human operator to extract a hypothesis and validate the statistical significance; **d)** proposed network architecture: the feature extraction block learns feature representation that maximizes class prediction accuracy but minimizes domain prediction accuracy by using a gradient reversal layer; the learned feature is assumed to be domain-independent because the feature is incapable of distinguishing between the domains, while the attention computation block computes an attention value for each time slice in the domain-independent manner by using a gradient reversal layer; **e)** highlighted trajectories of normal and PD mice by attention values of the neural network that is trained to extract cross-species locomotion features of worms and mice; **f)** example time-series of speed of normal/PD mice highlighted by our network, where the attention level is color-coded according to the scale bar shown on the right.

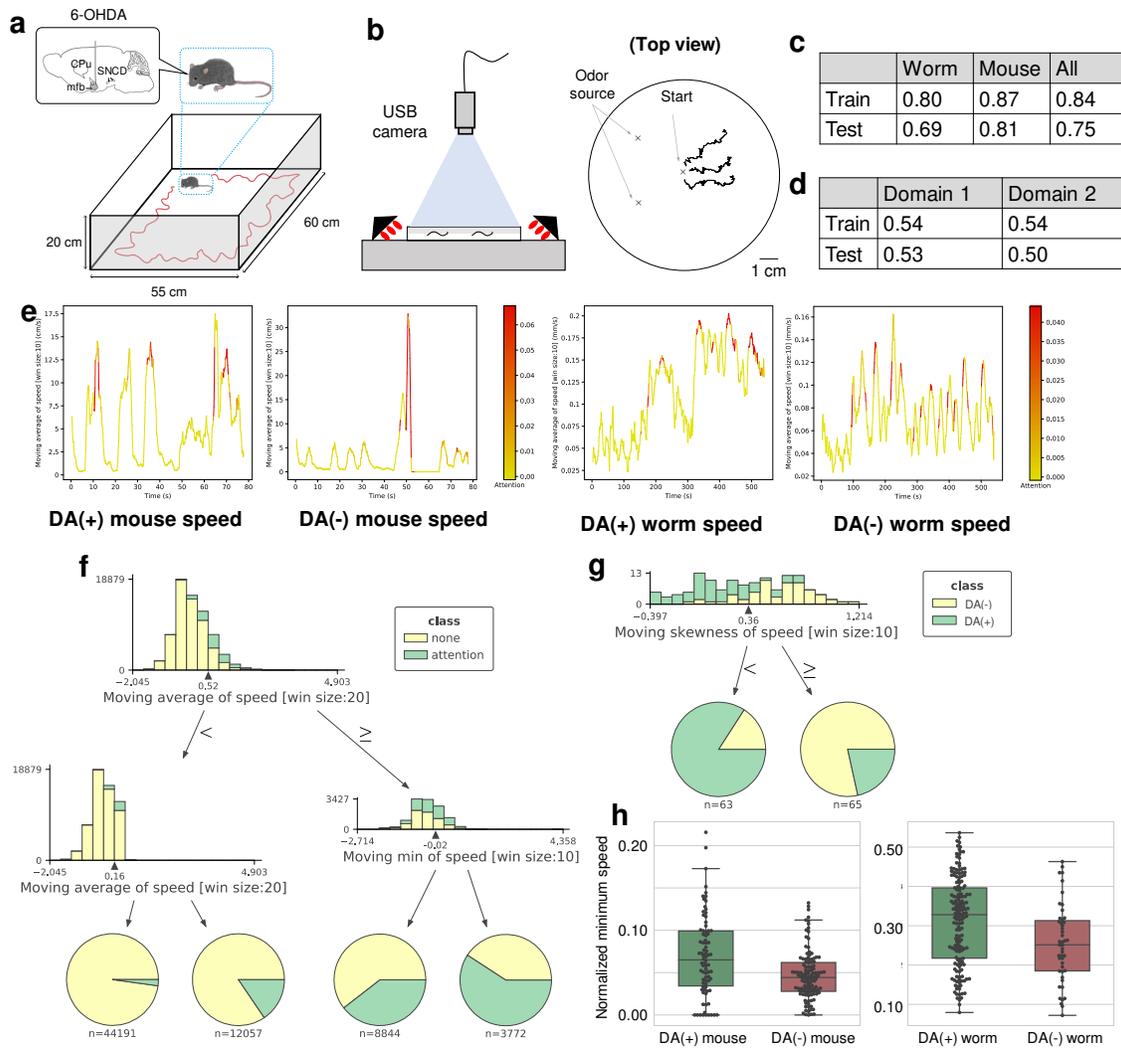
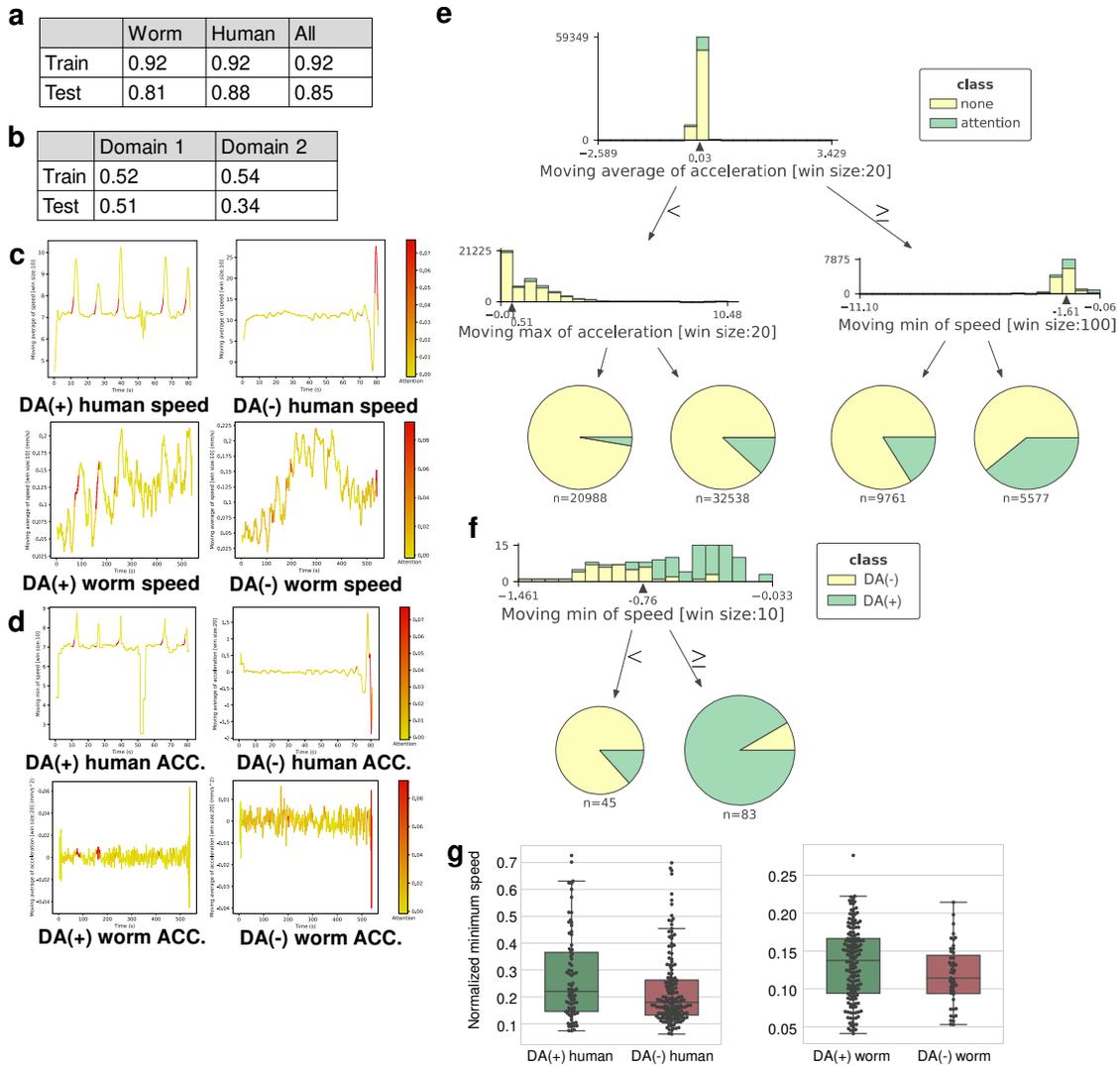


Figure 2: Analysis of DA(+)/DA(-) mice and worms: **a**) a PD mouse walks in the experimental setup, where dopaminergic neurons in the substantia nigra pars compacta were unilaterally lesioned with 6-hydroxydopamine; **b**) the experimental setup (left) for monitoring the worm's trajectory (right); **c**) classification accuracies of our network for DA(+)/DA(-) classes; **d**) classification accuracies of our network for the 1st and 2nd domain estimates, where the random guess ratio is 0.5 (1/2); **e**) example time-series of speed of DA(+)/DA(-) mice and worms highlighted by our network; **f**) a decision tree for explaining the meaning of attention by the neural network, i.e., classifying attended ('attention') and non-attended ('none') segments, where features were normalized for each trajectory (min-max normalization); we show only the 1st and 2nd layers and each histogram shows a tree node of the tree, and a feature used in the node is indicated at the bottom of the histogram; a histogram of each node is constructed from training instances' values of a feature used in the node, instances having feature values smaller than a threshold go to the left child, the threshold value is shown at the bottom of each histogram associated with a black arrow and pie charts at the bottom of the tree (leaf nodes) show distributions of training instances classified by the tree; **g**) a decision tree for explaining classification, where we show only the

1st layer; **h**) distributions of averaged minimum speed within a time window when the speed is high for mice/worms; to compute the averaged minimum speed, we compute the rolling minimum of the normalized speed within segments with high speed (top 20% average speed) for each trajectory (see Supplementary Information for details) and the box plot shows the 25-75% quartile, with embedded bar representing the median and the lower/upper whiskers show $Q1-1.5*IQR$ and $Q3+1.5*IQR$, respectively, where IQR is the interquartile range, Q1 is 25% quartile, and Q3 is 75% quartile; significant differences between DA(+) and DA(-) for both the mice and worms are observed (worm by the Brunner-Munzel test: $p = 7.61 \times 10^{-4}$; $w = 3.48$; $df = 90.41$; effect size = 0.65; mouse by the Welch's t-test: $p = 2.53 \times 10^{-3}$; $t = 3.07$; $df = 131.45$; effect size(r) = 0.45); see Supplementary Information for the normality tests of the distributions, which were used to select methods of statistical test for comparing two groups.



280

Figure 3: Analysis of DA(+)/DA(-) humans and worms: **a**) classification accuracies of our network for DA(+)/DA(-) classes; **b**) classification accuracies of our network for the 1st and 2nd domain estimates; **c**) example time-series of speed of DA(+)/DA(-) humans and worms highlighted by our network; **d**) example time-series of acceleration of DA(+)/DA(-) humans and worms highlighted by our network; **e**) a decision tree for explaining the meaning of attention; **f**) a decision tree for explaining classification; **g**) distributions of minimum speed within a time window when the average acceleration is high for humans and worms (see Supplementary Information for details); significant differences between DA(+) and DA(-) for both humans and worms are observed (worm by the Brunner-Munzel test: $p = 0.03$; $w = 2.26$; $df = 101.91$; effect size = 0.60; human by the Welch's t-test: $p = 0.02$; $t = -2.35$; $df = 126.89$; effect size(r) = -0.34); the p -value is two sided. The box plot shows the 25-75% quartile, with embedded bar representing the median and the lower/upper whiskers show $Q1-1.5*IQR$ and $Q3+1.5*IQR$, respectively, where IQR is the interquartile range, Q1 is 25% quartile, and Q3 is 75% quartile.

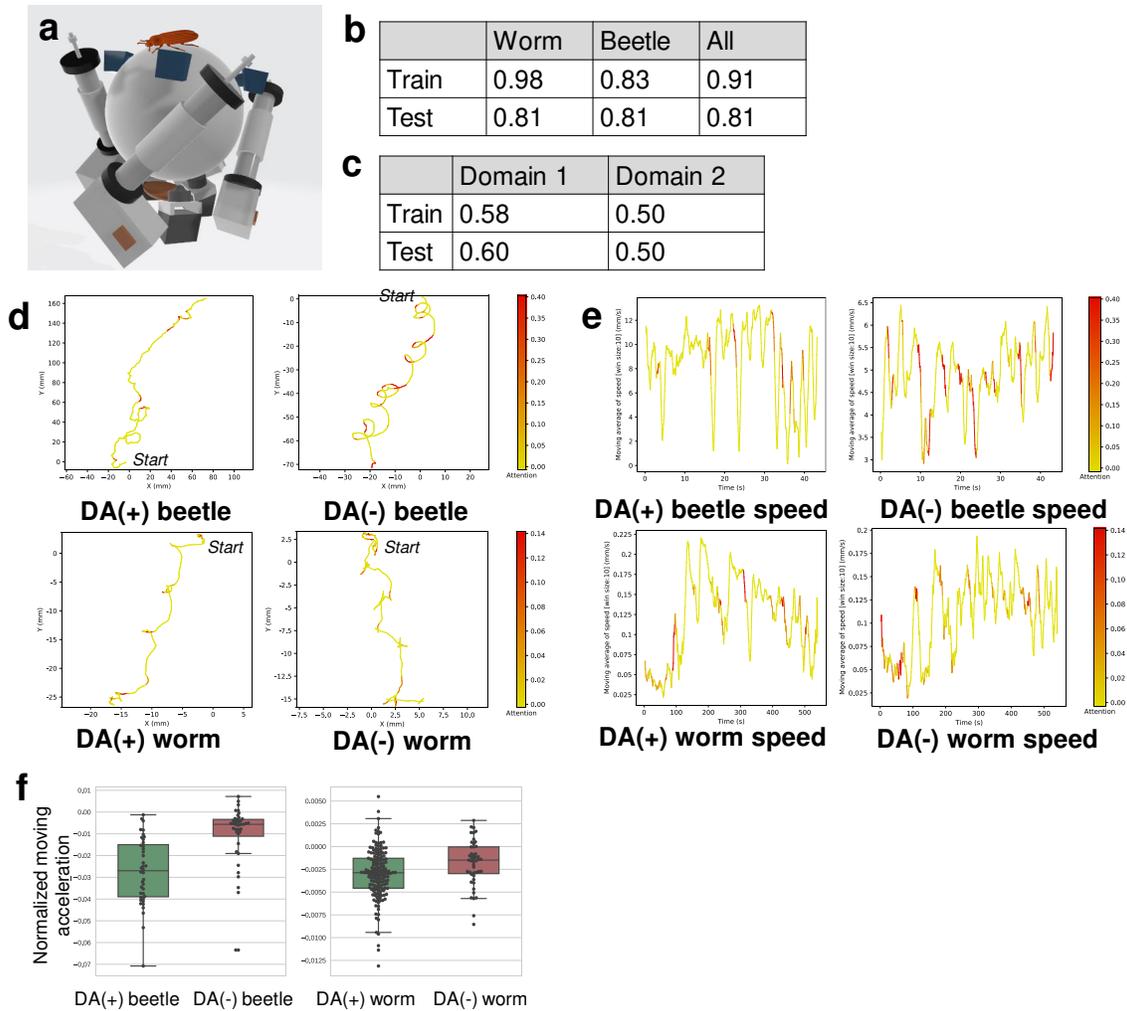


Figure 4: Analysis of DA(+)/DA(-) beetles and worms: **a**) experimental apparatus for beetle study (treadmill); **b**) classification accuracies of our network for DA(+)/DA(-) classes; **c**) classification accuracies of our network for the 1st and 2nd domain estimates; **d**) example trajectories of DA(+)/DA(-) beetles and worms highlighted by our network; **e**) example time-series of speed of DA(+)/DA(-) beetles and worms highlighted by our network; **f**) distributions of acceleration before turns for beetles and worms (see Supplementary Information for details); significant differences between DA(+) and DA(-) for both beetles and worms are observed (worm by the Brunner-Munzel test: $p = 0.005$; $w = -2.90$; $df = 72.76$; effect size = 0.36; beetle by the Brunner-Munzel test: $p = 2.99 \times 10^{-8}$; $w = -6.22$; $df = 71.50$; effect size = 0.18); the p -value is two sided. The box plot shows the 25-75% quartile, with embedded bar representing the median and the lower/upper whiskers show $Q1-1.5 \cdot IQR$ and $Q3+1.5 \cdot IQR$, respectively, where IQR is the interquartile range, $Q1$ is 25% quartile, and $Q3$ is 75% quartile.

282 **Methods**

283 **Worm Data** Young adult wild-type hermaphrodite *Caenorhabditis elegans* (*C. elegans*) were cul-
284 tivated with the bacteria OP50 and handled as described previously²⁸. The *C. elegans* wild-type
285 Bristol strain (N2) were obtained from the Caenorhabditis Genetics Center (University of Min-
286 nesota, USA), and *dop-3(tm1356)* was obtained from National Bioresource Project (Japan) and
287 backcrossed with N2 two times to generate KDK1.

288 The behavioral trajectories of wild-type (“normal”) and *dop-3* worms were monitored during
289 avoidance behavior to the repulsive odor 2-nonanone (normal: 162; *dop-3*: 47)²⁹. During the
290 odor avoidance, naive wild-type and *dop-3* worms exhibit essentially similar behavioral responses
291 although they behave differently after pre-exposure to the odor, indicating that their naive response
292 to the odor is not dependent of dopamine signaling, while their response after pre-exposure is²⁹.
293 Thus, in this study, we focused on the behavior of wild-type and *dop-3* worms after pre-exposure
294 to the odor. The odor avoidance behavior was monitored with a high resolution USB camera for
295 12 min at 1 Hz. Because the worms do not exhibit odor avoidance behavior during the first 2 min
296 because of the rapid increase in odor concentration³⁰, only the data from the following 10 min (i.e.,
297 600 s) was used¹³. A part of the original data had already been analyzed and published^{13,14}, being
298 re-analyzed in this study.

299 We trimmed /undersampled the trajectory data in order to make data lengths of all trajectories
300 (of humans, mice, and beetles) identical before feeding them into the neural network. For more
301 details about trajectory data of the four animals, see Supplementary Table 2.

302 **Mouse Data** Nine C57BL/6J mice purchased from Shimizu Laboratory Supplies (Kyoto, Japan)
303 (male or female; 6-17 months old at the beginning of the experiment) were housed in groups at
304 23°C, with food and water provided ad libitum in a 12h light and 12h dark cycle (day starting at
305 09:00). All tests were performed during the light period. For PD mice, under isoflurane anesthesia,
306 6-OHDA (4 mg/ml; Sigma) was injected through the implanted cannulae (AP -1.2 mm, ML 1.1
307 mm, DV 5.0 mm, $2\mu\text{l}$). The PD mice were allowed to recover for at least one week before post-
308 lesion behavioral testing.

309 We collected 52 trajectories of five normal mice and four unilateral 6-OHDA lesioned mouse
310 models of PD while they freely walked for 10 min in an open arena (60×55 cm, wall height=20
311 cm; normal: 22, PD: 30). The trajectories were tracked from the animal's head position extracted
312 from images captured by a digital video camera (60 fps) mounted on the ceiling of the enclosure.
313 Two sets of small red and green light-emitting diodes mounted above the animal's head were used
314 to track its location in each frame. We then created 150 s segments by splitting each trajectory
315 because training a neural network requires a number of trajectories. We used 201 segments in total
316 (normal: 88, PD: 113) collected from the mice. Note that we excluded 150 s segments that contain
317 no movements of a mouse.

318 **Human Data** We employed a publicly available gait data set of normal and Parkinson's disease
319 humans (Gait in Parkinson's Disease Dataset)³¹. In brief, this data set contains measures of gait
320 from 93 patients with idiopathic PD and 73 healthy controls. The data set includes the vertical
321 ground reaction force records from force 8 sensors underneath each foot with the 100 Hz sampling
322 rate as the subjects walked at their usual for approximately 2 minutes. Note that our study did not

323 use data collected during dual tasking (serial 7 subtractions). Because the durations of almost all
324 of the data were 82 s, we also did not use data shorter than 82 s.

325 **Beetle Data** We analyzed 80 walking trails of beetles collected from short and long selection
326 regimes strain beetles on a treadmill system³². The stock population of *Tribolium castaneum* used
327 in the present study has been maintained in laboratories for more than 25 years. The beetles are
328 fed a blend of whole wheat flour and brewer’s yeast at a 19:1 ratio. They are kept in an incubator
329 (Sanyo, Tokyo, Japan) maintained at 25°C under a 16 h light:8 h dark cycle. The selection regimes
330 with short and long duration of tonic immobility were used³³. The number of beetles derived from
331 the short selection regime (long selection regime) is 20, consisting of 10 males and 10 females.

332 **Preprocessing of Trajectories** We first convert a trajectory into a time-series of speed. Let P be
333 an input trajectory that consists of a sequence of two-dimensional positions with timestamps:

$$P = [P_1, P_2, \dots, P_T]$$
$$= [(t_1, x_1, y_1), (t_2, x_2, y_2), \dots, (t_T, x_T, y_T)]$$

334 For the trajectories of animals whose absolute coordinates are meaningless, such as those of an-
335 imals that freely move on an agar plate, the relative position analysis is required. Therefore, we
336 convert P into S , which is a sequence of speeds:

$$S = [s_2, s_3, \dots, s_T]$$

337 where s_i is speed at time i and described as

$$s_i = \frac{Dist(P_i, P_{i-1})}{t_i - t_{i-1}},$$

338 where $Dist(., .)$ computes the Euclidean distance between two coordinates. After this, we normal-
339 ize speed time-series for each trajectory. Each speed time-series is associated with a class label
340 and domain label.

341 **Processing of Human Gait Data** Because the stride time (i.e., the time elapsed between the first
342 contact of two consecutive footsteps of the same foot) is proportionate to the walking speed³⁴, we
343 first compute time-series of stride time for each foot. We then combine the two time-series, i.e.,
344 by sorting data points of stride times from the two time-series by their time stamps. After that, we
345 standardize a set of speed time-series. Each speed time-series is associated with a class label and
346 domain label.

347 **Attention-based Domain-adversarial Deep Neural Network** Here, we explain the proposed
348 deep neural network model shown in Fig. 1d in detail. The input of the model is a speed time-series
349 S with the length of l . In each 1D convolutional layer of the feature extraction block, we extract
350 features by convolving the input time-series through the time dimension using a filter with a width
351 of F_t . We use a stride (step size) of 1 sample in terms of the time axis. In addition, to reduce over-
352 fitting, we employ dropout, which is a simple regularization technique in which randomly selected
353 neurons are dropped during training³⁵. The dropout rate used in this study is 0.5.

354 We also employ 1D convolutional layers in the attention computation block to compute at-
355 tention time-series from the input time-series S . The attention layer in the attention computation
356 block computes attention from an output matrix of the 2nd convolutional layer Z as follows.

$$\mathbf{a} = \text{softmax}(W_{a2} \cdot \tanh(W_{a1} Z^T)),$$

357 where \mathbf{a} is an attention vector of length l that shows the importance (i.e., attention) of each data
358 point in the input time-series. Because the attention layer is implemented as two densely connected
359 layers with no bias, W_{a1} and W_{a2} show the weight matrices of the 1st and 2nd densely connected
360 layers, respectively. The softmax function ensures that the output values sum to 1, and the tanh
361 function limits the output value of its input to a value between -1 and 1. The attention is multiplied
362 by the outputs of the 2nd convolutional layer in the feature extraction block to contrast the segments
363 to which the network pays attention.

364 The extracted features multiplied by the attention are fed into the class predictor and the
365 1st domain predictor, which are composed of two densely connected layers. The 1st densely
366 connected layers employ the tanh activation function. The 2nd layers (output layers) employ the
367 softmax function. The class predictor and 1st domain predictor output class and domain estimates,
368 respectively. With the gradient reversal layer in the 1st domain predictor, we make the network
369 incapable of estimating domains from extracted features, which is described in detail later.

370 The output of the 2nd convolutional layer of the attention computation block is fed into the
371 2nd domain predictor for each time step. The 2nd domain predictor also consists of two densely
372 connected layers. The 1st densely connected layer employs the tanh activation function. The 2nd
373 layer (output layer) employs the softmax function. The 2nd domain predictor outputs a domain
374 estimate for each time step. With the gradient reversal layer in the 2nd domain predictor, we make
375 the way of attention computation domain-independent.

376 **Network Training** We train the neural network using the backpropagation algorithm¹¹. The gradi-
 377 ent reversal layer in the network multiplies the gradient with a negative constant value ($-\mu$) when
 378 we train the network using the backpropagation algorithm. Because the parameters in the fea-
 379 ture extraction and attention computation blocks are updated so that the domain estimates become
 380 worse, the gradient reversal layer makes the neural network incapable of distinguishing between
 381 the two domains. Note that achieving these different goals at the same time makes it difficult for
 382 the neural network to converge. Therefore, we iterate the three procedures described in the main
 383 text to train the network. Here we explain the procedures in detail.

384 The first procedure minimizes the error of the class estimates as well as maximizes the error
 385 of the domain estimates by the 2nd domain predictor using Equation (1). $\mathcal{L}_c^i(\theta_f, \theta_a, \theta_c)$ in Equation
 386 (1) shows the binary cross entropy loss for class estimates described as follows

$$\mathcal{L}_c^i(\theta_f, \theta_a, \theta_c) = -(y_i \cdot \log p^C(S_i|\theta_f, \theta_a, \theta_c) + (1 - y_i) \cdot \log(1 - p^C(S_i|\theta_f, \theta_a, \theta_c))),$$

387 where y_i is a ground truth class label (0 or 1) of the i -th time-series S_i and $p^C(S_i|\theta_f, \theta_a, \theta_c)$ is a
 388 class estimate by the class predictor. $\mathcal{L}_{d2}^i(\theta_a, \theta_{d2})$ shows the binary cross entropy loss for attention
 389 for each time step described as follows

$$\mathcal{L}_{d2}^i(\theta_a, \theta_{d2}) = -\frac{1}{T-1} \sum_{t=2}^T d_i \cdot \log p^D(s_{i,t}|\theta_a, \theta_{d2}) + (1 - d_i) \cdot \log(1 - p^D(s_{i,t}|\theta_a, \theta_{d2})),$$

390 where d_i is a ground truth domain label (0 or 1) of the i -th time-series, $s_{i,t}$ is a data point in the
 391 i -th time-series at time t and $p^D(s_{i,t}|\theta_a, \theta_{d2})$ is a domain estimate by the 2nd domain predictor.
 392 The second procedure maximizes the error of the domain estimates by the 1st domain predictor
 393 using Equation (2). $\mathcal{L}_{d1}^i(\theta_f, \theta_a, \theta_{d1})$ in Equation (2) shows the binary cross entropy loss for domain

394 estimates described as follows

$$\mathcal{L}_{d1}^i(\theta_f, \theta_a, \theta_{d1}) = -(d_i \cdot \log p^D(S_i|\theta_f, \theta_a, \theta_{d1}) + (1 - d_i) \cdot \log(1 - p^D(S_i|\theta_f, \theta_a, \theta_{d1}))),$$

395 where $p^D(S_i|\theta_f, \theta_a, \theta_{d1})$ is a domain estimate by the 1st domain predictor. We employ the algo-
 396 rithm Adam³⁶ in each procedure to minimize the loss functions. Note that, because parameters in
 397 the network are unstable in the earlier epochs, using large μ makes it difficult for the network to
 398 converge. Therefore, in the earlier epochs, we use small μ to properly train the feature extraction
 399 block and then gradually increase μ as follows.

$$\mu = \begin{cases} 0 & (0 \leq i < T_1) \\ \frac{L+1}{L \cdot (\alpha^{-\beta \frac{i-T_1}{T_2-T_1}}) + 1} - 1 & (T_1 \leq i < T_2) \\ L & (T_2 \leq i) \end{cases}$$

400 where i is the epoch number, L is the upper bound of μ , $\alpha = 1.4$, and $\beta = 10$.

401 **Decision Tree for Explaining Attention** We build a decision tree that explains the meaning of
 402 attention by the network using attention outputs by the network as training labels. We first extract
 403 a feature vector from input time-series for each time slice. We extract interpretable features for
 404 each data point described in Supplementary Table 1. We then label the feature vectors according to
 405 attention outputs by the network. When an attention value at time t is higher than a given threshold,
 406 we label a feature vector at time t as ‘‘attended.’’ Otherwise, we label as ‘‘none.’’ Note that, because
 407 the softmax function in the attention computation block ensures that all attention values in an input
 408 time-series sum to 1, we set the threshold as $\frac{1}{T-1}$. Then we train a binary classifier using the
 409 labeled feature vectors. With the trained decision tree, the user can understand the meaning of the
 410 extracted attention.

411 **Decision Tree for Explaining Classification** We build a decision tree that explains the meaning
412 of classification by the attention-based neural network. We first construct a feature vector for
413 each input time-series by averaging a feature vector prepared for each sliding time window, which
414 is extracted in the same way as that of training a decision tree for explaining attention. Note that,
415 when averaging, we calculate the weighted average according to the attention value by the network.
416 For example, when an attention value of an input time-series at time t is a_t , we multiply a_t to a
417 feature vector at time t . By doing so, we can build a rule mainly focusing on attended segments.
418 After that, we train a decision tree using the averaged feature vectors with their class labels (e.g.,
419 DA(+) or DA(-) class). With the trained decision tree, the user can understand the meaning of the
420 classification by taking into account attention by the network.

421 **References**

- 422 1. Kravitz, A. V. *et al.* Regulation of parkinsonian motor behaviors by optogenetic control of
424 basal ganglia circuitry. *Nature* **466**, 622–626 (2013).
- 425 2. Boix, J., Padel, T. & Paul, G. A partial lesion model of Parkinson’s disease in mice - Charac-
426 terization of a 6-OHDA-induced medial forebrain bundle lesion. *Behavioural Brain Research*
427 **284**, 196–206 (2015).
- 428 3. van der Staay, F. J. Animal models of behavioral dysfunctions: basic concepts and classifica-
429 tions, and an evaluation strategy. *Brain Research Reviews* **52**, 131–159 (2006).
- 430 4. Greenberg, G., Partridge, T., Weiss, E. & Pisula, W. Comparative psychology, a new perspec-
431 tive for the 21st century: Up the spiral staircase. *Developmental Psychobiology: The Journal*

- 432 *of the International Society for Developmental Psychobiology* **44**, 1–15 (2004).
- 433 5. Panksepp, J. & Panksepp, J. B. The seven sins of evolutionary psychology. *Evolution and*
434 *Cognition* **6**, 108–131 (2000).
- 435 6. Panksepp, J., Moskal, J., Panksepp, J. B. & Kroes, R. Comparative approaches in evolutionary
436 psychology: Molecular neuroscience meets the mind. *Neuroendocrinology Letters* **23**, 105–15
437 (2002).
- 438 7. Lickliter, R. The aims and accomplishments of comparative psychology. *Developmental*
439 *Psychobiology: The Journal of the International Society for Developmental Psychobiology*
440 **44**, 26–30 (2004).
- 441 8. Ganin, Y. *et al.* Domain-adversarial training of neural networks. *The Journal of Machine*
442 *Learning Research* **17**, 2096–2030 (2016).
- 443 9. Lin, Z. *et al.* A structured self-attentive sentence embedding. *arXiv preprint arXiv:1703.03130*
444 (2017).
- 445 10. Vaswani, A. *et al.* Attention is all you need. In *Advances in Neural Information Processing*
446 *Systems*, 5998–6008 (2017).
- 447 11. Rumelhart, D. E., Hinton, G. E. & Williams, R. J. Learning representations by back-
448 propagating errors. *Nature* **323**, 533–536 (1986).
- 449 12. Suo, S., Ishiura, S. & Van Tol, H. H. Dopamine receptors in *C. elegans*. *European Journal of*
450 *Pharmacology* **500**, 159–166 (2004).

- 451 13. Yamazoe-Umemoto, A., Fujita, K., Iino, Y., Iwasaki, Y. & Kimura, K. D. Modulation of
452 different behavioral components by neuropeptide and dopamine signalings in non-associative
453 odor learning of *Caenorhabditis elegans*. *Neuroscience Research* **99**, 22–33 (2015).
- 454 14. Yamazaki, S. J. *et al.* A hybrid versatile method for state estimation and feature extraction
455 from the trajectory of animal behavior. *Frontiers in Neuroscience* **13**, 626 (2019).
- 456 15. Albin, R. L., Young, A. B. & Penney, J. B. The functional anatomy of basal ganglia disorders
457 (1989).
- 458 16. DeLong, M. R. Primate models of movement disorders of basal ganglia origin. *Trends in*
459 *Neurosciences* **13**, 281–285 (1990).
- 460 17. Shen, W., Flajolet, M., Greengard, P. & Surmeier, D. J. Dichotomous dopaminergic control of
461 striatal synaptic plasticity. *Science* **321**, 848–851 (2008).
- 462 18. Jankovic, J. & Aguilar, L. G. Current approaches to the treatment of Parkinson's disease.
463 *Neuropsychiatric Disease and Treatment* **4**, 743 (2008).
- 464 19. Mazzoni, P., Shabbott, B. & Cortés, J. C. Motor control abnormalities in Parkinson's disease.
465 *Cold Spring Harbor Perspectives in Medicine* **2**, a009282 (2012).
- 466 20. Ruxton, G. D., Allen, W. L., Sherratt, T. N. & Speed, M. P. *Avoiding attack: the evolutionary*
467 *ecology of crypsis, aposematism, and mimicry* (Oxford University Press, 2019).

- 468 21. Miyatake, T. *et al.* Is death–feigning adaptive? Heritable variation in fitness difference of
469 death–feigning behaviour. *Proceedings of the Royal Society of London. Series B: Biological*
470 *Sciences* **271**, 2293–2296 (2004).
- 471 22. Miyatake, T. *et al.* Pleiotropic antipredator strategies, fleeing and feigning death, correlated
472 with dopamine levels in *Tribolium castaneum*. *Animal Behaviour* **75**, 113–121 (2008).
- 473 23. Uchiyama, H. *et al.* Transcriptomic comparison between beetle strains selected for short and
474 long durations of death feigning. *Scientific Reports* **9** (2019).
- 475 24. Cui, G. *et al.* Concurrent activation of striatal direct and indirect pathways during action
476 initiation. *Nature* **494**, 238–242 (2013).
- 477 25. Chesler, E. J., Wilson, S. G., Lariviere, W. R., Rodriguez-Zas, S. L. & Mogil, J. S. Identi-
478 fication and ranking of genetic and laboratory environment factors influencing a behavioral
479 trait, thermal nociception, via computational analysis of a large data archive. *Neuroscience &*
480 *Biobehavioral Reviews* **26**, 907–923 (2002).
- 481 26. Valletta, J. J., Torney, C., Kings, M., Thornton, A. & Madden, J. Applications of machine
482 learning in animal behaviour studies. *Animal Behaviour* **124**, 203–220 (2017).
- 483 27. Maekawa, T. *et al.* Deep learning-assisted comparative analysis of animal trajectories with
484 DeepHL. *Nature communications* **11** (2020).
- 485 28. Brenner, S. The genetics of *Caenorhabditis elegans*. *Genetics* **77**, 71–94 (1974).

- 486 29. Kimura, K. D., Fujita, K. & Katsura, I. Enhancement of Odor Avoidance Regulated by
487 Dopamine Signaling in *Caenorhabditis elegans*. *Journal of Neuroscience* **30**, 16365–16375
488 (2010).
- 489 30. Tanimoto, Y. *et al.* Calcium dynamics regulating the timing of decision-making in *C. elegans*.
490 *eLife* **6**, e21629 (2017).
- 491 31. Goldberger, A. L. *et al.* PhysioBank, PhysioToolkit, and PhysioNet: components of a new
492 research resource for complex physiologic signals. *Circulation* **101**, e215–e220 (2000).
- 493 32. Nagaya, N. *et al.* Anomalous diffusion on the servosphere: A potential tool for detecting
494 inherent organismal movement patterns. *PLoS ONE* **12**, e0177480 (2017).
- 495 33. Matsumura, K. & Miyatake, T. Responses to relaxed and reverse selection in strains artificially
496 selected for duration of death-feigning behavior in the red flour beetle, *Tribolium castaneum*.
497 *Journal of Ethology* **36**, 161–168 (2018).
- 498 34. Beauchet, O. *et al.* Walking speed-related changes in stride time variability: effects of de-
499 creased speed. *Journal of Neuroengineering and Rehabilitation* **6**, 32 (2009).
- 500 35. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I. & Salakhutdinov, R. Dropout: A sim-
501 ple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*
502 **15**, 1929–1958 (2014).
- 503 36. Kingma, D. P. & Ba, J. L. Adam: a method for stochastic optimization. *arXiv preprint*
504 *arXiv:1412.6980* 1–15 (2014).

505 **Acknowledgements** This work was supported by JSPS Kakenhi JP16H06539, JP16H06545, JP16H06543,
506 and JP17H05976.

507 **Ethics statement** The study on mice was approved by the Doshisha University Institutional Animal Care
508 and Use Committees.

509 **Statistics** We used the Scipy package (v. 1.2.1) of Python (v. 3.6.0) for statistical tests.

510 **Code availability** The software used in this study are available as Supplementary Software.

511 **Data availability** The data of mice, beetles, and worms are available with Supplementary Software.

512 **Competing Interests** The authors declare no competing financial interests.

513 **Contributions** T.M. conceived and directed the study, and performed method design, software implemen-
514 tation, data analysis, and manuscript writing. D.H. performed method design, software implementation,
515 and data analysis. T.H., K.I., S.T., K.D.K., T.M., and K.M. performed data collection, data analysis, and
516 manuscript writing.

517 **Correspondence** Correspondence to Takuya Maekawa.

Figures

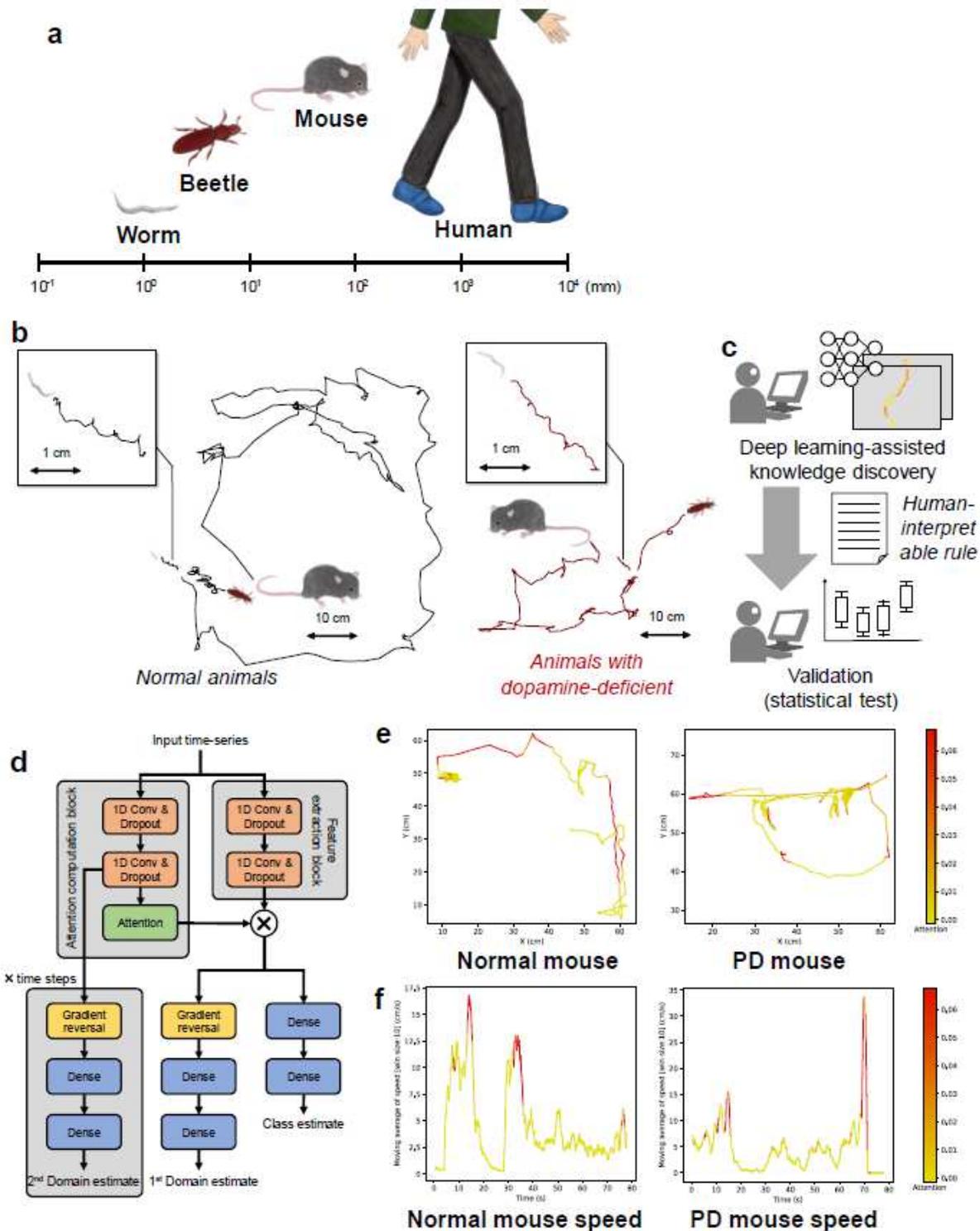


Figure 1

Cross-species behavior analysis using a domain-adversarial neural network with attention mechanism: a) differences in body scales among different species; b) locomotion trajectories of different animals (worm, beetle, and mouse) differ at the spatial scale, but show similar patterns, left, black lines depict

locomotion trajectories of normal animals, right, red lines show trajectories of dopamine-deficient individuals, inset, expanded trajectories of worms; c) our proposed procedure automatically finds a locomotion feature shared by different animals using deep learning, and exhibits the learned locomotion features, enabling the human operator to extract a hypothesis and validate the statistical significance; d) proposed network architecture: the feature extraction block learns feature representation that maximizes class prediction accuracy but minimizes domain prediction accuracy by using a gradient reversal layer; the learned feature is assumed to be domain-independent because the feature is incapable of distinguishing between the domains, while the attention computation block computes an attention value for each time slice in the domain independent manner by using a gradient reversal layer; e) highlighted trajectories of normal and PD mice by attention values of the neural network that is trained to extract cross-species locomotion features of worms and mice; f) example time-series of speed of normal/PD mice highlighted by our network, where the attention level is color-coded according to the scale bar shown on the right.

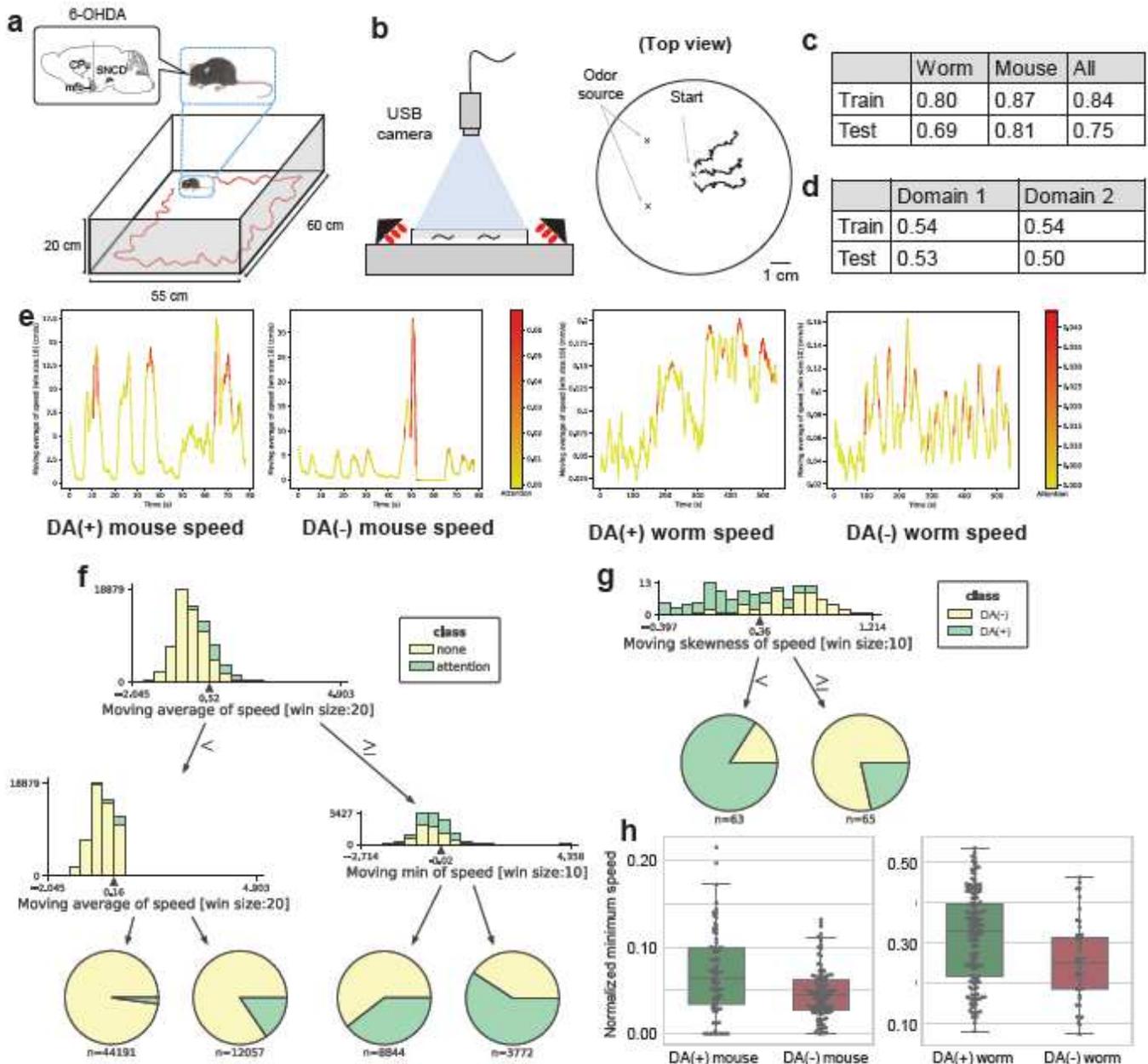


Figure 2

Analysis of DA(+)/DA(-) mice and worms: a) a PD mouse walks in the experimental setup, where dopaminergic neurons in the substantia nigra pars compacta were unilaterally lesioned with 6-hydroxydopamine; b) the experimental setup (left) for monitoring the worm's trajectory (right); c) classification accuracies of our network for DA(+)/DA(-) classes; d) classification accuracies of our network for the 1st and 2nd domain estimates, where the random guess ratio is 0.5 (1/2); e) example time-series of speed of DA(+)/DA(-) mice and worms highlighted by our network; f) a decision tree for explaining the meaning of attention by the neural network, i.e., classifying attended ('attention') and non-attended ('none') segments, where features were normalized for each trajectory (min-max normalization); we show only the 1st and 2nd layers and each histogram shows a tree node of the tree, and a feature used in the node is indicated at the bottom of the histogram; a histogram of each node is constructed from training instances' values of a feature used in the node, instances having feature values smaller than a threshold go to the left child, the threshold value is shown at the bottom of each histogram associated with a black arrow and pie charts at the bottom of the tree (leaf nodes) show distributions of training instances classified by the tree; g) a decision tree for explaining classification, where we show only the 1st layer; h) distributions of averaged minimum speed within a time window when the speed is high for mice/worms; to compute the averaged minimum speed, we compute the rolling minimum of the normalized speed within segments with high speed (top 20% average speed) for each trajectory (see Supplementary Information for details) and the box plot shows the 25-75% quartile, with embedded bar representing the median and the lower/upper whiskers show $Q1-1.5 \times IQR$ and $Q3+1.5 \times IQR$, respectively, where IQR is the interquartile range, Q1 is 25% quartile, and Q3 is 75% quartile; significant differences between DA(+) and DA(-) for both the mice and worms are observed (worm by the Brunner-Munzel test: $p = 7.61 \times 10^{-4}$; $w = 3.48$; $df = 90.41$; effect size = 0.65; mouse by the Welch's t-test: $p = 2.53 \times 10^{-3}$; $t = 3.07$; $df = 131.45$; effect size(r) = 0.45); see Supplementary Information for the normality tests of the distributions, which were used to select methods of statistical test for comparing two groups.

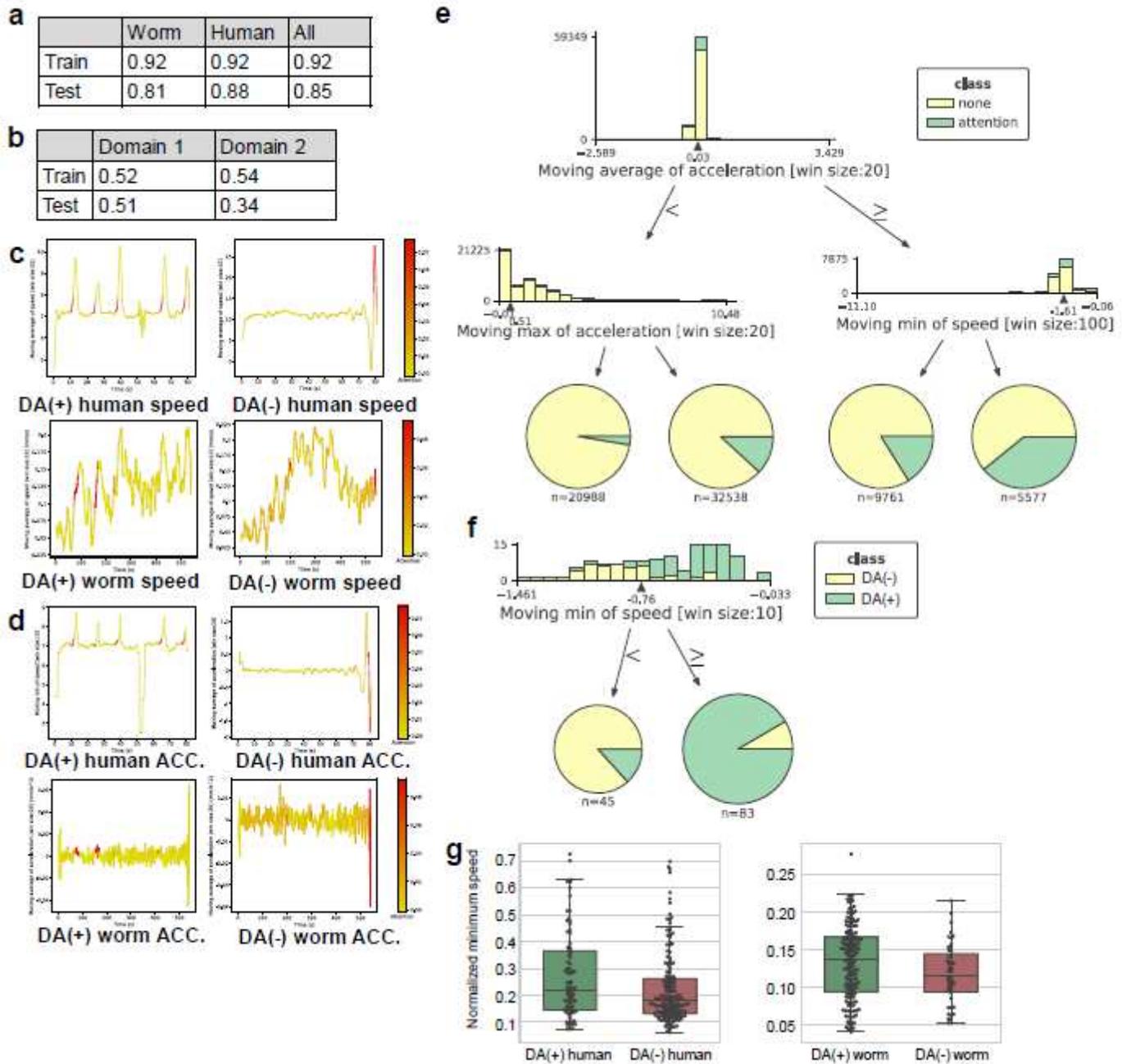


Figure 3

Analysis of DA(+)/DA(-) humans and worms: a) classification accuracies of our network for DA(+)/DA(-) classes; b) classification accuracies of our network for the 1st and 2nd domain estimates; c) example time-series of speed of DA(+)/DA(-) humans and worms highlighted by our network; d) example time-series of acceleration of DA(+)/DA(-) humans and worms highlighted by our network; e) a decision tree for explaining the meaning of attention; f) a decision tree for explaining classification; g) distributions of minimum speed within a time window when the average acceleration is high for humans and worms (see Supplementary Information for details); significant differences between DA(+) and DA(-) for both humans and worms are observed (worm by the Brunner-Munzel test: $p = 0.03$; $w = 2.26$; $df = 101.91$; effect size = 0.60; human by the Welch's t-test: $p = 0.02$; $t = -2.35$; $df = 126.89$; effect size(r) = -0.34); the p-value is two

sided. The box plot shows the 25-75% quartile, with embedded bar representing the median and the lower/upper whiskers show $Q1-1.5 \times IQR$ and $Q3+1.5 \times IQR$, respectively, where IQR is the interquartile range, Q1 is 25% quartile, and Q3 is 75% quartile.

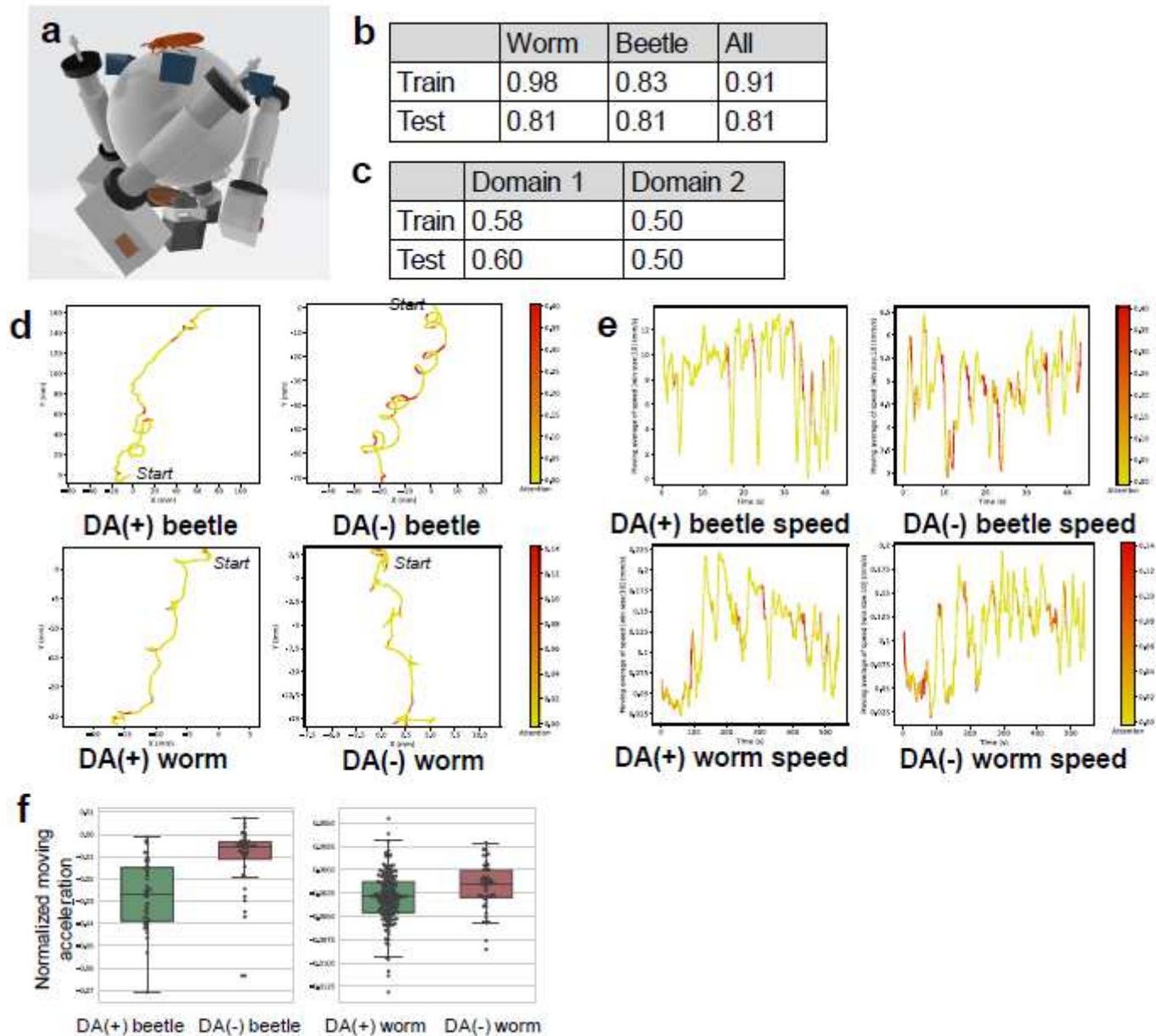


Figure 4

Analysis of DA(+)/DA(-) beetles and worms: a) experimental apparatus for beetle study (treadmill); b) classification accuracies of our network for DA(+)/DA(-) classes; c) classification accuracies of our network for the 1st and 2nd domain estimates; d) example trajectories of DA(+)/DA(-) beetles and worms highlighted by our network; e) example time-series of speed of DA(+)/DA(-) beetles and worms highlighted by our network; f) distributions of acceleration before turns for beetles and worms (see Supplementary Information for details); significant differences between DA(+) and DA(-) for both beetles and worms are observed (worm by the Brunner-Munzel test: $p = 0.005$; $w = -2.90$; $df = 72.76$; effect size = 0.36; beetle by the Brunner-Munzel test: $p = 2.99 \times 10^{-8}$; $w = -6.22$; $df = 71.50$; effect size = 0.18); the p -

value is two sided. The box plot shows the 25-75% quartile, with embedded bar representing the median and the lower/upper whiskers show $Q1-1.5*IQR$ and $Q3+1.5*IQR$, respectively, where IQR is the interquartile range, Q1 is 25% quartile, and Q3 is 75% quartile.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [crosspsi.pdf](#)
- [SupplementarySoftware.zip](#)